

Atypical regions in large genomic DNA sequences

(DNA sequence analysis/Markov chain/chromosome structure)

STEWART SCHERER*^{†‡}, MARY SARA MCPEEK^{§¶}, AND TERENCE P. SPEED[§]

*Human Genome Center, Lawrence Berkeley Laboratory 74-157, Berkeley, CA 94720; [†]Department of Microbiology, Box 196, University of Minnesota School of Medicine, Minneapolis, MN 55455; and [§]Department of Statistics, University of California, Berkeley, CA 94720

Communicated by Ronald W. Davis, March 30, 1994

ABSTRACT Large genomic DNA sequences contain regions with distinctive patterns of sequence organization. We describe a method using logarithms of probabilities based on seventh-order Markov chains to rapidly identify genomic sequences that do not resemble models of genome organization built from compilations of octanucleotide usage. Data bases have been constructed from *Escherichia coli* and *Saccharomyces cerevisiae* DNA sequences of >1000 nt and human sequences of >10,000 nt. Atypical genes and clusters of genes have been located in bacteriophage, yeast, and primate DNA sequences. We consider criteria for statistical significance of the results, offer possible explanations for the observed variation in genome organization, and give additional applications of these methods in DNA sequence analysis.

Large contiguous genomic DNA sequences have been determined for a number of species. The Human Genome Project and similar efforts for important model systems will be producing such sequences at an increasing rate. The most widely used computational technique for the analysis of new sequence information is the comparison of the sequence with all other known sequences for regions of similarity. With the growth of the public data bases, this approach has become increasingly successful, with about one-third of all large open reading frames in new genomic sequences or new cDNA sequences being related to known genes (1, 2). The genome projects will have a major impact on this approach because a decreasing fraction of database entries will be sequences of known function with the wealth of biological information associated with them.

Existing computational tools for the analysis of DNA sequences were generally developed for the analysis of individual genes. Most sequence analysis software reports tables of regions of interest based on a scoring system. This tabular presentation becomes unwieldy when faced with sequences of tens or hundreds of kilobases.

Statistical analyses of DNA sequences have been used since the origins of molecular biology. It was first noted that nearest-neighbor frequencies of different species varied greatly and that these differences could not be explained by the known variation in base composition (3). With the advent of efficient DNA sequencing techniques, it became clear that there were wide variations in codon usage (4), tetranucleotide (5), and hexanucleotide (6) frequencies among species. The increase in available sequence data as a result of the genome projects will permit the statistical evaluation of frequencies of longer oligonucleotides.

Building on these observations, we have developed a flexible approach to the analysis of large genomic DNA sequences based on Markov chain models. A data base is built from a large collection of sequences and the frequency of occurrence, termed usage, of all possible sequences of

length 8 is tabulated. This length was chosen for several reasons. First, it is the largest size for which effective statistics could be developed when the complete genomes of the important model microorganisms become available. Second, it is about the smallest size where protein–DNA interactions might be reflected. Finally, this size represents 16 bits of information and offers certain computational economies. Probabilities are calculated for a window of constant length at every possible position along a query sequence and the results are presented graphically. The query can be tested against a variety of models constructed by the selection of different sets of sequences for the underlying data base. The display indicates how well regions of the sequence resemble the model.

Our statistical model assumes stationary behavior of the query sequence but makes no biological assumptions about the nature of that sequence or those used to build the data base. For example, coding and noncoding regions are treated equivalently. For long genomic sequences, one would expect substantial local statistical variations. We discuss the validity of these assumptions and criteria for the determination of significance of the extreme values that are observed.

By building a data base from all known sequences from a single species, we have found regions in large query sequences with atypical organization in both prokaryotes and eukaryotes. These results agree well with our understanding of DNA organization of these regions and highlight areas for future biological investigations. We discuss possible explanations for the presence of genes and gene clusters with unusual organization and other applications of this approach to DNA sequence analysis.

METHODS

DNA Sequence Data. All DNA sequences used are derived from those in GenBank release 75. Compilations of sequences from a particular species used the name in the organism field of the entry. Sequences less than 1000 bases in length (10,000 for human DNA) were not included and the reverse complement of each sequence was added to the compilations.

