

Information Content of Binding Sites on Nucleotide Sequences

Thomas D. Schneider, Gary D. Stormo, Larry Gold

*Department of Molecular, Cellular and Developmental Biology
University of Colorado, Boulder, CO 80309, U.S.A.*

and

Andrzej Ehrenfeucht

*Department of Computer Science
University of Colorado, Boulder, CO 80309, U.S.A.*

(Received 22 June 1984, and in revised form 28 August 1985)

Repressors, polymerases, ribosomes and other macromolecules bind to specific nucleic acid sequences. They can find a binding site only if the sequence has a recognizable pattern. We define a measure of the information (R_{sequence}) in the sequence patterns at binding sites. It allows one to investigate how information is distributed across the sites and to compare one site to another. One can also calculate the amount of information ($R_{\text{frequency}}$) that would be required to locate the sites, given that they occur with some frequency in the genome. Several *Escherichia coli* binding sites were analyzed using these two independent empirical measurements.

The two amounts of information are similar for most of the sites we analyzed. In contrast, bacteriophage T7 RNA polymerase binding sites contain about twice as much information as is necessary for recognition by the T7 polymerase, suggesting that a second protein may bind at T7 promoters. The extra information can be accounted for by a strong symmetry element found at the T7 promoters. This element may be an operator. If this model is correct, these promoters and operators do not share much information. The comparisons between R_{sequence} and $R_{\text{frequency}}$ suggest that the information at binding sites is just sufficient for the sites to be distinguished from the rest of the genome.

1. Introduction

When studying molecular binding sites in DNA or RNA, it is conventional practice to align the sequences of several sites recognized by the same macromolecular recognizer† and then to choose the most common bases at each position to create a consensus sequence (e.g. see Davidson *et al.*, 1983). Consensus sequences are difficult to work with and are not reliable when searching for new sites (Sadler *et al.*, 1983b; Hawley & McClure, 1983). This is partly because information is lost when the relative frequency of specific bases at each position is ignored. For example, the first position of *Escherichia coli* translational initiation codons has

94% A, 5% G, 1% U and 0% C, which is not represented precisely by the consensus "A". To avoid this problem, four histograms can be made that record the frequencies of each base at each position of the aligned sequences. Such histograms can be compressed into a single curve by the use of a χ^2 function (Gold *et al.*, 1981; Stormo *et al.*, 1982b). Although these curves show where information lies in the site, they have several disadvantages: the χ^2 scale is not easily understood in simple terms; it is difficult to compare the overall information content of two different kinds of sites, such as ribosome binding sites and restriction enzyme sites; and χ^2 histograms are not directly useful in searching for new sites (Stormo *et al.*, 1982a).

We present here a method for evaluating the information content of sites recognized by one kind of macromolecule. The method begins with an alignment of known sites, just as with the

† We use the term recognizer to mean a macromolecule that locates specific sites on nucleic acids. These include repressors, activators, polymerases and ribosomes.

evaluation of consensus sequences or χ^2 histograms. However, the calculation of the information content (called R_{sequence}) does not ignore variability of individual positions within a set of sites, as do consensus sequences. Furthermore, R_{sequence} is a measure that encourages direct comparisons between sites recognized by different macromolecules, which is an improvement over χ^2 histograms. R_{sequence} has units of bits per site. The values obtained precisely describe how different the sequences are from all possible sequences in the genome of the organism, in a manner that clearly delineates the important features of the site.

An independent approach is to measure the information needed to find sites in the genome. This relies on the size of the genome and the number of sites in the genome rather than nucleotide sequence information. There is at least one *lac* operator in *E. coli*, while there are thousands of ribosome binding sites. We have defined another measure, $R_{\text{frequency}}$, that is a function of the frequency of sites in the genome. More information would be necessary to identify a single site than any one in a set of thousands. Thus $R_{\text{frequency}}$ is greater for the *lac* operator than for ribosome binding sites. $R_{\text{frequency}}$, like R_{sequence} , is expressed in bits per site.

R_{sequence} , which measures the information in binding site sequences, should be related to the specific binding interaction between the recognizer and the site. $R_{\text{frequency}}$, based only on the frequency of sites, is related to the amount of information required for the sites to be distinguished from all sites in the genome. The problem of how proteins can find their required binding sites among a huge excess of non-sites has been discussed (Lin & Riggs, 1975; von Hippel, 1979). R_{sequence} and $R_{\text{frequency}}$ give us quantitative tools for addressing this problem. Thus we compare R_{sequence} and $R_{\text{frequency}}$ and come to the pleasing conclusion that the values are similar for each site studied. This result was not necessarily expected.

2. Materials and Methods

(a) Calculation of R_{sequence}

(i) Formula for R_{sequence}

Data for calculating R_{sequence} come from 2 sources. One is the nucleotide sequences at which a recognizer has been shown to bind; the other is the nucleotide composition of the genome in which the recognizer functions. The sequences are aligned by a base (the zero base) to give the largest possible homology between them (see Fig. 9 for an example). Some positions have little variation, while others have more. We tabulate the frequency of each base B at each position L in the site, to make a table called $f(B, L)$. Focusing on 1 position at a time, we want to measure the possible variations. For this we have chosen the "uncertainty" measure introduced by Shannon in 1948 (Shannon, 1948; Shannon & Weaver, 1949; Weaver, 1949; Abramson, 1963; Singh, 1966; Gatlin, 1972; Sampson, 1976; Pierce, 1980; Campbell, 1982; Schneider, 1984).

When there are M possible symbols, with probabilities P_i (such that $\sum_{i=1}^M P_i = 1$), the general formula for

uncertainty is:

$$H = - \sum_{i=1}^M P_i \log_2 P_i \text{ (bits per symbol).} \quad (1)$$

One bit of information resolves the uncertainty of choice between 2 equally likely symbols. For nucleotide sequences, there are $M = 4$ possible bases. Using the frequencies of bases as estimates for probabilities, the uncertainty is calculated as:

$$H_s(L) = - \sum_{B=A}^T f(B, L) \log_2 f(B, L) \text{ (bits per base).} \quad (2)$$

(B is either A, C, G or T). The formula gives sensible results for 3 simple cases. (1) If only 1 base appears in the sequences, such as an A, then $f(A, L) = 1$, while the other frequencies are zero. $H_s(L)$ gives zero bits ($0 \log 0 = 0$), meaning that if we were to sequence another site, we would have no uncertainty that the next base will be an A. (2) If 2 bases appeared with equal frequency (as in $f(A, L) = 0.5$, $f(C, L) = 0$, $f(G, L) = 0.5$ and $f(T, L) = 0$), our uncertainty would be 1 bit. (3) If all 4 bases appeared with equal frequencies, then $f(B, L) = 0.25$ and the uncertainty is 2 bits.

If we sequenced randomly in the genome, and aligned sequences arbitrarily, we would see all 4 bases, with probabilities $P(B)$ and our uncertainty about what base we would see next would be:

$$H_g = - \sum_{B=A}^T P(B) \log_2 P(B) \text{ (bits per base).} \quad (3)$$

This number is close to 2 bits for the organism *E. coli*, considered in this paper. In contrast, when sequences are aligned at binding sites (as in typical consensus alignments), a pattern appears that decreases the uncertainty below that of randomly aligned fragments (eqn (2)). For each position L , the decrease would be:

$$R_{\text{sequence}}(L) = H_g - H_s(L) \text{ (bits per base).} \quad (4)$$

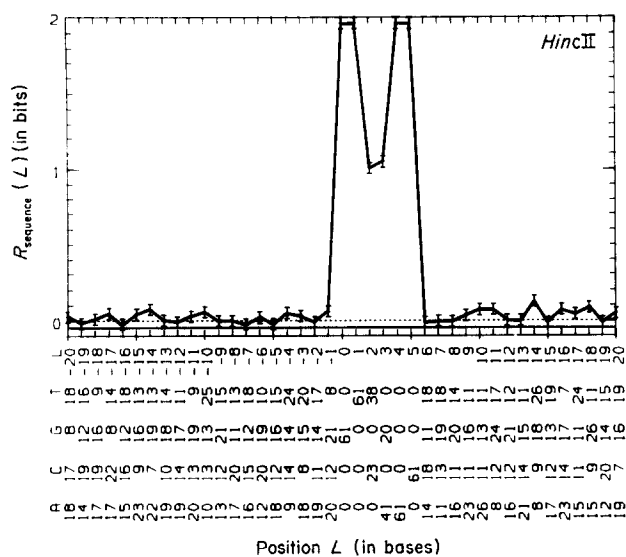
This is a measure of the sequence information gained by aligning the sites. The total information gained will be the total decrease in uncertainty:

$$R_{\text{sequence}} = \sum_L R_{\text{sequence}}(L) \text{ (bits per site).} \quad (5)$$

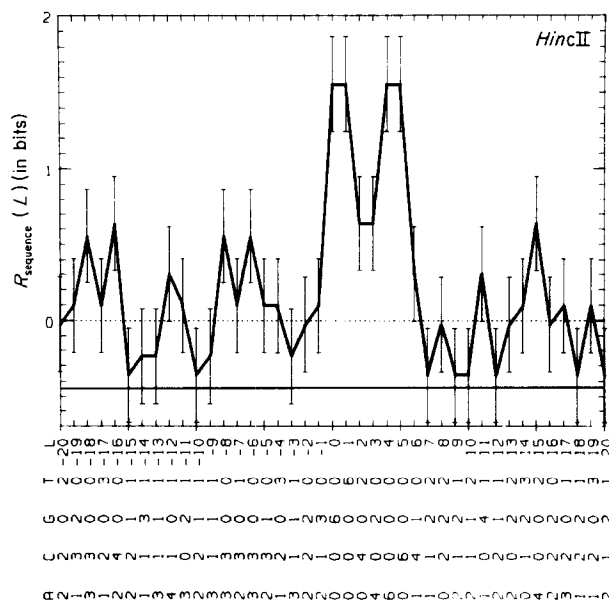
By summing, we make the simplifying assumption that the frequencies at one position are not influenced by those at another position. It is possible also to calculate R_{sequence} from dinucleotides or oligonucleotides (Shannon, 1951; Gatlin, 1972; Lipman & Maizel, 1982). When dinucleotides were used for ribosome binding sites, the total information content was not different from that given in Results (unpublished results). Unfortunately, sampling error prevents one from making the calculation in most cases.

(ii) Graphs of R_{sequence} and correction for sampling error

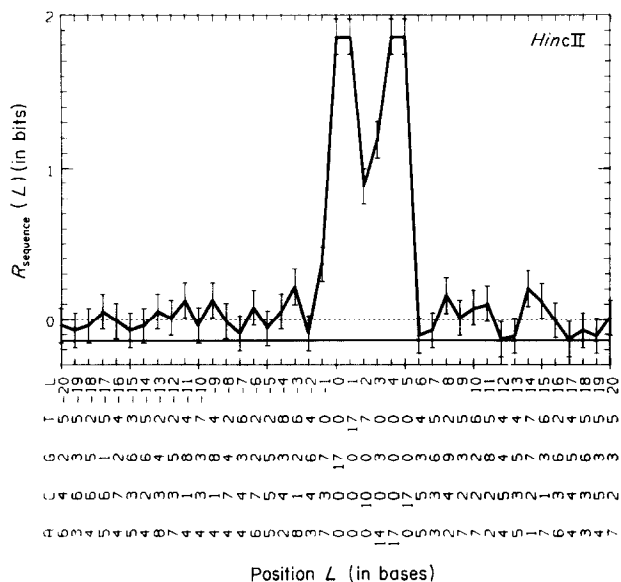
In Fig. 1, we show the curve $R_{\text{sequence}}(L)$ for either (a) 61, (b) 17 or (c) 6 *HincII* sites (G-T-Py-Pu-A-C; Roberts, 1983) chosen from the left end of bacteriophage T7 (Dunn & Studier, 1983). Here, the G residues in the *HincII* sites have been placed at position $L = 0$, and $R_{\text{sequence}}(L)$ was calculated for 20 bases on either side. There are 2 major 2-bit peaks of information content surrounding a 1-bit valley in Fig. 1(a). None of the curves goes to zero (the continuous straight line) outside the sites, although they come close at several points. This effect is not small: for 6 sites (Fig. 1(c)), the background is at 0.44 bit per base, so that with sequences 41 bases long R_{sequence} will be overestimated by 18 bits. A sampling error correction for



(a)



(c)



(b)

Figure 1. Information content, $R_{\text{sequence}}(L)$ in bits/base, at various positions (L) in and around *HincII* sites (G-T-(T/C)-(A/G)-A-C). The numbers of bases at each position, $n(B, L)$, are given. The sites were obtained starting at the left end of the bacteriophage T7 DNA sequence (Dunn & Studier, 1983) and only 1 orientation of each site was used. The left-most base in each site (G) was placed at position 0 in each case, and the sequence examined for 20 nucleotides in each direction from this base. The continuous lines are the zero, without sampling error correction. The broken lines are the zero, when the correction is made. The bars show 1 standard deviation above or below $R_{\text{sequence}}(L)$. They show the variation of the sampling error correction. (a) 61 sites, $R_{\text{sequence}} = 10.7(\pm 0.2)$ bits; (b) 17 sites, $R_{\text{sequence}} = 9.9(\pm 0.7)$ bits; (c) 6 sites, $R_{\text{sequence}} = 8.3(\pm 2.0)$ bits.

