# Selective amplification of variants of a complex repeating unit in DNA of a crustacean

(evolution/amplified DNA/DNA nucleotide sequence determination/divergence measurements/crabs)

NELWYN T. CHRISTIE* AND DOROTHY M. SKINNER†‡

*University of Tennessee–Oak Ridge Graduate School of Biomedical Sciences; and †Biology Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee 37830

**ABSTRACT** The nucleotide sequence of the repeating unit of a fraction of the highly repetitive DNA of the red crab, *Geryon quinquedens,* is reported. Treatment of total DNA with *Hind*III nuclease produced an 81-base-pair monomer and multimers to the size of an octamer. Several of the multimers contained large amounts of fragments of variant sequences, which cannot easily be explained by random mutation alone. That the alterations were not random was corroborated by divergence measurements made on the distribution of *Hha* I nuclease sites within several multimers. The analyses showed that a fraction of each of them is characterized by 4% divergence, while the amounts of dimer, tetramer, and octamer suggest that they have undergone 2–4 times more divergence than that. These results, coupled with the data on sequence variants that are more prevalent in the dimer, indicate that amplification of divergent repeating units could easily explain enhanced amounts of selected multimers.

The structural diversity of satellites and other highly repeated DNAs suggests that several mechanisms may be responsible for their formation. The organization of both mouse satellite (1) and a (G+C)-rich satellite of the hermit crab (2) is compatible with unequal crossing-over of sister chromatids in highly repeated DNAs as postulated (3). Repetitive DNAs may also arise by a mechanism of saltatory replication (4). Increases in genome size have been correlated with increases of particular fractions of repetitive DNAs in some organisms (5). That satellites may be formed by selective amplification is suggested by the presence of specific satellite sequences in only one or a few species in a group of related species (6, 7). Within the Crustacea, homology of repetitive sequences has been demonstrated for widely divergent species (8). Similarly, sequences homologous to the satellite of the mouse *Mus musculus* are found in related species, *M. cervicolor* and *M. caroli* (9). Retention of sequences over long evolutionary periods might be explained by frequent amplifications of highly repetitive DNA.

The genome of the red crab, *Geryon quinquedens*, while lacking satellites, contains 40% highly repetitive DNA (10). In a subset of this repetitive DNA representing 5% of the genome, the distribution of *Hind*III restriction endonuclease sites cannot be explained by random mutation alone and is indicative of amplification of selected regions of the subset (11). We report here the nucleotide sequence of the repeating unit. Some sequence variants occur too frequently to be explained by random mutation alone and may instead be derived from amplified regions of the genome. The complexity of the repeating unit [81 base pairs (bp)] is considerably greater than that of other crab satellite DNAs (2, 12, 13), with the exception of a (G+C)-rich satellite in the Bermuda land crab and a pair of cryptic satellites in the hermit crab (ref. 14 and unpublished

observations). The presence of complex repeats in crab DNAs as well as in the satellites of a number of vertebrates (15–18) suggests a universality in the mechanisms of origin for some highly repeated DNAs and possibly also in their functions. The range of sequence complexities observed in these DNAs, coupled with the localization of both simple and complex repetitive DNAs in the same centromeric regions (19), indicates that several mechanisms may be acting to produce a continuum of complexities. A plausible pathway for evolutionary change in the repetitive DNA of *Geryon* is amplification of a complex repeat unit, divergence from that sequence, and subsequent amplification of divergent subsets.

## MATERIALS AND METHODS

DNA was isolated as described (10). Restriction endonuclease fragments were isolated and sized on polyacrylamide gels, using fragments produced by digestion of phage φX174 DNA with *Hin*fI as size markers. This procedure gave a more accurate size of the 81-bp fragment previously sized as 75 bp by comparison to dye markers (10). After treatment with bacterial alkaline phosphatase, fragments were labeled at the 5′ end (11).

Nucleotide sequence analyses were performed as described (20, 21).

*Hha* I digestion products of each *Hind*III multimer were labeled at the 5′ end, electrophoresed on 7% polyacrylamide gels, frozen, sliced, and assayed for radioactivity as described (2).