Octamer Usage Statistics. The usage of all octamers $f(b_1, b_2, \dots, b_8)$ (where b is any base and f is the tabulated count) and total usage $N = \sum_{\text{all 8-mers}} f$ were determined. From these data, conditional probabilities were calculated: $P(b_8|b_1, \dots, b_7) = [f(b_1, b_2, \dots, b_8) + k] / \sum_{b_8} [f(b_1, b_2, \dots, b_8) + k]$ where $k = 0.01 N / 4^8$. This flattening constant k is approximately unity for the data bases used below that were built from large collections of genomic sequences. See ref. 7 for a discussion of this topic. Translated base 2 logarithms of these values were determined as $Z_i = \log_2 P(b_i|b_{i-7}, \dots, b_{i-1}) + 2$. Dinucleotide and trinucleotide frequencies were determined

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.

[‡]To whom reprint requests should be addressed.

[¶]Present address: Department of Statistics, University of Chicago, Chicago, IL 60637.

from the octamer usage. These values were similarly zero adjusted and used to determine values of Z .

Confidence Limits. DNA sequences of a length equal to that of the query sequence were generated randomly by using the usage of heptamers in the data base to provide an initial point and using the conditional probabilities described above. One hundred such sequences were used to establish 95% confidence limits for the upper limit of the plots.

Graphical Presentation. Scores for sequence blocks of length $j + 1$ were calculated for all values of i where $i + j$ is less than the sequence length; $y_{i,j} = -\sum_{k=i}^{i+j} Z_k$ was plotted at $x = i + (j/2)$. The y axis was divided by the mean $\mu_j = [(j + 1)/N] \sum_{\text{all } 8\text{-mers}} Z_f$, so that the values presented are multiples of μ . The entire sequence is plotted in a single frame with the x axis being the position on the sequence. Because of variation in the lengths of the query sequences, widths of the peaks for similar sized structures vary between figures. Note that the plots are $-\log_2 P$. Therefore, sequences that do not resemble the underlying model will point upwards and those overrepresented in the data base will point downwards. Three horizontal lines are shown. The upper line is the position where all values would plot if the data base had a uniform distribution of all sequences. The central line is the mean of the distribution, μ_j . If the sequence resembled the underlying model, one would expect a noise band centered on this position. The third line is twice the mean.

RESULTS

All figures in this work present the comparison of a DNA sequence with an octamer usage data base constructed from genomic DNA sequences from a particular species. Generally all known sequences from a given species were used as described in *Methods*. In some cases, the counts were tabulated from a single DNA sequence. Probabilities were determined for all subsequences of a given length (512 or 1024 bases) in the query sequence and values for each position along the sequence were plotted as described in *Methods*.

Bacteriophage λ . Among large bacterial DNA sequences, phage λ is one of the most intensively characterized (8). In particular, it exhibits a modular organization with adjacent genes often having related functions (9). Fig. 1 presents an analysis with the entire phage genome as query and *Escherichia coli* sequences used to construct the model. Fig. 1*a* uses a seventh-order Markov chain analysis, whereas Fig. 1*b* uses first-order Markov statistics (i.e., dinucleotide frequencies). While the overall pattern is similar, it is clear that the first-order chain, which is built by using nearest-neighbor frequencies, does not effectively reveal the block structure present in the phage genome. Several features stand out in the plot. First, the late genes A–J clearly are representative of typical *E. coli* sequences. The first significant peak near the end of the late transcript is the *lom* gene (nt 18,695–19,582; ref. 10). This gene is transcribed during lytic growth but is not essential for plaque formation. Other regions that produce the downward peaks in the right half of the genome include the recombination, replication, and lysis genes. The regions not essential for lytic growth, most notably the b2 region (leftward from 27,731), deviate most from the model.

Saccharomyces Chromosome 3. This sequence was the first complete eukaryotic chromosome to be determined (2). Because of the extreme length of the yeast sequences being tested, the size of the window was doubled to assure statistical significance. This also makes the remaining larger features easier to detect given the compression of the x -axis. Fig. 2*a* shows the pattern produced using the entire chromosome as query with the model built from all other yeast sequences of >1 kb. While the extent of variation from the model is not as extensive as that observed with phage λ , there is a large region centered about coordinate 210,000 that does not