$H_s(L)$ ($e(n)$; see Appendix) can be joined with H_g to give the final formula:

$$R_{\text{sequence}} = \sum_L \{E(H_{nb}) - H_s(L)\} \text{ (bits per site).} \quad (6)$$

With this correction, the information content measured at various positions of an aligned set of random sequences will vary above and below zero. On average, it should be zero outside a binding site. The information content inside a site will rise above zero. These features can be seen in all Figures, where the corrected zero is shown as a broken line.

The standard deviation value reported for each R_{sequence} is based on the variance of H_{nb} (Appendix), which is sensitive to the number of sequence examples, but not to the actual sequences. It is only a measure of variance in the correction for small sample sizes; the variation in the information content of individual sites will be described elsewhere. The variance of the sampling correction is shown in all Figures as a bar extending 1

standard deviation above and below the $R_{\text{sequence}}(L)$ curve.

(iii) *Determining the binding site size*

The range is the nucleic acid region over which the sum of $R_{\text{sequence}}(L)$ is taken. If the range is larger than the binding site, the $R_{\text{sequence}}(L)$ fluctuations outside the site will cancel each other on average. On the other hand, if the range is too small, information content will be lost. That is, one must be sure not to delete part of the site.

Determining the range of a site is difficult because experimental methods, such as deletion analysis, chemical protection or footprinting, do not define the exact region contacted. It is dangerous to judge the range by eye from the sequences themselves or the $R_{\text{sequence}}(L)$ curves derived from a small sequence collection (note that some positions of Fig. 1(c) show the same information content as the 1-bit valley). To avoid these difficulties, we have added 5 bases to both sides of the largest range suggested by experimental data. Consequently, the results will be

more variable than they may have been, but it is unlikely that part of a site will be lost. On average, the background will be cancelled, although in specific cases it may not be. In the cases where 2 sites are adjacent, we extend the range to just before the point of overlap. If adjacent sites do interpenetrate, then some of the information content is lost.

When it is likely that a site is symmetrical, both the sequence and its complement are used in the analysis. This doubles the number of sequences available, and refines the answer. If we had arbitrarily chosen an orientation for each sequence, we might have biased the results.

(iv) Variable spacing

When a recognition site has 2 or more parts with various spacings between them, alignment by 1 part may blur out information in the other part. For example, if the 4 variants of this site:

```
ACGTACGTACGTnnnnnnnGGCC
nACGTACGTACGTnnnnnnnGGCC
nnACGTACGTACGTnnnnnnnGGCC
nnnACGTACGTACGTnnnnnnnGGCC
```

●●●●●●●●

occurred with equal frequency, then the positions marked by dots would have zero information content, even though these sequences would give a large information content if they were aligned with each other. To handle this, one may align each part separately and add the information contents together. However, this leads to an overestimate of the information, because the variable spacing is not taken into account. To take it into account, one may calculate how uncertain the spacing is from a tabulation of the frequency of each spacing and subtract this from the total information from the 2 parts. (This is equivalent to increasing the uncertainty of the site, H_s .) For the example above, $R_{\text{sequence}} = 24$ (A-C-G-T-A-C-G-T-A-C-G-T) + 8 (G-G-C-C) - 2 (spacing) = 30 bits. When this was done for ribosome binding sites, the total information content was not different from that given in Results (unpublished results).

(b) Formula for $R_{\text{frequency}}$

If a genome contains G bases, there are $M = G$ ways that its sequence can be aligned or G potential binding sites. If these are all equally likely, then $P_i = 1/G$ and eqn (1) reduces to:

$$H_{gf} = \log_2 G \text{ (bits)}. \quad (7)$$

If the genome contains γ sites, we assume that the probabilities of binding to each site are equal and that the probability of significant binding to other sequences is zero. This allows eqn (1) to be reduced to:

$$H_{sf} = \log_2 \gamma \text{ (bits)} \quad (8)$$

(One property of H is that it is at a maximum when the probabilities are equal. Thus both H_{gf} and H_{sf} are maxima.)

The decrease in positional uncertainty during binding or alignment is the difference:

$$\begin{aligned} R_{\text{frequency}} &= H_{gf} - H_{sf} = -\log_2 \frac{\gamma}{G} \\ &= -\log_2 f \text{ (bits per site)}, \end{aligned} \quad (9)$$

where f is the frequency of sites in the genome.

$R_{\text{frequency}}$ is the amount of information needed to pick γ sites out of G possible sites. As the number of sites in the genome increases, the information needed to find a site

decreases. As long as the simplifying assumption for eqn (8) holds and γ is restricted to the number of known sites (i.e. γ is not an estimate), eqn (9) gives an upper bound on $R_{\text{frequency}}$, since some sites may exist that are not now known. A second property of this formula is that $R_{\text{frequency}}$ is insensitive to small changes in G or γ . The frequency of sites must change by a factor of 2 to alter $R_{\text{frequency}}$ by only 1 bit. The largest possible value of $R_{\text{frequency}}$ occurs for a single site in the genome: $\log_2 G$. (For *E. coli*, $R_{\text{frequency}} = 22.9$ bits in this case.) On the other hand, if all positions in the genome were sites, one would not need any information to find them, and $R_{\text{frequency}}$ would be zero.

The number of potential binding sites G is twice the number of base-pairs in a DNA genome, because there are 2 orientations for a recognizer to bind at each base-pair. A symmetrical recognizer on DNA has 2 ways to bind each base-pair, and both ways are used at a binding site. Here, γ is twice the number of conventional binding sites. An asymmetric recognizer on DNA will use only 1 orientation at any particular base-pair. In this case, γ is equal to the number of binding sites. On RNA, there is only one possible orientation. Thus G and γ reflect both the genome size and number of binding sites, and the symmetries of the recognizer and nucleic acid.

(c) Skewed genomes

This paper considers the relationship between R_{sequence} and $R_{\text{frequency}}$. For restriction enzymes cutting genomes with equal numbers of the 4 bases randomly distributed, R_{sequence} and $R_{\text{frequency}}$ are equal. For example, one commonly assumes that *HaeIII* (G-G-C-C; Roberts, 1983; $R_{\text{sequence}} = 8$ bits) cuts once in 256 bases ($R_{\text{frequency}} = 8$ bits). This is not true for skewed genomes, in which the frequencies of each base are significantly unequal. For example, in a genome like that of bacteriophage T4, which is $\frac{2}{3}$ A + T, R_{sequence} for any tetramer is 7.7 bits. Yet G-G-C-C should occur once in every 1296 bases ($(1/6)^4$; $R_{\text{frequency}} = 10.3$ bits) and conversely, A-A-T-T should occur once in every 81 bases ($(1/3)^4$; $R_{\text{frequency}} = 6.3$ bits). An alternative formula:

$$R_{\text{sequence}}^*(L) = \sum_{B=A}^T f(B,L) \log_2 \frac{f(B,L)}{P(B)} \quad (10)$$

matches $R_{\text{frequency}}$ in examples of this type. When the genomes are equiprobable, as they are in this paper, the 2 R_{sequence} formulae give the same values. We suggest that both be tried for sites in skewed genomes.

(d) Programs and computers

All programs used for analyses were written in Pascal (Jensen & Wirth, 1978; Schneider et al., 1982, 1984). The major programs used were:

Name	Version	Purpose
CalHnb	2-15	Calculate statistics of H_{nb} : $E(H_{nb})$, $AE(H_{nb})$ and $Var(H_{nb})$ (generates (Fig. A2)).
Rseq	4-46	Information content of sequences, R_{sequence} as calculated in this paper (with correction for sampling error).
RsGra	2-45	A non-standard FORTRAN program using device-independent graphics (Warner, 1979) for drawing the Figures on microfilm.

Most of the work was performed on a CDC Cyber 170/720 computer. Figures were generated on a CDC 280/284 microfilm recorder.

(e) Sequence data

We used 2 large prokaryotic sequence databases called LIB1 (bacteriophage) and LIB2 (*E. coli* and *Salmonella typhimurium*: Stormo *et al.*, 1982b) for the sequences of ribosome binding sites. In all, 25 new sites were included: T4 gene 67 (Völker *et al.*, 1982), T4 lysozyme, *IPIII* (Owen *et al.*, 1983); *E. coli* genes *thrB*, *thrC* (Cossart *et al.*, 1981), *rpsT* (Mackie, 1981), *rpsB*, *tsf* (An *et al.*, 1981), *ndh* (Young *et al.*, 1981), *aroH* (Zurawski *et al.*, 1981), *alaS* (Putney *et al.*, 1981), *rpoD* (Burton *et al.*, 1981), *tufA* (Yokota *et al.*, 1980), *unc1*, *unc6*, *uncC*, *uncB*, *uncD*, *uncA* (Gay & Walker, 1981a,b; Kanazawa *et al.*, 1981), *tufB* (An & Friesen, 1980), *lexA* (Horii *et al.*, 1981; Miki *et al.*, 1981; Markham *et al.*, 1981), *ampC* (Jaurin & Grundström, 1981; Jaurin *et al.*, 1981); *EcoRI* endonuclease, methylase (Greene *et al.*, 1981; Newman *et al.*, 1981), DHFR (Swift *et al.*, 1981; Zolg & Hänggi, 1981). Sequences other than ribosome binding sites were stored in a library called SITELI. The corresponding Delila instructions were stored as modules in a single file called SITEIN, and the Module program was used to extract the instructions for each analysis. The sequences for *carAB*, *argI* and *argR* were from Cunin *et al.* (1983). The *lacZ* "pseudo"-operator sequence was from Kalnins *et al.* (1983). The remaining SITELI sequences described in Results were from the GenBank (TM) magnetic tape, release 14.0 (November, 1983), which is available from Bolt Beranek and Newman Inc., Cambridge, Mass., U.S.A.

3. Results

(a) Ribosomes and ribosome binding sites

We aligned the sequence of 149 *E. coli* and coliphage ribosome binding sites by their initiation codons because the process of initiation requires that the *fMet*-tRNA_f binds there. Since ribosomes search mRNA, we used the composition of the transcript library (Stormo *et al.*, 1982b) to calculate H_q : A = 29,526, C = 25,853, G = 27,800, T = 28,951 for which $H_q = 1.99817$ bits per base. The frequencies of bases at each position of the sites were used to find the information content, $R_{\text{sequence}}(L)$, as a function of position (eqns (2), (3) and (A8)). Figure 2 shows that the largest peak is for the initiation codon. The second largest peak represents the Shine-Dalgarno sequence (Shine & Dalgarno, 1974). There are at least five other distinct peaks.