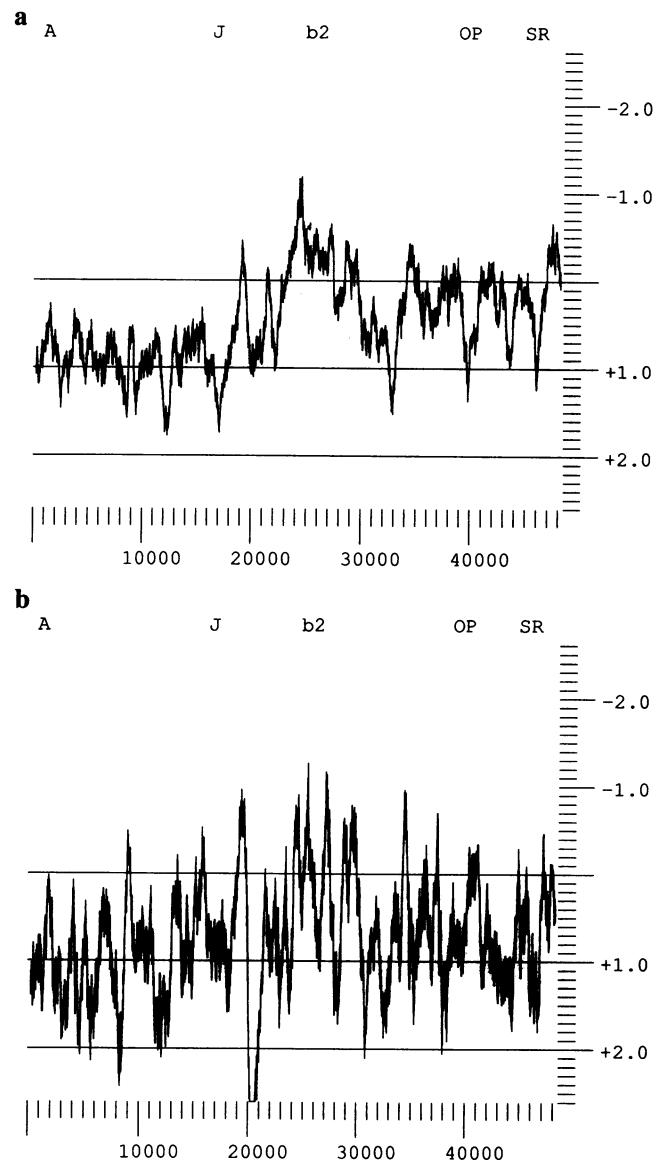


Fig. 1. Seventh-order (*a*) and first-order (*b*) Markov chain analysis of the bacteriophage λ sequence. The block size was 512 nt. In *a*, the 95% confidence limit for significant deviation from the model is -0.28 . Locations of λ markers are given at the top of each panel.

strongly resemble yeast sequences. Several other statistically significant variations are observed. Notable is the lowest point in the graph, which is at the location of a Ty element. This is expected for sequences that are highly represented in the data base.

When yeast chromosome 3 is tested against an *E. coli* model, no part of the sequence strongly resembles bacterial DNA. The mean of the noise band is at -0.55 and none of the low points on this plot (not shown) align with the peaks in Fig. 2*a* or *b*.

Fig. 2*b* is a plot where the data base is built solely from the actual sequence of yeast chromosome 3. Most of the sequence shows no significant variation from other sequences on the chromosome. Several features are of interest. The lowest point on the plot is again the Ty element. The other three pronounced downward peaks are *HML* (12,000), *MAT* (199,000), and *HMR* (292,000). This indicates that any sequence representing about 1% of a data base will be readily detected without resorting to a conventional homology search. The two highest points are the region around 210,000 which was found not to resemble yeast sequences in general and the glucokinase gene (50,000).

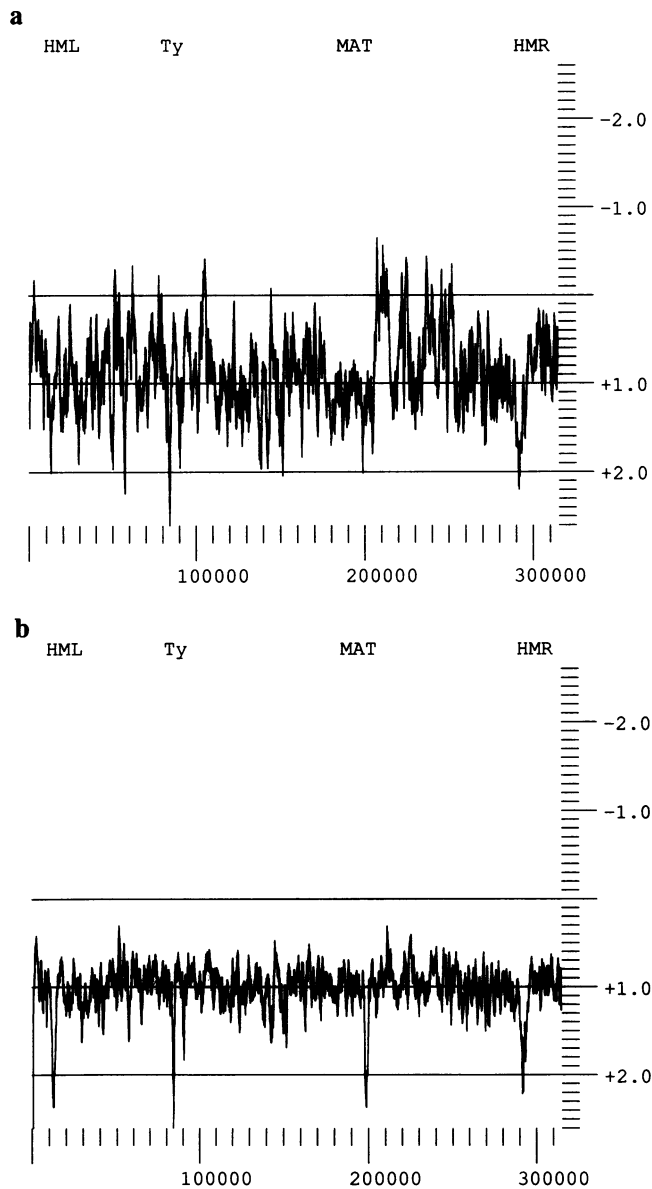


FIG. 2. Analysis of yeast chromosome 3 with the yeast model (*a*) or a yeast chromosome 3 model (*b*). The block size was 1024 nt. The 95% confidence limit for the upper bound in *a* is -0.15 . The locations of several chromosome 3 landmarks are shown at the top of each panel.

Fig. 3 shows an enlargement of the region near coordinate 210,000 on chromosome 3. The window size has been reduced to show smaller features of the plot. Note that the most atypical regions align well with the largest open reading frames of the region; however, many large open reading frames nearby, such as *TSM1* (201,000–204,000) and *THR4* (216,000), show more typical organization.

Saccharomyces Chromosome 1. Among several large yeast genomic sequences examined by this approach, including a total of about 100 kb from chromosome 5, the most significant variation from the general yeast model was noted for this sequence from chromosome 1 (Fig. 4). The highest peak in the plot again aligns with an open reading frame. While the function of this gene in the region of this peak is not completely understood, it is homologous to the hamster *RCC1* locus (11). Given the role of this gene in the cell cycle, one would expect it to have some type of yeast counterpart; however, it is unclear why such a gene would require atypical organization.

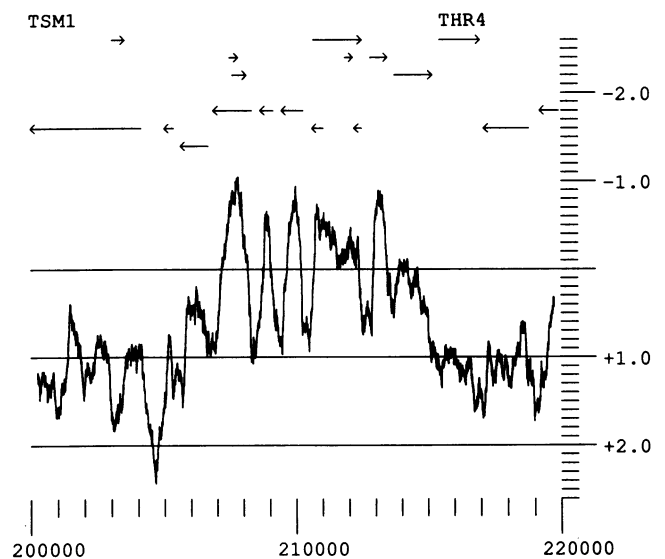


FIG. 3. An enlargement of part of yeast chromosome 3 with the block size reduced to 512 nt. Open reading frames (>100 amino acids starting with methionine) are shown at top with arrows indicating direction of translation. Note the alignment of the peaks with the largest open reading frame in those regions.