R_{sequence} : the total information content of the site, is found by adding together the individual information contents from each position (eqn (6)). Previous statistical analyses showed a range of -21 to +13 (zero is the first base of the initiation codon), which corresponds well to the regions of RNA protected by ribosomes from ribonucleases (Gold *et al.*, 1981). This range was extended by five bases on both sides. For this range, we calculate an R_{sequence} value of 11.0 bits per site. Alignment by the Shine-Dalgarno sequence gives less than

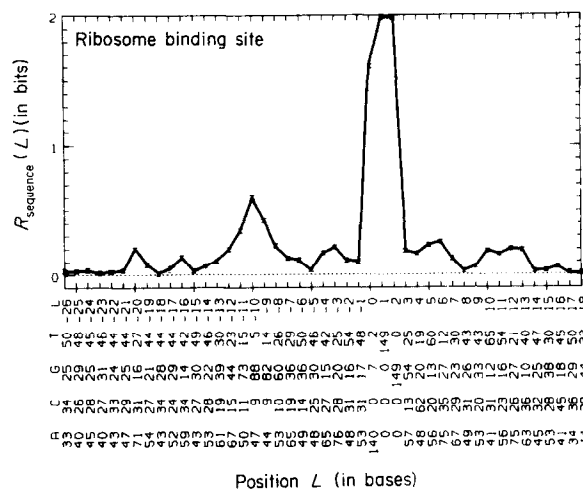


Figure 2. Ribosome binding site information content, determined as for Fig. 1. Position 0 is the first base of the initiation codon.

8.3 bits (data not shown), which suggests that this is not a good alignment.

A good estimate for the size of the *E. coli* genome is 3.9×10^6 base-pairs (Bachmann & Low, 1980). In determining $R_{\text{frequency}}$, we assume that almost all of the genome is transcribed into messages and that, for the most part, only one strand is transcribed. The number of potential ribosome binding sites is therefore 3.9×10^6 . On the basis of the coding capacity versus DNA insert size of 24 plasmids selected at random from the Clark-Carbon bank (P. Bloch, personal communication; Neidhardt *et al.*, 1983), and a genome size of 3.9×10^6 base-pairs, we estimate the number of proteins encoded by *E. coli*, and therefore the number of ribosome binding sites, to be 2574. Equation (9) therefore gives an $R_{\text{frequency}}$ value of 10.6 bits per site. The data for all analyses are presented in Table 1.

(b) *lexA* and SOS boxes

In response to DNA damage, a set of unlinked *E. coli* genes are expressed (Kenyon *et al.*, 1982; Little, 1983; for a review, see Little & Mount, 1982). The genes of the SOS regulatory system are controlled in part at the level of transcription by the direct binding of the *lexA* gene product to the promoters. Five binding sites are well characterized. Two sites are linked to *lexA*, one is linked to each of *recA* (Little *et al.*, 1981; Brent & Ptashne, 1981; Uhlin *et al.*, 1982), *wvrA* (the same site as for *ssb*: Sancar *et al.*, 1982a; Brandsma *et al.*, 1983; Backendorf *et al.*, 1983) and *wvrB* (Sancar *et al.*, 1982b). Two others have been identified reasonably well, at *sulA* (= *sfiA*: Cole, 1983) and on the plasmid *cloDF13* (van den Elzen *et al.*, 1982). Several plasmid promoters may have two deeply overlapping *lexA* sites (Ebina *et al.*, 1981; van den Elzen *et al.*, 1982; Morlon *et al.*, 1983). Since it is possible

Table 1
Information content of several molecular binding sites

Organism	Recognizer	Type	<i>n</i>	Range	R_s	S.D.	γ	$G \times 10^{-6}$	R_f	R_s/R_f	$R_s - R_f$
<i>E. coli</i>	Ribosome	A	149	-26 to 18	11.0	0.1	2574	3.9	10.6	1.0	0.4
<i>E. coli</i>	LexA	E	14	-9 to 10	21.1	0.6	22	7.8	18.4	1.1	2.7
<i>E. coli</i>	TrpR	E	6	-18 to 19	23.4	1.9	6	7.8	20.3	1.1	3.0
<i>E. coli</i>	LacI	O	2	-21 to 21	19.2	2.8	2	7.8	21.9	0.9	-2.6
<i>E. coli</i>	ArgR	E	16	-9 to 10	16.4	0.5	22	7.8	18.4	0.9	-2.0
λ	cI/Cro	O	12	-9 to 9	17.1	0.7	12	7.8	19.3	0.9	-2.2
T7	RNA Pol	A	17	-29 to 12	35.4	0.7	83	7.8	16.5	2.1	18.9
T7	Symmetry	E	34	-6 to 7	16.4	0.2	34	7.8	17.8	0.9	-1.4

Type of site: A, asymmetric, E, symmetric without a central base (even), O, symmetric with a central base (odd). *n*, Number of sequenced sites (for symmetric sites, both strands are counted). The range is the region over which R_{sequence} is calculated. R_s stands for R_{sequence} . S.D. is the standard deviation of R_{sequence} owing to small sample size; the variance of information content for individual sites will be present elsewhere. γ is the number of distinct binding sites in the genome. For symmetrical sites, there are 2 possible ways to bind, so γ is twice the number of conventional sites. G is the number of potential binding sites on the genome. R_f stands for $R_{\text{frequency}}$. Calculations were carried out to 5 decimal places and then rounded.

that one of these is not functional, which would confuse the analysis, we did not use these sites. Since there are two adjacent sites upstream from the *lexA* gene, the range was limited to 20 bases. This is approximately the region protected by LexA protein from digestion by DNase I (Little *et al.*, 1981; Brent & Ptashne, 1981). For both the R_{sequence} and $R_{\text{frequency}}$ calculations, we assumed that LexA repressor binds to its operators symmetrically (Little & Mount, 1982), and that the center of the symmetry is between bases 0 and 1 (Fig. 3). For the 14 example sequences, $R_{\text{sequence}} = 21.1$ bits per site. The nucleotide composition used for this and all

remaining recognizers was from *E. coli* chromosomal DNA (LIB2): A = T = 21,260, C = G = 21,644 (Stormo *et al.*, 1982b). $H_g = 1.99994$ bits per base.

The damage-inducible (*din*) genes are spread around the *E. coli* genome (Little & Mount, 1982), so the size of *E. coli* DNA determines G . There are at least 11 chromosomal genes under *lexA* control (Little & Mount, 1982), giving a minimum estimate for the number of sites γ , and an upper bound on $R_{\text{frequency}}$ of 18.4.

(c) *trp* aporepressor and *trp* operators

At least three operons of *E. coli* are transcriptionally controlled by the *trp* aporepressor: the tryptophan biosynthetic operon *trpEDCBA*, the aromatic amino acid biosynthesis operon *aroH* and the gene for *trp* aporepressor itself, *trpR* (Bennett *et al.*, 1976; Gunsalus & Yanofsky, 1980; Singleton *et al.*, 1980; Bogosian *et al.*, 1981; Zurawski *et al.*, 1981; Joachimiak *et al.*, 1983).

A single dimer of aporepressor binds to the operator in the presence of L-tryptophan (Joachimiak *et al.*, 1983). Likewise, each binding site contains a 2-fold symmetry protected by aporepressor from nucleases. We define the center of this symmetry to be between positions 0 and 1 (Fig. 4). A deletion ending at one end of the *trp* operator, *trp* Δ LC145, is thought to define the range of the sites, since it does not affect repression (Bertrand *et al.*, 1976; Bennett & Yanofsky, 1978). However, when *E. coli* *trp* aporepressor is bound to *trp* operator DNA of *Salmonella typhimurium* and the methylation of unprotected purine residues is measured (Oppenheim *et al.*, 1980), the aporepressor protects the region -13 to +14 rather than -11 to +12. We used the range covering five bases on either side of this protected area, giving $R_{\text{sequence}} = 23.4$ bits per site. If one uses the exact range defined by deletion *trp* Δ LC145, R_{sequence} would be 20.6.

Although non-physiologically high concentrations of *trp* aporepressor can regulate several other

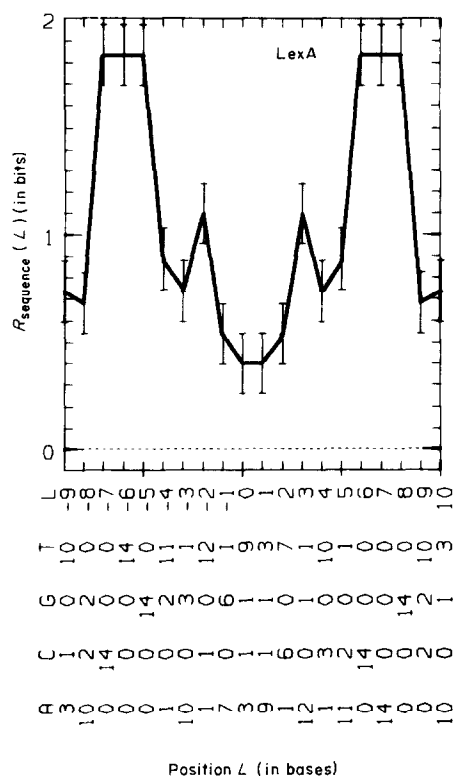


Figure 3. LexA operator information content, determined as for Fig. 1.

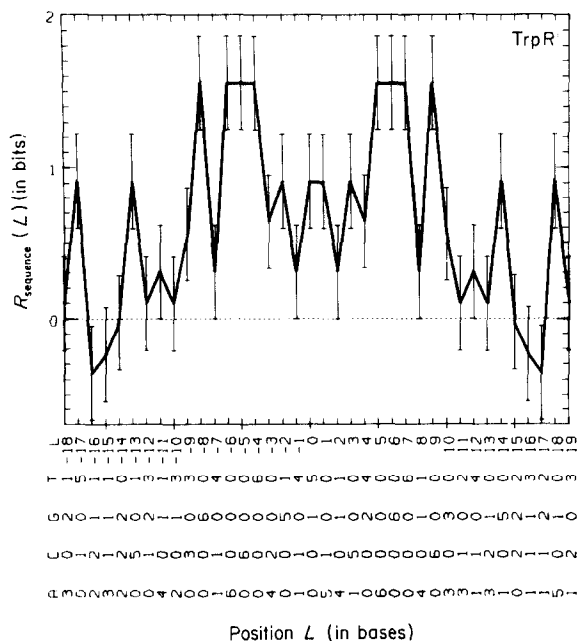


Figure 4. TrpR operator information content, determined as for Fig. 1.

operons (Johnson & Somerville, 1983; Bogosian & Somerville, 1983), we calculate $R_{\text{frequency}}$ for only three sites. The relevant genome is that of *E. coli*, so $R_{\text{frequency}} = 20.3$ bits per site.

(d) lac repressor and the lac operator

One cannot measure information content from a single sequence. Dyad symmetries in DNA (palindromes) are an exception, because the sequence of both the palindrome and its complement are available. This enables us to estimate how much information appears in the lac operator (Beckwith, 1978; Goeddel *et al.*, 1978; Sadler *et al.*, 1983a). Gilbert & Maxam (1973) found that the tetrameric lac repressor protein protects 24 base-pairs from DNase digestion. This is a region from -13 to +10, where the zero is the central base. More recently, exonuclease III digestion give the range from -14 to +16 (Shalloway *et al.*, 1980). To analyze the site, we extended the range from -16 to +16 by five bases on both sides (Fig. 5). This range includes the "extended operator" (Dickson *et al.*, 1975; Heyneker *et al.*, 1976). As with other operators, the sequence was compared to its complement using the program Rseq. The central position was included, giving $R_{\text{sequence}} = 19.2$ bits per site. Because there are only two examples, there is a large sampling error. If there is only one functional lac repressor binding site in the *E. coli* genome, then $R_{\text{frequency}} = 21.9$ bits per site. "Pseudo"-operator sequences exist for which there is no known function (Reznikoff *et al.*, 1974; Winter & von Hippel, 1981). If we include the strong secondary pseudo-operator, $R_{\text{sequence}} = 16.2 \pm 2.6$ and $R_{\text{frequency}} = 20.9$ bits.

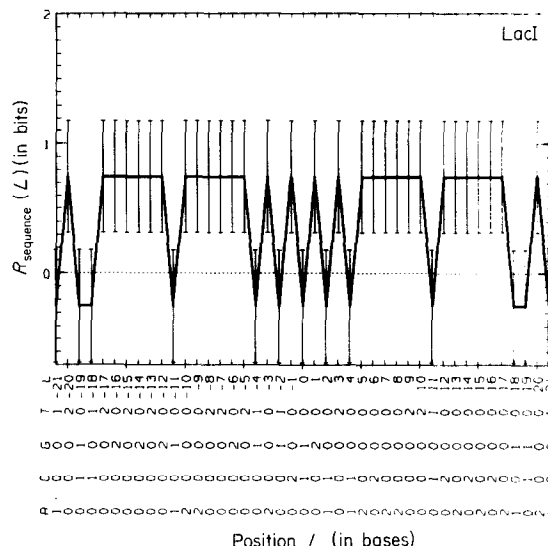


Figure 5. LacI operator information content, determined as for Fig. 1.

(e) argR and arg boxes

The gene *argR* encodes a repressor that controls the synthesis of enzymes of arginine biosynthesis (Maas & Clark, 1964; Maas *et al.*, 1964). Several symmetrical binding sites have been identified tentatively by a few mutations and similarities in sequence (Cunin *et al.*, 1983). Since some sites are adjacent, the range covered only 20 base-pairs (Fig. 6). Also, we used an alignment for the *argR*

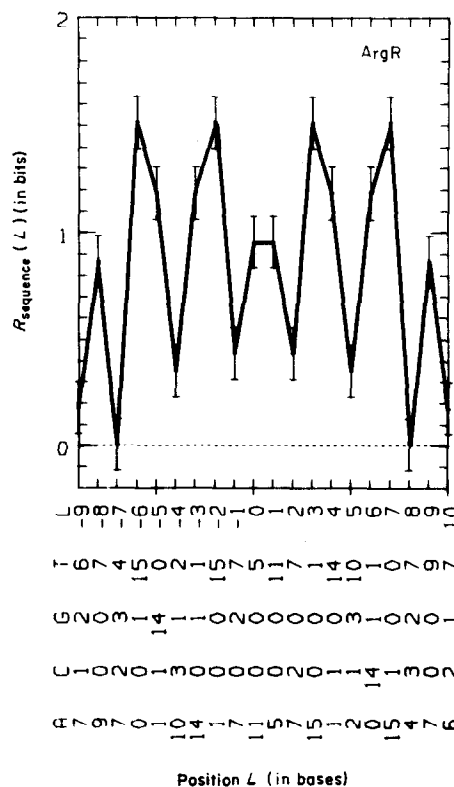


Figure 6. ArgR operator information content, determined as for Fig. 1.

sequence that was shifted one base to the left of that given by Cunin *et al.* (1983). This is presumably a better alignment, because it increased R_{sequence} by 1.5 bits. (It would also improve the "consensus".) $R_{\text{sequence}} = 16.4$, while $R_{\text{frequency}} = 18.4$ bits per site.

By avoiding overlapping sites, we may have deleted part of the arginine boxes. It is possible that two neighboring sites can interpenetrate, if the recognizers bind to different faces of a DNA helix (Hochschild *et al.*, 1983). If the sites are extended to a range -15 to $+16$, R_{sequence} becomes 18.6. In any case, the sites of the arginine regulon have not been characterized by DNase footprinting, chemical protection or other experiments, and several more sites remain to be sequenced.

(f) *cI* repressor, *Cro* and λ operators

All six symmetrical operators of bacteriophage λ are bound by both of the dimeric proteins, repressor and Cro (Ptashne *et al.*, 1976, 1980; Johnson *et al.*, 1981; Matthews *et al.*, 1983). Maniatis *et al.* (1975) originally suggested that the sites are 17 base-pairs wide, separated by A+T-rich "spacers". Since then it has been thought that these regions are not part of the sites. However, a non-random sequence contains information. Chemical protection experiments that probed for guanine residues (Humayun *et al.*, 1977a,b; Johnson *et al.*, 1978; Pabo *et al.*, 1982) did not address the issue, since the region is almost completely devoid of G residues and contacts in the region may not be directed to G·C pairs. Adenine residues were unprotected either because the proteins do not cover that region or because the proteins bind to the opposite side of the DNA from the modifiable group. Two promoter mutations in these regions increase the A+T-richness and do not affect repressor binding (Ptashne *et al.*, 1976; *prm116*, Meyer *et al.*, 1975; *sex1*, Kleid *et al.*, 1976). One mutation, *prm up-1*, decreases the A+T-richness. The effect of *prm up-1* on repressor binding is said to be small (Johnson *et al.*, 1979; Meyer *et al.*, 1980). In contrast to this mutant, nuclease-protection experiments show the sites to be 25 base-pairs wide (Humayun *et al.*, 1977b). Thus it is possible that a portion or all of the spacers are part of the binding sites. However, in keeping with the rules defined in Materials and Methods, we used a range 19 bases wide to avoid overlap between O_{L3} and O_{L2} (Fig. 7). (This also avoids the *prm up-1* site.) Most of the information content of the spacers was lost by this procedure; $R_{\text{sequence}} = 17.1$, $R_{\text{frequency}} = 19.3$ bits per site. If overlaps are ignored, and the sites are extended to the size protected from DNase (25 base-pairs wide, -12 to $+12$), R_{sequence} becomes 19.0.

(g) *T7* RNA polymerase and *T7* promoters

One of the early bacteriophage T7 proteins, encoded by gene 1, is a new RNA polymerase (Chamberlin *et al.*, 1970). This polymerase tran-

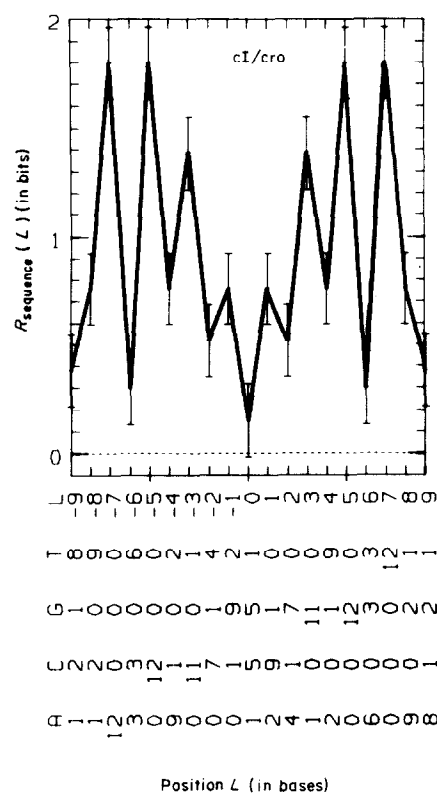


Figure 7. λ *cI/Cro* operator information content, determined as for Fig. 1.

scribes the middle and late genes of the phage genome. Concurrently, the T7 proteins encoded by genes 0.7 and 2 inactivate the host RNA polymerase, so that transcription is directed to the T7 genome rather than that of the host (Hesselbach & Nakada 1977a,b; for reviews on T7, see Studier, 1969, 1972; Krüger & Schroeder, 1981; Dunn & Studier, 1983).

All 17 T7 RNA polymerase promoters have been sequenced (Dunn & Studier, 1983). Deletion experiments *in vitro* and the homology among the promoters suggest that a functional promoter is at least 32 base-pairs long. Five bases beyond the range -24 to $+7$ was used to calculate R_{sequence} (Fig. 8). (The zero base is thought to be the start of each transcript, see Fig. 9 for the alignment.) $R_{\text{sequence}} = 35.4$ bits per site.

To calculate $R_{\text{frequency}}$, we must determine both G and γ . There are two genomes that can contribute to the potential binding sites: the host and the phage. The host DNA is destroyed by the products of gene 3 (endonuclease; Center *et al.*, 1970) and gene 6 (exonuclease; Sadowski & Kerr, 1970), which are synthesized from T7 RNA polymerase-dependent transcripts. They are therefore made following the synthesis of the T7 RNA polymerase. This means that the gene 1 product may search both the *E. coli* and T7 genomes. The T7 genome is only one-hundredth of the size of the host genome, so it does not contribute much. The relevant genome is probably the host DNA. Because promoters are asymmetric, there are twice as many

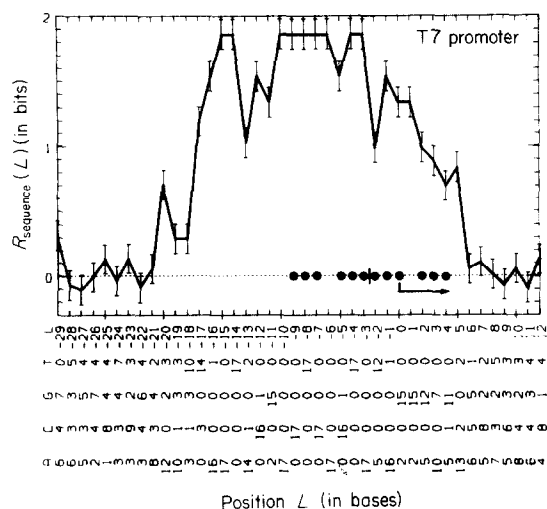


Figure 8. T7 promoter information content, determined as for Fig. 1. The center of the symmetry element is marked by a bar and the points of symmetry by dots. The start of transcription at base zero is shown by an arrow.

potential binding sites on the genome as there are base-pairs, so G is twice the genomic size of *E. coli* (Table 1).

The transcriptional map of T7 is known in great detail (Carter *et al.*, 1981); there are almost certainly no more than 17 T7 polymerase sites (Dunn & Studier, 1983). The activity of T7 RNA polymerase on *E. coli* DNA is 4% of its activity on T7 DNA (Chamberlin & Ring, 1973; see also Summers & Siegel, 1970). Therefore, the total number of sites on *E. coli* DNA could be (17 sites/

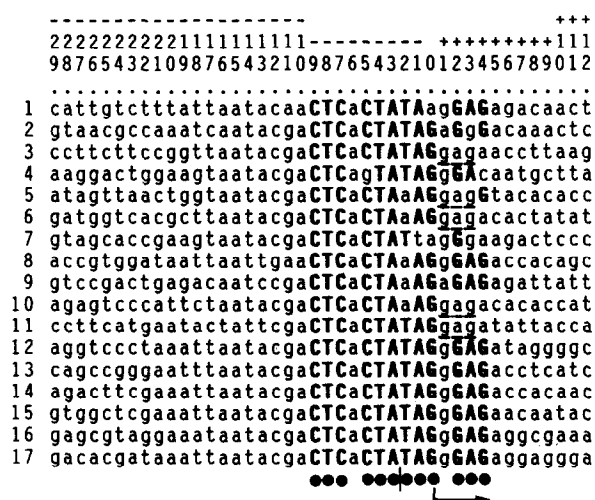


Figure 9. T7 promoter symmetry element. The sequences of the 17 T7 polymerase binding sites are shown. Position zero is presumed to be the start point for transcription (Dunn & Studier, 1983). The position numbers are written vertically. The positions found to be part of the symmetry (Table 2) are shown as capital letters printed in bold-face. The GAGs that may be shifted to the left by 1 base are indicated by an underline.

Table 2
Matches between the left and right halves of the T7 promoter symmetry

Left position	Right position	Number of matches	Probability of matches
-3	-2	12	8.8×10^{-5}
-4	-1	16	3.0×10^{-9}
-5	0	14	1.1×10^{-6}
-6	1	0	7.5×10^{-3}
-7	2	12	8.8×10^{-5}
-8	3	10	2.5×10^{-3}
-9	4	11	5.3×10^{-4}
-10	5	2	0.11
-11	6	4	0.22
-12	7	3	0.19
-13	8	4	0.22
-14	9	5	0.19
-15	10	3	0.19
-16	11	4	0.22
-17	12	3	0.19

The probability of each number of matches is calculated from a binomial distribution, where $p(\text{match}) = 0.25$ and $n = 17$.