β -Globin Genes. Fig. 5 *a* and *b* show seventh-order and second-order analyses of the human β -globin gene cluster (12) using all other human sequences >10 kb to build the data base. This higher cutoff is necessary because a high proportion of the data base is cDNA sequences. The presence of highly repeated DNA sequences such as *Alu* repeats clearly influence the seventh-order results but not the second-order plot. All of the strong downward peaks are at the position of *Alu* elements. When the model is built solely with a related sequence such as the 40-kb *Galago crassicaudatus* (bush baby) β -globin gene cluster (13), the only significant downward peaks are at the regions of homology encoding the β -globin-like genes (data not shown).

Fig. 6*a* shows the results of a seventh-order plot with the bush baby β -globin cluster as query and a human model built from all sequences >10 kb. Note the large region around position 7000 that differs greatly from the human model. A portion of this region was known to be highly homologous to

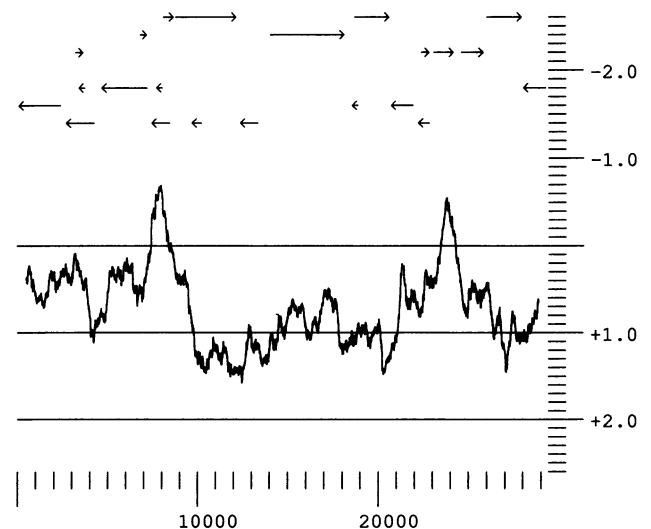


FIG. 4. Analysis of a region of yeast chromosome 1 with the yeast model. The query sequence is GenBank entry YSCLTESPO. The block size and confidence limits are as described in Fig. 2. Open reading frames are shown at top.

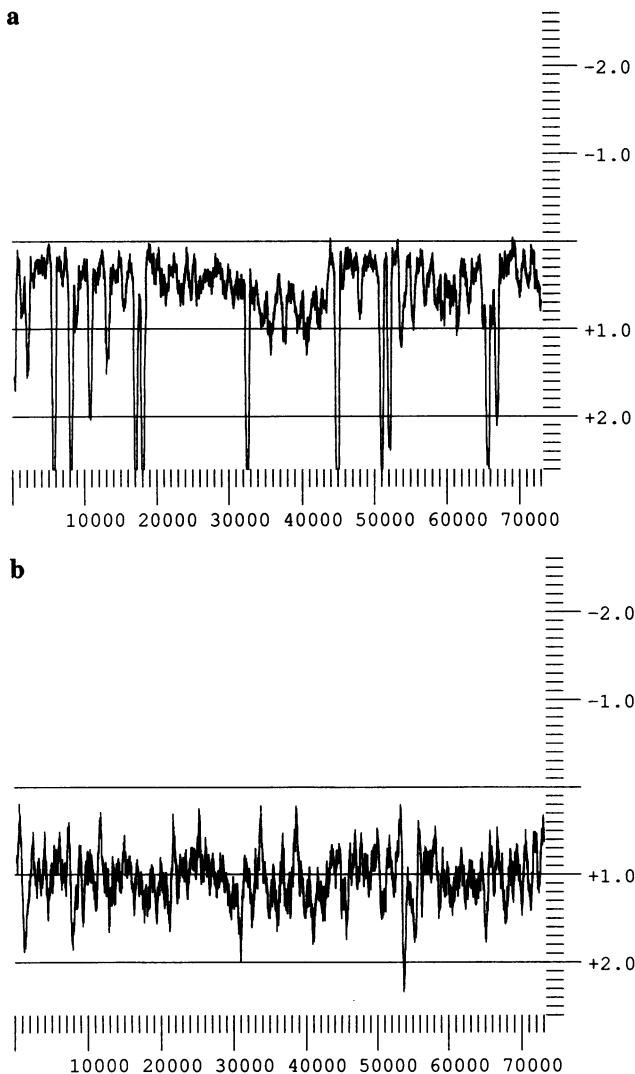


FIG. 5. Human β -globin gene cluster. The query sequence used is HUMHBB with a block size of 512 nt. *a* presents a seventh-order plot; *b* presents a second-order plot. The strong downward peaks in *a* are all at the location of *Alu* elements in the sequence.

the *E. coli* insertion sequence IS186. Our statistical analysis detected a much larger region of unusual organization. The larger region is highly homologous to the *E. coli entD* gene (14) with the left end homologous to IS421. Fig 6*b* shows the same sequence plotted with an *E. coli* model. While the *entD* homology fits well with the *E. coli* model, the other peak in Fig. 6*a* does not. This region (near 32,000) contains a segment homologous to the *E. coli lac* operon.

DISCUSSION

We have developed a simple method, using seventh-order Markov chains, to locate genes and clusters of genes in bacterial and yeast DNA sequences that are quite dissimilar from other sequences in those species. Our approach makes no assumptions about the nature of the sequences being tested and can scan tens of kilobases per second. Similar analyses can also reveal repeated elements and rapidly locate regions of homology between two sequences. This approach should permit the rapid location of interesting regions in large DNA sequences for subsequent computational analyses and biological experimentation.

If a stationary seventh-order Markov model were appropriate for the query sequences, then for most subsequences

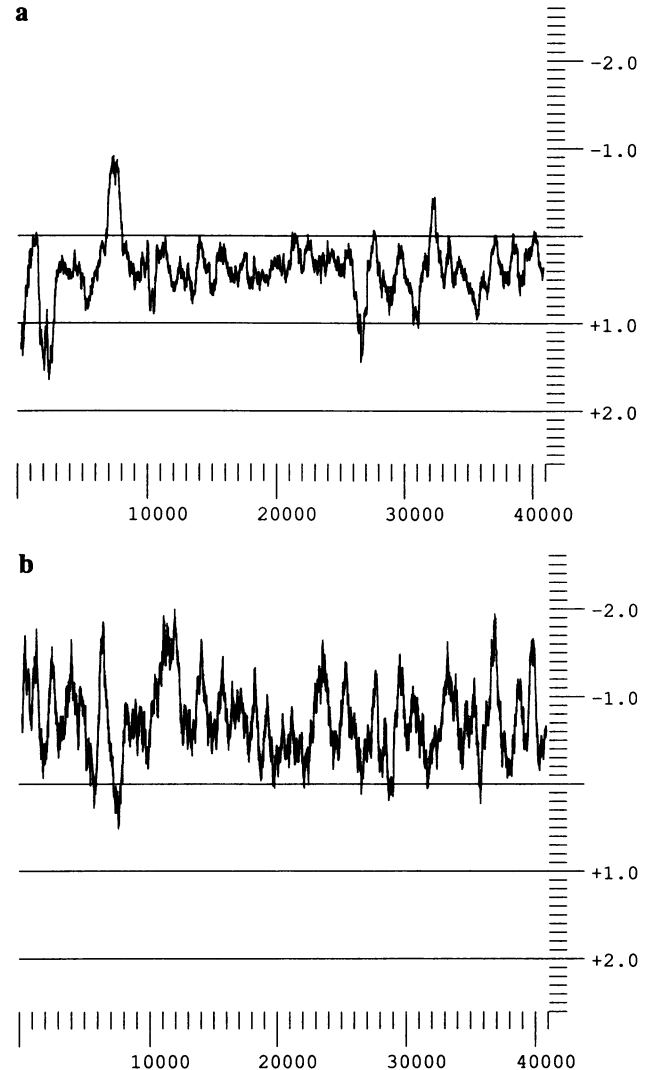


FIG. 6. Bush baby β -globin gene cluster. The query sequence is GCRHBEGEB with a block size of 512 nt. A seventh-order Markov chain analysis is shown using the human (*a*) or *E. coli* (*b*) model. Both of the upward peaks in *a* contain sequences homologous to *E. coli* DNA sequences.

of length $j + 1$ the value $-\log_2 P(b_i, b_{i+1}, \dots, b_{i+j})$ should be approximately $(j + 1)H$, where H is the entropy of the Markov chain. Our analysis differs in that our $-\log_2 P$ values are derived from a data base of other sequences, zero adjusted, and conditional on $b_1 \dots b_7$. The evidence we have accumulated argues strongly that the seventh-order Markov model does not hold globally and forms the basis for identifying biologically interesting regions in large sequences.