39,936 base-pairs T7) $\times (3.9 \times 10^6$ base-pairs *E. coli*) $\times 0.04 = 66$. On infection by T7, there could be as many as 83 sites in the cell. This gives a lower bound for $R_{\text{frequency}}$ of 16.5 bits per site. If there are no sites in the *E. coli* genome, and thus only 17 sites in the cell, $R_{\text{frequency}}$ would be 18.8 bits per site. This is the first case for which R_{sequence} is much bigger than $R_{\text{frequency}}$, so we studied the sequences more closely.

Oakley & Coleman (1977) and Oakley *et al.* (1979) observed that several of the T7 promoters contain a symmetric element centered between bases -3 and -2. The 17 promoter sequences are presented in Figure 9. The extent of the symmetry in all 17 promoters was found by counting the numbers of complementary matches between the two halves. For example, position -14 matches the corresponding position +9 in only five of the 17 sites. This number is likely to occur if the bases were not correlated. The rest of the complementary matches are tabulated in Table 2. Twelve positions have a significantly high number of matches ($p < 0.005$), and these are taken to represent the symmetry. (The positions -6 and 1 are presumably not involved, because they have exceptionally few complementary matches.) Several of the sites contain C-T-C-n-C-T-A : T-A-G-n-G-A-G, while in a few the GAG is shifted to the left by one position.

The information content of these palindromes was determined from the 17 sequences and their complements (34 sequences total) centered as described above (Fig. 10). The R_{sequence} value given in Table 1 is for the 12 positions of the symmetry. R_{sequence} is 16.4 bits per site. There are at least 17 sites in an infected cell, so $R_{\text{frequency}}$ is less than or equal to 17.8 bits per site.

(h) *E. coli* RNA polymerase and *E. coli* promoters
We measured R_{sequence} for sites recognized by *E. coli* RNA polymerase. Hawley & McClure (1983)

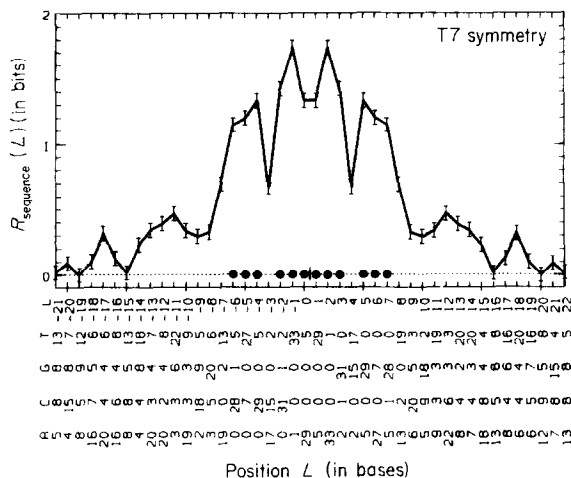


Figure 10. T7 symmetry element information content, determined as for Fig. 1. The information content outside the 12 positions of the symmetry element is from the asymmetric promoter sequences.

compiled data on 112 well-characterized *E. coli* promoters. For these promoters, aligned by the -35 and -10 regions and using the range given by Hawley & McClure (1983), R_{sequence} is only 11.1 bits. There are two difficulties with this analysis. First, a variable gap was introduced between the two regions, which will increase the uncertainty H_s and decrease R_{sequence} substantially, perhaps as much as 2 bits (unpublished results). The other difficulty is that a reasonable estimate for the number of promoters in *E. coli* does not exist, so $R_{\text{frequency}}$ cannot be estimated. Nevertheless, promoters may be more frequent in *E. coli* (one per 500 base-pairs) than is commonly assumed (see Discussion).

4. Discussion

(a) Measurement of R_{sequence}

Many authors have estimated the frequency of a binding site by considering the site size (Gilbert & Müller-Hill, 1970; Riggs *et al.*, 1970; Müller-Hill *et al.*, 1977; Nei & Li, 1979; Pribnow, 1979; von Hippel, 1979; Harel, 1980). R_{sequence} , the sum of $R_{\text{sequence}}(L)$ over a binding site, is similar to the number of bases recognized by macromolecule. In addition, it takes into account the variation of individual sequences. The sampling error correction prevents overestimation of the amount of information in the sequences, but can lead to underestimation in some circumstances (see Fig. 1 and Appendix).

R_{sequence} does not tell us anything about the physical mechanisms a recognizer uses to contact the nucleic acid. For example, the ribosome prefers a particular base composition in the Shine-Dalgarno region. The mechanism is an RNA/RNA contact. *regA*, the translational repressor of bacteriophage T4 (Wiberg & Karam, 1983) uses protein/RNA contacts. It is possible for two such

recognizers to have the same base preferences. Since we use sequences to estimate the probabilities of bases at each position, the analysis will give the same information content for two entirely distinct mechanisms. That is, not only is the mechanism irrelevant to the analysis, but one cannot infer anything about the mechanism from the sequence data, the frequency of bases or the information content, because several mechanisms may give the same results. How physical and chemical contacts determine the preferred base frequencies is a separate question (Pabo & Sauer, 1984).

(b) R_{sequence} for different recognizers

R_{sequence} can be used to investigate relationships between different sites. First, one may ask which binding site has more information than another. For example, ribosome binding sites contain, on the average, less information (11 bits) than do *EcoRI* sites (12 bits). When repressors are compared, R_{sequence} varies between 16 and 23 bits (Table 1), in every case representing a higher information content than that for ribosome binding sites. Indeed, individual repressors regulate transcription at a subset of the *E. coli* genes.

Secondly, the information patterns are different for the various repressors. LexA and TrpR have high peaks three bases wide, while ArgR has double spikes and cI/Cro have single spikes. These distinctive morphological differences probably reflect the location and strength of structural contacts between the different repressors and their cognate sites.

(c) The relationship between R_{sequence} and $R_{\text{frequency}}$

We showed how to estimate the information contained in several binding sites (R_{sequence}), and we determined values for different kinds of sites. But what determines how much information is in a site? One way to approach this question is to make a different measurement, based on "how much information should be needed to locate the sites?" ($R_{\text{frequency}}$) and then compare this to the first measurement. The results of each analysis are summarized by the ratio of R_{sequence} to $R_{\text{frequency}}$ and their difference (Table 1). For ribosomes, LexA, TrpR, LacI, ArgR and cI/Cro, the ratio is close to 1. The sum of the differences for the same six systems is -0.7 bit (out of more than 100 bits of total R_{sequence}).

The large amount of information at T7 polymerase promoters is surprising. We cannot account for this result by using a different sized genome, by changing the number of sites, by sampling error, by overspecification to avoid host sites, or by comparison with *E. coli* promoters. However, there is a simple explanation. The sites have twice as much information as is necessary to locate them in a genome the size of *E. coli*. Therefore, a second recognizer could be using the

extra bits. The sites have symmetry elements that by themselves contain roughly half the information of the entire site. Since T7 RNA polymerase transcribes T7 DNA strictly in one direction (Chamberlin *et al.*, 1970; Summers & Siegel, 1970; Carter *et al.*, 1981; Zavriev & Shemyakin, 1982), it is surprising to find such strong symmetry elements in the promoter sequences. Because the polymerase acts asymmetrically, we assign it to the asymmetric portion of the site.

The symmetric elements could then be the binding sites for the second recognizer. Symmetric elements in promoters suggest the presence of operators (Chamberlin, 1974; Dickson *et al.*, 1975; Dykes *et al.*, 1975; Smith, 1979; Ptashne *et al.*, 1980; Gicquel-Sanzey & Cossart, 1982; Joachimiak *et al.*, 1983). With this in mind, it is intriguing that wild-type T7 bacteriophage decreases late mRNA synthesis around ten minutes after infection, while an amber mutation in gene 3.5 prevents the shutoff; therefore the product of gene 3.5 is a candidate repressor of late T7 transcription (McAllister & Wu, 1978; McAllister *et al.*, 1981; Studier, 1972; Inouye *et al.*, 1973; Jensen & Pryme, 1974; Kerr & Sadowski, 1975; Silberstein *et al.*, 1975; Kleppe *et al.*, 1977; Miyazaki *et al.*, 1978; Krüger & Schroeder, 1981; Dunn & Studier, 1983).

The $R_{\text{sequence}}/R_{\text{frequency}}$ ratio of 2 suggests that there are likely to be two sites at T7 late promoters. In almost all the examples other than T7, a ratio of 1 for $R_{\text{sequence}}/R_{\text{frequency}}$ suggested one site. The exceptional case now becomes the λ operators, where we know that two different proteins bind: *cI* repressor and Cro. (The effects of the third protein that binds these regions, *E. coli* RNA polymerase, are probably blurred out when R_{sequence} is measured.) The existing biochemical and genetic data show that *cI* repressor and Cro bind to the same nucleotides (Johnson *et al.*, 1981). Both λ repressor and Cro are dimers that can bind symmetrically and so may share binding site information. If the two proteins used identical information, the ratio would be 1. If they had used different information the ratio could have been as high as 2, as occurs in the T7 promoter/operator sites. In T7, the proposed repressor would bind symmetrically, and so it could not depend only on information in the asymmetric promoter. Conversely, the polymerase could not depend entirely on symmetrical patterns. That is, asymmetric and symmetric sites must have some separate information.

(d) *How are secondary sites avoided?*

Sequences that are "similar" to true sites might compete with the true sites for binding to the recognizer. For example, the *E. coli* genome should contain about 1000 *EcoRI* restriction enzyme sites (G-A-A-T-T-C), but that same genome should also contain about 18,000 sequences one nucleotide removed from an *EcoRI* site. Site recognition by and action of *EcoRI* within *E. coli* must include

enough discrimination against the more abundant similar sites to avoid a fragmented genome (Pingoud, 1985). Restriction enzymes have enough specificity to do this. It seems that many recognizers do not because operator mutations may decrease binding by only 20-fold (Flashman, 1978). Most single-base changes in promoters and ribosome binding sites decrease synthesis by 2- to 20-fold (Mulligan *et al.*, 1984; Stormo, 1986). Binding to similar sites would degrade the function of the entire system. For repressors, binding to pseudo-operators would increase the chances of gratuitously inhibiting transcription and may also serve as a sink for the recognizer. For ribosomes, binding sites within mRNAs would lead to the expression of inactive protein fragments.

There are several solutions to the problem of avoiding similar sites when the recognizer has limited specificity (Lin & Riggs, 1975). It is possible that similar sites are hidden so that they do not interfere. For example, mRNA secondary structure could prevent ribosomes from inspecting sites similar to ribosome binding sites (Gold *et al.*, 1981). Chromatin structure may occlude the DNA, so that repressors do not actually have as many potential binding sites as the number of base-pairs. A related possibility is that similar sites do not exist in the genome. For example, if a repressor's binding site is composed of oligonucleotides that are relatively rare in the genome, the number of similar sites could be many fewer than expected just from mononucleotide information. Any such special effects constrain the genome to particular oligonucleotide patterns. Discrimination against some oligonucleotides might account for the observed non-random distribution of oligonucleotides in the genome (Grantham *et al.*, 1981; Stormo *et al.*, 1982b; Fickett, 1982; Nussinov, 1984). Finally, von Hippel (1979) pointed out that recognizers could enhance site selectivity by binding to longer sites. If a repressor were to recognize a 15 base-pair long sequence in *E. coli*, not only could its site be unique, but there might not be any sites with one mismatch. When this strategy is used, one expects R_{sequence} to exceed $R_{\text{frequency}}$. The sampling error correction we made may have led to an underestimate of R_{sequence} (see Fig. 1). It is possible that R_{sequence} would be larger if it were calculated from longer oligonucleotides, rather than mononucleotides. We are usually prevented from making that measurement, because the sampling error variance increases rapidly. Still, our results suggest that R_{sequence} is usually close to $R_{\text{frequency}}$.