We have established our confidence limits by using a simulation procedure (see *Methods*). Other efforts to determine the probability of the plots crossing any value assumed a normal distribution for the values being plotted because each would represent the sum of many identically distributed random variables. $Q-Q$ plots (not shown) of these values for large sequences clearly indicate that this assumption is not correct and that the distributions of the plotted values lack the tails that would expand the range of values that might be expected. The variance determined by fitting a line to these plots was used in efforts to calculate the probability of extreme values arising in Markov chains of the length under study. These latter values gave 95% confidence limits far outside the range of the plots shown. While it is important not to equate biological and statistical significance (15), the

confidence limits determined by the simulation procedure are quite effective in identifying interesting regions in the plots.

Markov chains have been used extensively to study DNA sequences (16, 17), particularly with shorter sequences. While some features mentioned above would have been revealed through analysis of base composition and nearest-neighbor frequencies, much additional useful information is present in the higher-order Markov chains. With bacterial sequences, the general outline of the plot emerges in a second-order chain, reflecting the high fraction of coding sequence. This is not true with human genomic sequences, as coding sequences do not produce the predominant features of the plot. Fifth-order plots begin to indicate the presence of highly repeated elements such as *Alu* sequences.

The sequence under study is not included when the data base is constructed; however, multiple entries or similar sequences may be present for subsets of the sequence. At present, there is no systematic way to resolve these two possibilities for particular regions of a large DNA sequence. This type of overrepresentation does not significantly affect our conclusions. By use of data bases the size of known yeast and *E. coli* sequences, an extra copy of a sequence will lower its position in the plot by $\approx 0.1 \mu_j$ (data not shown).

There are many possible explanations for why a portion of a chromosome might have drastically different sequence composition from the remainder of the chromosome. Some variation would be expected by chance alone in the stationary model. Unknown selective pressures or structural features of chromosomes may be operating. Additional possibilities include horizontal transmission from an unrelated organism, viral infection, and incorporation of sequences from an organelle into the nuclear genome. The demonstration of transfer of plasmids from bacteria into yeast is one example of the first possibility. Many viruses are known to integrate their DNA into human chromosomes, and fractions of the yeast mitochondrial genome have been detected in nuclear DNA. The methods presented here will readily detect hepatitis B virus in the background of human genomic DNA sequences (data not shown).

The modular structure of bacterial genomes is built on the clustering of genes with related functions and the genetic exchanges that can take place between species. The central (b2) region of phage λ has long been known to have a base composition that is markedly different from the composition of other parts of the phage genome. This region is the part substituted in many λ transducing phages with *E. coli* DNA, making it surprising that this is the part of the genome least like *E. coli* DNA. Similarly, the *lom* gene and the sequences at the extreme right end of λ are expressed in lysogens (10). All of the regions we find different from typical *E. coli* organization are known to be dispensable to the phage. Such regions of the phage appear to be under different selective pressures than the rest of the genome and might be optimized for function in another host.

The sequence of yeast chromosome 3 has been extensively analyzed by a variety of statistical tests. The open reading frames near the glucokinase gene and to the right of *MAT* have been noted to have unusual base compositions in the third position of their codons (18). This agrees well with our findings. While there is much evidence for movement of blocks of genes at the telomeres, much less is known about internal insertions or substitutions in yeast chromosomes.

The results in Fig. 6 indicate that this approach will permit the rapid scanning of large sequences for the presence of vector or other bacterial contaminants that might have been included as a result of some step of the cloning or sequencing. Other reports have indicated the presence of *E. coli* IS elements in eukaryotic database entries as the result of

homology searches (19, 20). The approach taken here does not rely on homology to known sequences to detect suspect regions.

The usage tables need not be constructed from all sequences of an organism. Any subset of adequate size could be used such as transcribed or translated regions.