(e) *Why is R_{sequence} approximately equal to $R_{\text{frequency}}$?*

$R_{\text{frequency}}$ is a function of genome size and the number of sites. Both of these quantities are fixed by factors that have little to do with recognition: genome size is essentially invariant within a species, and the number of sites required by the organism is fixed by physiology and genetics. For example, a

ribosome binding site must precede every gene and the number of genes is determined by physiology and evolutionary history. Unless the population of organisms is undergoing speciation or rapid change in a new environment (Gould, 1977), there is a reasonably fixed frequency of sites and thus $R_{\text{frequency}}$ is approximately fixed. To account for our results, we focus attention on R_{sequence} . Sequence drift will keep R_{sequence} from being larger than is needed for the regulatory process to function properly. If an organism were to have a collection of sites that were more conserved in sequence than was required, mutations in some of the positions of the sites could be tolerated. This would mean an increase in the uncertainty H_s at those positions in the site and a decrease in R_{sequence} . Uncertainty is related to thermodynamic entropy (Shannon, 1948; Tribus & McIrvine, 1971). Just as the entropy of an isolated system tends to increase, excess binding-site information content should tend to atrophy. The lower limit to the drift would be the point at which proper function of the regulatory circuit is diminished.

We are left with many puzzles. How does the information content of sites evolve to equal that needed to find the sites? How is binding energy related to information content? How are chemical contacts related to the base frequencies? What happens in skewed genomes? Lastly, are there situations in biology capable of sustaining large $R_{\text{sequence}}/R_{\text{frequency}}$ ratios, similar to those observed for the T7 late promoters, but for which there is really only one macromolecular recognizer? That is, could a high information content be advantageous for reasons not encountered in the systems studied thus far?

We thank many friends and colleagues for their suggestions, criticisms and patience during the years that this work evolved. We also thank Phil Bloch for a current estimate of the coding capacity of *E. coli*; F. W. Studier for sending us the sequence of T7 before publication; Michael Perry for a general proof for formula (A5); and Kathie Piekarski for typing the manuscript. Computer resources were generously provided by the University of Colorado Academic Computing Services. This work was supported by NIH grant GM28755.

APPENDIX

Calculation of Sampling Uncertainty and Variance

Thomas D. Schneider, Jeffrey S. Haemer and Gary D. Stormo

Using sampling frequencies in place of population probabilities leads to a bias in the uncertainty measurement H (Basharin, 1959). Here we discuss two methods of finding the correction factor when estimating H from a few examples. The first method uses an exact calculation of the average uncertainty for small samples. The probability of obtaining a particular combination of n bases, nb , can be found from a multinomial distribution. The information for the combination, H_{nb} , is calculated and weighted by the probability of obtaining the combination. The weighted information summed for all combinations is the desired result, the expectation of H_{nb} , $E(H_{nb})$. The second method uses a formula to approximate the correction factor.

(a) Exact method

For the exact calculation of $E(H_{nb})$, there are four choices for each base at a position of a site. If one were to calculate H for each possible combination, and then average them, there would be 4^n calculations to perform, where n is the number of sites sequenced. The exact calculation would be impractical for all but the smallest values of n : note that $n = 17$ implies 10^{10} calculations.

Fortunately, the formula for a multinomial distribution permits calculation of many combinations at once (Breiman, 1969). If na , nc , ng and nt are the numbers of A, C, G and T residues in a site and P_a , P_c , P_g , P_t are the frequencies of each base in the genome, then the probability of obtaining a particular combination of na to nt (nb) is estimated by:

$$P_{nb} = \frac{n!}{na!nc!ng!nt!} P_a^{na} P_c^{nc} P_g^{ng} P_t^{nt}, \quad (\text{A1})$$

where $n = na + nc + ng + nt$. The factorial portion on the left is the number of ways that each combination can be arranged. P_{nb} is the probability of obtaining the uncertainty H_{nb} :

$$H_{nb} = - \sum_{b=A}^T \left(\frac{nb}{n} \right) \log_2 \left(\frac{nb}{n} \right). \quad (\text{A2})$$

Finally, to obtain the average uncertainty as decreased owing to sampling:

$$E(H_{nb}) = \sum_{\text{all } nb} P_{nb} H_{nb}. \quad (\text{A3})$$

As a practical matter, one should note that equation (A1) can be calculated quickly by taking the logarithm of the right side and spreading out all

```

NA := N; NC := 0; NG := 0; NT := 0; DONE := FALSE;
REPEAT
  (* CALCULATE EQUATIONS A.1 TO A.3 HERE *)
  IF NT > 0
  THEN BEGIN (* ENDING ON A T - DO OUTER LOOPS *)
    IF NG > 0
    THEN BEGIN (* TURN G INTO T *)
      NG := NG - 1; NT := NT + 1
    END
    ELSE IF NC > 0
    THEN BEGIN (* TURN ONE C INTO G,
      AND ALL T TO G (NOTE NG=NC=0 INITIALLY) *)
      NC := NC - 1; NG := NT + 1; NT := 0
    END
    ELSE IF NA > 0
    THEN BEGIN (* TURN ONE A INTO C AND
      ALL G AND T TO C. (NOTE NG=NC=0 INITIALLY) *)
      NA := NA - 1; NC := NT + 1; NT := 0
    END
    ELSE DONE := TRUE (* SINCE NT = N *)
  END
  ELSE BEGIN (* NO T - INCREMENT INNERMOST LOOP *)
    IF NG > 0
    THEN BEGIN (* TURN G INTO T *)
      NG := NG - 1; NT := NT + 1
    END
    ELSE IF NC > 0
    THEN BEGIN (* TURN C INTO G *)
      NC := NC - 1; NG := NG + 1
    END
    ELSE BEGIN (* NA > 0; TURN A INTO C *)
      NA := NA - 1; NC := NC + 1
    END
  END
UNTIL DONE;

```

Figure A1. Algorithm corresponding to eqn (A4).

the components (including the factorials) into a set of precalculated sums (followed by exponentiation).

The catch in formula (A3) is to avoid calculating all 4^n combinations. A nested series of sums will cover all the required combinations in alphabetical order:

$$\sum_{\text{all } nb} y = \sum_{b_1=A}^T \sum_{b_2=b_1}^T \sum_{b_3=b_2}^T \cdots \sum_{b_n=b_{n-1}}^T y. \quad (\text{A4})$$

At y , in the center of all these sums (nested loops in a computer program), the number of index variables that have the value A must be tallied up to obtain na . This must be done also for nc , ng and nt . Several algorithms to simulate these sums are possible. In Figure A1, we show an algorithm written in Pascal that uses only the variables na , nc , ng and nt to simulate nested loops. The algorithm begins with all A values by setting na to n and the other nb to zero. At each pass through the loop, the sum of $na + nc + ng + nt$ remains invariant. The loop is repeated until the variable **DONE** is set to true after the combination with all T values has been calculated. Since the combinations are covered in alphabetical order, two combinations such as $\{A, A, A, T, C, G\}$ and $\{T, C, G, A, A, A\}$ will be counted only once. The factorial portion of equation (A1) accounts for the actual number of combinations. It can be shown that the loop is entered only

$$(n+1)(n+2)(n+3)/6 \quad (\text{A5})$$

times. Since this is polynomial in n , the direct calculation of $E(H_{nb})$ is practical.

With large numbers of sites, the exact calculation of $E(H_{nb})$ still becomes enormously expensive. For ribosome binding sites, n varies with position in the site. Even if the entire sequence around the site were available, there are sites at the 5' end of a

transcript, so there are regions in the aligned set that must be blank. It is not practical to calculate $E(H_{nb})$ exactly when n is between 108 and 149 (for the range -60 to $+40$).

(b) Approximate method

The second method to calculate the sampling error correction is from Miller (1955) and Basharin (1959), who derived an approximation for the expectation of a sampled uncertainty, $AE(H_{nb})$, that is good for large values of n :

$$AE(H_{nb}) = H_g - \frac{s-1}{2 \ln(2) n} \text{ (bits per base), } (\text{A6})$$

where s , the number of symbols, is 4 for mononucleotides. Figure A2 shows $E(H_{nb})$ and $AE(H_{nb})$ for several values of n . This Table helps one to choose between $AE(H_{nb})$ (a computationally cheap estimate that is inaccurate for small n values but accurate for large n values) and $E(H_{nb})$ (an exact calculation that is computationally costly for large n values). We use $AE(H_{nb})$ above $n = 50$ because the cumulative difference between $E(H_{nb})$ and $AE(H_{nb})$ in a site 100 positions wide would be, at most, 0.078 bit. The exact $E(H_{nb})$ is used for n values less than or equal to 50, since its computation is rapid in this range.

(c) Use of the correction factor

The two methods of calculation produce the expected uncertainty of n sample bases, $E(H_{nb})$:

$$E(H_{nb}) = H_g - e(n) \text{ (bits per base). } (\text{A7})$$

When $H_s(L)$ is calculated from a small sample, it is too small by the amount $e(n)$, on average. To correct $R_{\text{sequence}}(L)$, we use:

$$R_{\text{sequence}}(L) = H_g - [H_s(L) + e(n)] \text{ (bits per base). } (\text{A8})$$

That is, the uncertainty of the pattern is increased because there is only a small sample. Substituting equations (A7) and (A8) into (5) gives equation (6). H_g also could be corrected but the correction is negligible if H_g is calculated from a large sample of the organism's sequence.

The curve for $E(H_{nb})$ as a function of the number of example sites, n (Fig. A3), has several important general properties. As the number of example sites increases, $E(H_{nb})$ approaches H_g ($= 2$ bits/base in the Figs) since the error $e(n)$ becomes smaller. As the number of examples drops, $E(H_{nb})$ also drops (the error increases), until, at only one example, $E(H_{nb})$ is zero. With only one example, the uncertainty of what the sequence is, $H_s(L)$, is also zero. At this point, R_{sequence} is forced to zero (from eqn (6)): one cannot measure an information content from only one example.

The sampling error correction results in an interesting effect. If R_{sequence} could be measured for an infinite number of *HincII* sites (this would look

CALHNB 2.15 CALCULATE STATISTICS OF HNB

GENOMIC COMPOSITION: A = 1, C = 1, G = 1, T = 1
 GENOMIC ENTROPY, HG = 2.00000 BITS

N IS THE NUMBER OF SEQUENCE EXAMPLES
 E(HNB) IS THE EXPECTATION OF THE ENTROPY HNB
 CALCULATED FROM N EXAMPLES
 AE(HNB) AN APPROXIMATION OF E(HNB) THAT IS CALCULATED
 MORE RAPIDLY THAN E(HNB) FOR LARGE N
 E DIFF E(HNB)-AE(HNB)
 VAR(HNB) IS THE VARIANCE OF HNB
 E(N) HG - E(HNB), THE SAMPLING ERROR.

UNITS ARE BITS/BASE, EXCEPT FOR THE VARIANCES WHICH
 ARE THE SQUARE OF THESE.

N	E(HNB)	AE(HNB)	E DIFF	VAR(HNB)	E(N)
1	0.00000	-0.16404	0.16404	0.00000	2.00000
2	0.75000	0.91798	-0.16798	0.18750	1.25000
3	1.11090	1.27865	-0.16775	0.18227	0.88910
4	1.32399	1.45899	-0.13500	0.15171	0.67601
5	1.46291	1.56719	-0.10429	0.12148	0.53709
6	1.55923	1.63933	-0.08010	0.09639	0.44077
7	1.62900	1.69085	-0.06185	0.07661	0.37100
8	1.68129	1.72949	-0.04821	0.06129	0.31871
9	1.72155	1.75955	-0.03800	0.04947	0.27845
10	1.75328	1.78360	-0.03031	0.04034	0.24672
11	1.77879	1.80327	-0.02448	0.03325	0.22121
12	1.79966	1.81966	-0.02000	0.02769	0.20034
13	1.81699	1.83354	-0.01654	0.02331	0.18301
14	1.83159	1.84543	-0.01384	0.01982	0.16841
15	1.84403	1.85573	-0.01170	0.01701	0.15597
16	1.85475	1.86475	-0.00999	0.01473	0.14525
17	1.86408	1.87270	-0.00862	0.01287	0.13592
18	1.87227	1.87978	-0.00750	0.01133	0.12773
19	1.87952	1.88610	-0.00658	0.01004	0.12048
20	1.88598	1.89180	-0.00582	0.00897	0.11402
21	1.89177	1.89695	-0.00518	0.00805	0.10823
22	1.89699	1.90163	-0.00465	0.00727	0.10301
23	1.90172	1.90591	-0.00419	0.00660	0.09828
24	1.90604	1.90983	-0.00380	0.00601	0.09396
25	1.90998	1.91344	-0.00346	0.00551	0.09002
50	1.95594	1.95672	-0.00078	0.00130	0.04406
75	1.97081	1.97115	-0.00034	0.00057	0.02919
100	1.97817	1.97836	-0.00019	0.00032	0.02183
125	1.98257	1.98269	-0.00012	0.00020	0.01743
150	1.98549	1.98557	-0.00008	0.00014	0.01451
175	1.98757	1.98763	-0.00006	0.00010	0.01243
200	1.98913	1.98918	-0.00005	0.00008	0.01087

Figure A2. Statistics of H_{nb} for equiprobable genomic composition. Output of the program CalHnb.

something like Fig. 1(a)), the two peaks would be 2 bits per base. When the correction is made for a small sample, the peaks are less than 2 bits per base (Fig. 1(b) and (c)). This appears odd if we know exactly what *HincII* recognizes. However, given only six examples, we would not be so sure what the "real" pattern is. The sampling-error correction prevents us from assuming that we have more knowledge than can be obtained from the sequences alone. That is, the value $e(n)$ represents our

uncertainty of the pattern, owing to small sample size. In the extreme case of one sequence, we have no knowledge of what the pattern at the site is, even though we see a sequence. Because of the correction, R_{sequence} will be underestimated at truly conserved positions when only a few sites are known. R_{sequence} for six *HincII* sites in Figure 1(c) is estimated to be 8 bits, even though we "know" (by looking at more than 6 examples) that *HincII* recognizes 10 bits.

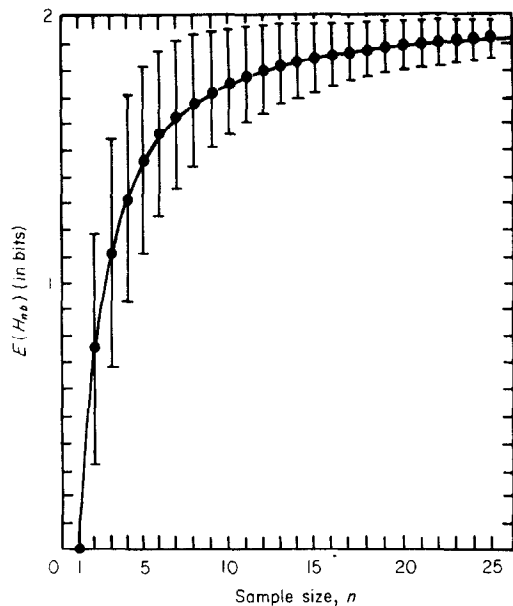


Figure A3. $E(H_{nb})$ versus number of sites, n . These data are for an equiprobable genomic composition. The curve is less than 1% lower for the composition of *E. coli*. Each bar represents 1 s.d. above and below the curve.

(d) *Variance of the correction factor*

$E(H_{nb})$ is the mean of the noisy estimate H_{nb} . The variance (calculated exactly) can be shown to be:

$$\text{Var}(H_{nb}) = \left\{ \sum_{\text{all } H_{nb}} P_{nb}(H_{nb})^2 \right\} - E(H_{nb})^2. \quad (\text{A9})$$

This can be used to estimate the standard deviation of R_{sequence} owing to sampling error. If a site is r bases wide, then the standard deviations is $\sqrt{r \text{Var}(H_{nb})}$.

References

- Abramson, N. (1963). *Information Theory and Coding*, McGraw-Hill Book Co., New York.
- An, G. & Friesen, J. D. (1980). *Gene*, **12**, 33–39.
- An, G., Bendiak, D. S., Mamelak, L. A. & Friesen, J. D. (1981). *Nucl. Acids Res.* **9**, 4163–4172.
- Bachmann, B. J. & Low, K. B. (1980). *Microbiol. Rev.* **44**, 1–56.
- Backendorf, C., Brandsma, J. A., Kartasova, T. & van de Putte, P. (1983). *Nucl. Acids Res.* **11**, 5795–5810.
- Basharin, G. P. (1959). *Theory Probability Appl.* **4**, 333–336.
- Beckwith, J. R. (1978). *The Operon* (Miller, J. H. & Reznikoff, W. S., eds), pp. 11–30, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y.
- Bennett, G. N. & Yanofsky, C. (1978). *J. Mol. Biol.* **121**, 179–192.
- Bennett, G. N., Schweingruber, M. E., Brown, K. D., Squires, C. & Yanofsky, C. (1976). *Proc. Nat. Acad. Sci., U.S.A.* **73**, 2351–2355.
- Bertrand, K., Squires, C. & Yanofsky, C. (1976). *J. Mol. Biol.* **103**, 319–337.
- Bogosian, G. & Somerville, R. (1983). *Mol. Gen. Genet.* **191**, 51–58.
- Bogosian, G., Bertrand, K. & Somerville, R. (1981). *J. Mol. Biol.* **149**, 821–825.
- Brandsma, J. A., Bosch, D., Backendorf, C. & van de Putte, P. (1983). *Nature (London)*, **305**, 243–245.
- Breiman, L. (1969). *Probability and Stochastic Processes, with a View Towards Applications*, Houghton Mifflin Co., Boston.
- Brent, R. & Ptashne, M. (1981). *Proc. Nat. Acad. Sci., U.S.A.* **78**, 4204–4208.
- Burton, Z., Burgess, R. R., Lin, J., Moore, D., Holder, S. & Gross, C. A. (1981). *Nucl. Acids Res.* **9**, 2889–2903.
- Campbell, J. (1982). *Grammatical Man: Information, Entropy, Language, and Life*, Simon and Schuster, New York.
- Carter, A. D., Morris, C. E. & McAllister, W. T. (1981). *J. Virol.* **37**, 636–642.
- Center, M. S., Studier, F. W. & Richardson, C. C. (1970). *Proc. Nat. Acad. Sci., U.S.A.* **65**, 242–248.
- Chamberlin, M. J. (1974). *Annu. Rev. Biochem.* **43**, 721–775.
- Chamberlin, M. J. & Ring, J. (1973). *J. Biol. Chem.* **248**, 2235–2244.
- Chamberlin, M., McGrath, J. & Waskell, L. (1970). *Nature (London)*, **228**, 227–231.
- Cole, S. T. (1983). *Mol. Gen. Genet.* **189**, 400–404.
- Cossart, P., Katinka, M. & Yaniv, M. (1981). *Nucl. Acids Res.* **9**, 339–347.
- Cunin, R., Eckhardt, T., Piette, J., Boyen, A., Piérard, A. & Glansdorff, N. (1983). *Nucl. Acids Res.* **11**, 5007–5019.
- Davidson, E. H., Jacobs, H. T. & Britten, R. J. (1983). *Nature (London)*, **301**, 468–470.
- Dickson, R. C., Abelson, J., Barnes, W. & Reznikoff, W. S. (1975). *Science*, **187**, 27–35.
- Dunn, J. J. & Studier, F. W. (1983). *J. Mol. Biol.* **166**, 477–535.
- Dykes, G., Bambara, R., Marians, K. & Wu, R. (1975). *Nucl. Acids Res.* **2**, 327–345.
- Ebina, Y., Kishi, F., Miki, T., Kagamiyama, H., Nakazawa, T. & Nakazawa, A. (1981). *Gene*, **15**, 119–126.
- Fickett, J. W. (1982). *Nucl. Acids Res.* **10**, 5303–5318.
- Flashman, S. M. (1978). *Mol. Gen. Genet.* **166**, 61–73.
- Gatlin, L. L. (1972). *Information Theory and the Living System*, Columbia University Press, New York.
- Gay, N. J. & Walker, J. E. (1981a). *Nucl. Acids Res.* **9**, 2187–2194.
- Gay, N. J. & Walker, J. E. (1981b). *Nucl. Acids Res.* **9**, 3919–3926.
- Gicquel-Sanzey, B. & Cossart, P. (1982). *EMBO J.* **1**, 591–595.
- Gilbert, W. & Maxam, A. (1973). *Proc. Nat. Acad. Sci., U.S.A.* **70**, 3581–3584.
- Gilbert, W. & Müller-Hill, B. (1970). *The Lactose Operon* (Beckwith, J. R. & Zipser, D., eds), p. 104, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y.
- Goeddel, D. V., Yansura, D. G. & Caruthers, M. H. (1978). *Proc. Nat. Acad. Sci., U.S.A.* **75**, 3578–3582.
- Gold, L., Pribnow, D., Schneider, T., Shinedling, S., Singer, B. S. & Stormo, G. (1981). *Annu. Rev. Microbiol.* **35**, 365–403.
- Gould, S. J. (1977). *Ever Since Darwin*, W. W. Norton & Co., Inc., New York.
- Grantham, R., Gautier, C., Gouy, M., Jacobzone, M. & Mercier, R. (1981). *Nucl. Acids Res.* **9**, r43–r74.
- Greene, P. J., Gupta, M., Boyer, H. W., Brown, W. E. & Rosenberg, J. M. (1981). *J. Biol. Chem.* **256**, 2143–2153.