Our results indicate that caution must be exercised in the construction of data sets for use in gene-finding algorithms. It is clear that several authentic *E. coli* and yeast genes differ significantly from the norm using our statistical measure. Using known genes to build the training set will significantly bias the set of reading frames selected by computer programs toward finding more genes like those we already know.

The highly nonrandom character of DNA sequences greatly complicates efforts at statistically modeling large genomic DNA sequences (21). Our studies provide a simple method to visually display this behavior. Graphical presentation of the results of sequence analysis offers several additional advantages over familiar tabular presentations. For example, the growth in size in public data bases and query sequences causes an increase in the score required for statistical significance in a variety of tests. Even with an incomplete understanding of the statistical behavior of DNA sequences, the presence of scores all along a large query rapidly gives the investigator a feel for background scores. We expect that this approach to the presentation of DNA sequence analysis results will have wide application.

We thank John Mulligan for providing yeast genomic DNA sequences prior to publication. This work was supported by Department of Energy Contract DE-AC03-76SF00098 and by National Science Foundation Grant DMS 9113527 to T.P.S.

- Green, P., Lipman, D., Hillier, L., Waterston, R., States, D. & Claverie, J. M. (1993) *Science* **259**, 1711–1716.
- Oliver, S. G., van der Aart, Q. J., Agostoni-Carbone, M. L., Aigle, M., Alberghina, L., Alexandraki, D., Antoine, G., Anwar, R., Ballesta, J. P. & Benit, P. (1992) *Nature (London)* **357**, 38–46.
- Josse, J., Kaiser, A. D. & Kornberg, A. (1961) *J. Biol. Chem.* **236**, 864–875.
- Wada, K., Wada, Y., Ishibashi, F., Gojobori, T. & Ikemura, T. (1992) *Nucleic Acids Res.* **20**, S2111–S2118.
- Burge, C., Campbell, A. M. & Karlin, S. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 1358–1362.
- Cuticchia, A. J., Ivarie, R. & Arnold, J. (1992) *Nucleic Acids Res.* **20**, 3651–3657.
- Good, I. J. (1965) *The Estimation of Probabilities: An Essay in Modern Bayesian Methods* (MIT Press, Cambridge, MA).
- Daniels, D., Schroeder, J., Szybalski, W., Sanger, F. & Blattner, F. (1983) in *Lambda II*, eds. Hendrix, R. W., Roberts, J. W., Stahl, F. W. & Weisberg, R. A. (Cold Spring Harbor Lab. Press, Plainview, NY), pp. 467–518.
- Campbell, A. & Botstein, D. (1983) in *Lambda II*, eds. Hendrix, R. W., Roberts, J. W., Stahl, F. W. & Weisberg, R. A. (Cold Spring Harbor Lab. Press, Plainview, NY), pp. 365–380.
- Barondess, J. J. & Beckwith, J. M. (1990) *Nature (London)* **346**, 871–874.
- Ouellette, B. F., Clark, M. W., Keng, T., Storms, R. K., Zhong, W., Zeng, B., Fortin, N., Delaney, S., Barton, A. & Kaback, D. B. (1993) *Genome* **36**, 32–42.
- Hardison, R. & Miller, W. (1993) *Mol. Biol. Evol.* **10**, 73–102.
- Tagle, D. A., Stanhope, M. J., Siemieniak, D. R., Benson, P., Goodman, M. & Slightom, J. L. (1992) *Genomics* **13**, 741–760.
- Coderre, P. E. & Earhart, C. F. (1989) *J. Gen. Microbiol.* **135**, 3043–3055.
- Karlin, S. & Brendel, V. (1992) *Science* **257**, 39–49.
- Pesole, G., Prunella, N., Liuni, S., Attimonelli, M. & Saccone, C. (1992) *Nucleic Acids Res.* **20**, 2871–2875.
- Kleffe, J. & Borodovsky, M. (1992) *Comp. Appl. Biosci.* **8**, 433–441.
- Sharp, P. M. & Lloyd, A. T. (1993) *Nucleic Acids Res.* **21**, 179–183.
- Lamperti, E. D., Kittelberger, J. M., Smith, T. F. & Villa-Komaroff, L. (1992) *Nucleic Acids Res.* **20**, 2741–2747.
- Binns, M. (1993) *Nucleic Acids Res.* **21**, 779.
- Karlin, S. & Brendel, V. (1993) *Science* **259**, 677–680.