- Gunsalus, R. P. & Yanofsky, C. (1980). *Proc. Nat. Acad. Sci., U.S.A.* **77**, 7117-7121.
- Harel, D. (1980). *Commun. ACM.* **23**, 379-389.
- Hawley, D. K. & McClure, W. R. (1983). *Nucl. Acids Res.* **11**, 2237-2255.
- Hesselbach, B. A. & Nakada, D. (1977a). *J. Virol.* **24**, 736-745.
- Hesselbach, B. A. & Nakada, D. (1977b). *J. Virol.* **24**, 746-760.
- Heyneker, H. L., Shine, J., Goodman, H. M., Boyer, H. W., Rosenberg, J., Dickerson, R. E., Narang, S. A., Itakura, K., Lin, S. & Riggs, A. D. (1976). *Nature (London)*, **263**, 748-752.
- Hochschild, A., Irwin, N. & Ptashne, M. (1983). *Cell*, **32**, 319-325.
- Horii, T., Ogawa, T. & Ogawa, H. (1981). *Cell*, **23**, 689-697.
- Humayun, Z., Kleid, D. & Ptashne, M. (1977a). *Nucl. Acids Res.* **4**, 1595-1607.
- Humayun, Z., Jeffrey, A. & Ptashne, M. (1977b). *J. Mol. Biol.* **112**, 265-277.
- Inouye, M., Arnheim, N. & Sternglanz, R. (1973). *J. Biol. Chem.* **248**, 7247-7252.
- Jaurin, B. & Grundström, T. (1981). *Proc. Nat. Acad. Sci., U.S.A.* **78**, 4897-4901.
- Jaurin, B., Grundström, T., Edlund, T. & Normark, S. (1981). *Nature (London)*, **290**, 221-225.
- Jensen, H. B. & Pryme, I. F. (1974). *Biochem. Biophys. Res. Commun.* **59**, 1117-1123.
- Jensen, K. & Wirth, N. (1978). *Pascal User Manual and Report*, 2nd edit., Springer-Verlag, New York.
- Joachimiak, A., Kelley, R. L., Gunsalus, R. P., Yanofsky, C. & Sigler, P. B. (1983). *Proc. Nat. Acad. Sci., U.S.A.* **80**, 668-672.
- Johnson, A. D., Meyer, B. J. & Ptashne, M. (1978). *Proc. Nat. Acad. Sci., U.S.A.* **75**, 1783-1787.
- Johnson, A. D., Meyer, B. J. & Ptashne, M. (1979). *Proc. Nat. Acad. Sci., U.S.A.* **76**, 5061-5065.
- Johnson, A. D., Poteete, A. R., Lauer, G., Sauer, R. T., Ackers, G. K. & Ptashne, M. (1981). *Nature (London)*, **294**, 217-223.
- Johnson, D. I. & Somerville, R. L. (1983). *J. Bacteriol.* **155**, 49-55.
- Kalnins, A., Otto, K., Rütther, U. & Müller-Hill, B. (1983). *EMBO J.* **2**, 593-597.
- Kanazawa, H., Mabuchi, K., Kayano, T., Tamura, F. & Futai, M. (1981). *Biochem. Biophys. Res. Commun.* **100**, 219-225.
- Kenyon, C. J., Brent, R., Ptashne, M. & Walker, G. C. (1982). *J. Mol. Biol.* **160**, 445-457.
- Kerr, C. & Sadowski, P. D. (1975). *Virology*, **65**, 281-285.
- Kleid, D., Humayun, Z., Jeffrey, A. & Ptashne, M. (1976). *Proc. Nat. Acad. Sci., U.S.A.* **73**, 293-297.
- Kleppe, G., Jensen, H. B. & Pryme, I. F. (1977). *Eur. J. Biochem.* **76**, 317-326.
- Krüger, D. H. & Schroeder, C. (1981). *Microbiol. Rev.* **45**, 9-51.
- Lin, S. & Riggs, A. D. (1975). *Cell*, **4**, 107-111.
- Lipman, D. J. & Maizel, J. (1982). *Nucl. Acids Res.* **10**, 2723-2739.
- Little, J. W. (1983). *J. Mol. Biol.* **167**, 791-808.
- Little, J. W. & Mount, D. W. (1982). *Cell*, **29**, 11-22.
- Little, J. W., Mount, D. W. & Yanisch-Perron, C. R. (1981). *Proc. Nat. Acad. Sci., U.S.A.* **78**, 4199-4203.
- Maas, W. K. & Clark, A. J. (1964). *J. Mol. Biol.* **8**, 365-370.
- Maas, W. K., Maas, R., Wiame, J. M. & Glansdorff, N. (1964). *J. Mol. Biol.* **8**, 359-364.
- Mackie, G. A. (1981). *J. Biol. Chem.* **256**, 8177-8182.
- Maniatis, T., Ptashne, M., Backman, K., Kleid, D., Flashman, S., Jeffrey, A. & Maurer, R. (1975). *Cell*, **5**, 109-113.
- Markham, B. E., Little, J. W. & Mount, D. W. (1981). *Nucl. Acids Res.* **9**, 4149-4161.
- Matthews, B. W., Ohlendorf, D. H., Anderson, W. F., Fisher, R. G. & Takeda, Y. (1983). *Trends Biochem. Sci.* **8**, 25-29.
- McAllister, W. T. & Wu, H. (1978). *Proc. Nat. Acad. Sci., U.S.A.* **75**, 804-808.
- McAllister, W. T., Morris, C., Rosenberg, A. H. & Studier, F. W. (1981). *J. Mol. Biol.* **153**, 527-544.
- Meyer, B. J., Kleid, D. G. & Ptashne, M. (1975). *Proc. Nat. Acad. Sci., U.S.A.* **72**, 4785-4789.
- Meyer, B. J., Maurer, R. & Ptashne, M. (1980). *J. Mol. Biol.* **139**, 163-194.
- Miki, T., Ebina, Y., Kishi, F. & Nakazawa, A. (1981). *Nucl. Acids Res.* **9**, 529-543.
- Miller, G. A. (1955). *Information Theory in Psychology* (Quastler, H., ed.), pp. 95-100, Free Press, Glencoe, Ill.
- Miyazaki, J., Ryo, Y., Fujisawa, H. & Minagawa, T. (1978). *Virology*, **89**, 327-329.
- Morlon, J., Llobes, R., Chartier, M., Bonicel, J. & Lazdunski, C. (1983). *EMBO J.* **2**, 787-789.
- Müller-Hill, B., Gronenborn, B., Kania, J., Schlotmann, M. & Beyreuther, K. (1977). *Nucleic Acid-Protein Recognition* (Vogel, H. J., ed.), p. 219, Academic Press, New York.
- Mulligan, M. E., Hawley, D. K., Enriken, R. & McClure, W. R. (1984). *Nucl. Acids Res.* **12**, 789-800.
- Nei, M. & Li, W. (1979). *Proc. Nat. Acad. Sci., U.S.A.* **76**, 5269-5273.
- Neidhardt, F. C., Vaughn, V., Phillips, T. A. & Bloch, P. L. (1983). *Microbiol. Rev.* **47**, 231-284.
- Newman, A. K., Rubin, R. A., Kim, S. & Modrich, P. (1981). *J. Biol. Chem.* **256**, 2131-2139.
- Nussinov, R. (1984). *Nucl. Acids Res.* **12**, 1749-1763.
- Oakley, J. L. & Coleman, J. E. (1977). *Proc. Nat. Acad. Sci., U.S.A.* **74**, 4266-4270.
- Oakley, J. L., Strothkamp, R. E., Sarris, A. H. & Coleman, J. E. (1979). *Biochemistry*, **18**, 528-537.
- Oppenheim, D. S., Bennett, G. N. & Yanofsky, C. (1980). *J. Mol. Biol.* **144**, 133-142.
- Owen, J. E., Schultz, D. W., Taylor, A. & Smith, G. R. (1983). *J. Mol. Biol.* **165**, 229-248.
- Pabo, C. O. & Sauer, R. T. (1984). *Annu. Rev. Biochem.* **53**, 293-321.
- Pabo, C. O., Krovatin, W., Jeffrey, A. & Sauer, R. T. (1982). *Nature (London)*, **298**, 441-443.
- Pierce, J. R. (1980). *An Introduction to Information Theory: Symbols, Signals and Noise*, 2nd edit., Dover Publications Inc., New York.
- Pingoud, A. (1985). *Eur. J. Biochem.* **147**, 105-109.
- Pribnow, D. (1979). *Biological Regulation and Development* (Goldberger, R. F., ed.), vol. 1, pp. 219-277, Plenum Press, New York.
- Ptashne, M., Backman, K., Humayun, M. Z., Jeffrey, A., Maurer, R., Meyer, B. & Sauer, R. T. (1976). *Science*, **194**, 156-161.
- Ptashne, M., Jeffrey, A., Johnson, A. D., Maurer, R., Meyer, B. J., Pabo, C. O., Roberts, T. M. & Sauer, R. T. (1980). *Cell*, **19**, 1-11.
- Putney, S. D., Meléndez, D. L. & Schimmel, P. R. (1981). *J. Biol. Chem.* **256**, 205-211.
- Reznikoff, W. S., Winter, R. B. & Hurley, C. K. (1974). *Proc. Nat. Acad. Sci., U.S.A.* **71**, 2314-2318.
- Riggs, A. D., Suzuki, H. & Bourgeois, S. (1970). *J. Mol. Biol.* **48**, 67-83.

- Roberts, R. J. (1983). *Nucl. Acids Res.* **11**, r135-r167.
- Sadler, J. R., Sasmor, H. & Betz, J. L. (1983a). *Proc. Nat. Acad. Sci., U.S.A.* **80**, 6785-6789.
- Sadler, J. R., Waterman, M. S. & Smith, T. F. (1983b). *Nucl. Acids Res.* **11**, 2221-2231.
- Sadowski, P. D. & Kerr, C. (1970). *J. Virol.* **6**, 149-155.
- Sampson, J. R. (1976). *Adaptive Information Processing*, Springer-Verlag, New York.
- Sancar, A., Sancar, G. B., Rupp, W. D., Little, J. W. & Mount, D. W. (1982a). *Nature (London)*, **298**, 96-98.
- Sancar, G. B., Sancar, A., Little, J. W. & Rupp, W. D. (1982b). *Cell*, **28**, 523-530.
- Schneider, T. D. (1984). Ph.D. thesis, University of Colorado.
- Schneider, T. D., Stormo, G. D., Haemer, J. S. & Gold, L. (1982). *Nucl. Acids Res.* **10**, 3013-3024.
- Schneider, T. D., Stormo, G. D., Yarus, M. A. & Gold, L. (1984). *Nucl. Acids Res.* **12**, 129-140.
- Shalloway, D., Kleinberger, T. & Livingston, D. M. (1980). *Cell*, **20**, 411-422.
- Shannon, C. E. (1948). *Bell System Tech. J.* **27**, 379-423, 623-656.
- Shannon, C. E. (1951). *Bell System Tech. J.* **30**, 50-64.
- Shannon, C. E. & Weaver, W. (1949). *The Mathematical Theory of Communication*, University of Illinois Press, Urbana.
- Shine, J. & Dalgarno, L. (1974). *Proc. Nat. Acad. Sci., U.S.A.* **71**, 1342-1346.
- Silberstein, S., Inouye, M. & Studier, F. W. (1975). *J. Mol. Biol.* **96**, 1-11.
- Singh, J. (1966). *Great Ideas in Information Theory, Language and Cybernetics*, Dover Publications, Inc., New York.
- Singleton, C. K., Roeder, W. D., Bogosian, G., Somerville, R. L. & Weith, H. L. (1980). *Nucl. Acids Res.* **8**, 1551-1560.
- Smith, H. O. (1979). *Science*, **205**, 455-462.
- Stormo, G. D. (1986). In *Maximizing Gene Expression*, (Gold, L. & Reznikoff, W., eds), pp. 270-328, Benjamin/Cummings Publishing Co., Inc.
- Stormo, G. D., Schneider, T. D., Gold, L. & Ehrenfeucht, A. (1982a). *Nucl. Acids Res.* **10**, 2997-3011.
- Stormo, G. D., Schneider, T. D. & Gold, L. M. (1982b). *Nucl. Acids Res.* **10**, 2971-2996.
- Studier, F. W. (1969). *Virology*, **39**, 562-574.
- Studier, F. W. (1972). *Science*, **176**, 367-376.
- Summers, W. C. & Siegel, R. B. (1970). *Nature (London)*, **228**, 1160-1162.
- Swift, G., McCarthy, B. J. & Heffron, F. (1981). *Mol. Gen. Genet.* **181**, 441-447.
- Tribus, M. & McIrvine, E. C. (1971). *Scient. Amer.* **225** (Sept.), 179-188.
- Uhlin, B. E., Völkert, M. R., Clark, A. J., Sancar, A. & Rupp, W. D. (1982). *Mol. Gen. Genet.* **185**, 251-254.
- van den Elzen, P. J. M., Maat, J., Walters, H. H. B., Veltkamp, E. & Nijkamp, H. J. J. (1982). *Nucl. Acids Res.* **10**, 1913-1928.
- Völker, T. A., Gafner, J., Bickle, T. A. & Showe, M. K. (1982). *J. Mol. Biol.* **161**, 479-489.
- von Hippel, P. H. (1979). *Biological Regulation and Development* (Goldberger, R. F., ed.), vol. 1, pp. 279-347, Plenum Press, New York.
- Warner, J. R. (1979). *DIGRAF: Device Independent Graphics from FORTRAN, User's Guide Version 2.0*, Graphics Development Group, University Computing Center, University of Colorado, Boulder.
- Weaver, W. (1949). *Scient. Amer.* **181**, 11-15.
- Wiberg, J. S. & Karam, J. D. (1983). In *Bacteriophage T4* (Mathews, C. K., Kutter, E. M., Mosig, G. & Berget, P. B., eds), pp. 193-201, American Society for Microbiology, Washington, D.C.
- Winter, R. B. & von Hippel, P. H. (1981). *Biochemistry*, **20**, 6948-6960.
- Yokota, T., Sugisaki, H., Takanami, M. & Kaziro, Y. (1980). *Gene*, **12**, 25-31.
- Young, I. G., Rogers, B. L., Campbell, H. D., Jaworoski, A. & Shaw, D. C. (1981). *Eur. J. Biochem.* **116**, 165-170.
- Zavriev, S. K. & Shemyakin, M. F. (1982). *Nat. Acids Res.* **10**, 1635-1652.
- Zolg, J. W. & Hänggi, U. J. (1981). *Nucl. Acids Res.* **9**, 697-710.
- Zurawski, G., Gunsalus, R. P., Brown, K. D. & Yanofsky, C. (1981). *J. Mol. Biol.* **145**, 47-73.