## RESULTS

**Sequences of the Basic Repeat Unit and Organization of the Multimers.** A multimeric series of *Hind*III fragments, which includes sizes from an 81-bp monomer to octamer length, can be obtained by digestion of total *Geryon* DNA (11). In *Hha* I digests of each multimer, the two major products are 47 and 40 bp. These are designated fragments *a* and *b*. Determining the sequences of these two fragments from the monomer, dimer, and tetramer indicated that the same-sized *Hha* I fragment from each multimer had essentially the same sequence. A repeat length of 81 bp was obtained because of the overlap of four nucleotides for the *Hind*III site and the overlap of two nucleotides for the *Hha* I site. Fragment *a* was read from the labeled adenine of the upper strand in Fig. 1; fragment *b*, from the labeled adenine of the lower strand. The sequences in the basic unit are more complex than repeated DNA in several other crustaceans (2, 13, 22) in that they lack short repeating units internal to the 81-bp repeat. Despite the considerable complexity of the sequence, there are several runs of pyrimidines. This DNA reassociates in the correct sequence register; following reassociation of dissociated DNA, fragments of
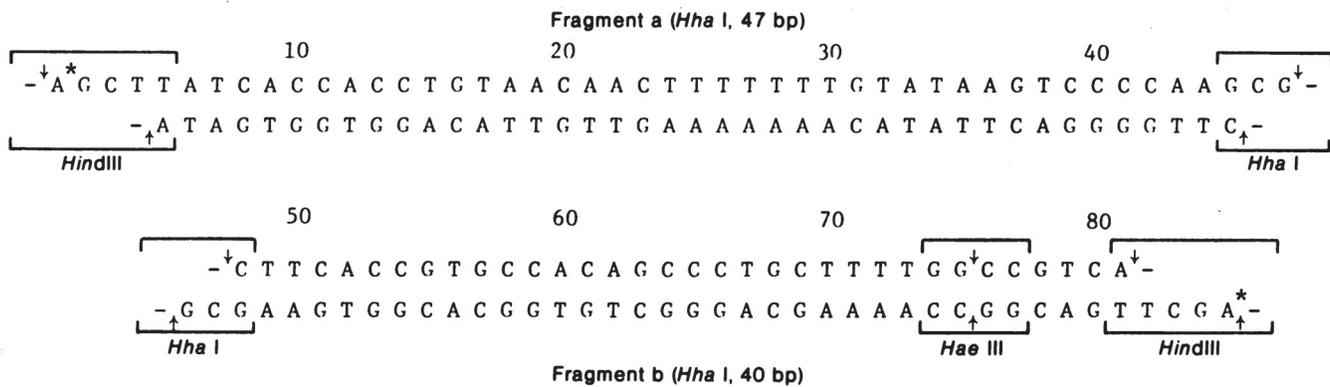
Abbreviation: bp, base pairs.

FIG. 1. Predominant nucleotide sequence of the monomer of the *Hin*dIII multimeric series. Sites for restriction endonuclease cleavage are shown. For sequence analyses performed on *Hha* I digestion products, sequences for each fragment were read from the indicated labeled nucleotide (A*).

monomer and dimer size are produced by *Hin*dIII digestion of both second-order repetitive components but not by digestion of the C$_0$t 10$^{-5}$ (moles of nucleotide per liter × sec) or foldback fraction (10). Minor fragments were produced by *Hha* I digestion of the *Hin*dIII dimer of native DNA (Fig. 2). Among these the 128- and 121-bp fragments are observed in digests of all the multimers that are clearly related to the basic repeating unit of 81 bp (11).

**Sequence Variants Indicating Amplification of Divergent Sequences.** Sequence analyses of individual multimers (Fig. 3) indicated that fragments of each size class are heterogeneous. Differences between the monomer and dimer indicated the presence of one or more sequence subsets either more abundant or unique to the dimer. A sequence variation was found at position 26 of the dimer (Fig. 3B). A G residue occurred in a minor sequence coincident with the T residue found in the major sequence. Whereas there was background radioactivity in the G+A channel for all the Ts between positions 24 and 30, a band in the G channel occurred only at position 26. Although this sequence variant has not been detected in the sequences of the monomer (Fig. 3C), it might be present in small amounts. The same sequence variant was detected in four different sets of sequencing reactions performed on three different preparations of dimers. That background fragments contribute significantly to the sequences seems improbable. The base-specific reactions were performed on *Hin*dIII fragments that had been further selected by digestion with a second enzyme and the bases we attribute to minor variants are distinct and few in number. Finally, the pattern observed in the sequencing gels near the
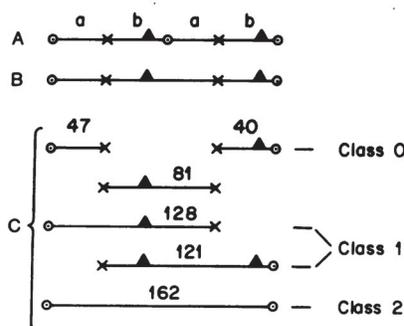


FIG. 2. Organization of sequences in the *Hin*dIII dimer. (A) The original sequence pattern is shown with restriction endonuclease sites given for *Hin*dIII (⊙), *Hha* I (×), and *Hae* III (▲). (B) A representative of the *Hin*dIII dimer fragments that contains all the original sites except the central *Hin*dIII site is presented. (C) *Hha* I digestion of
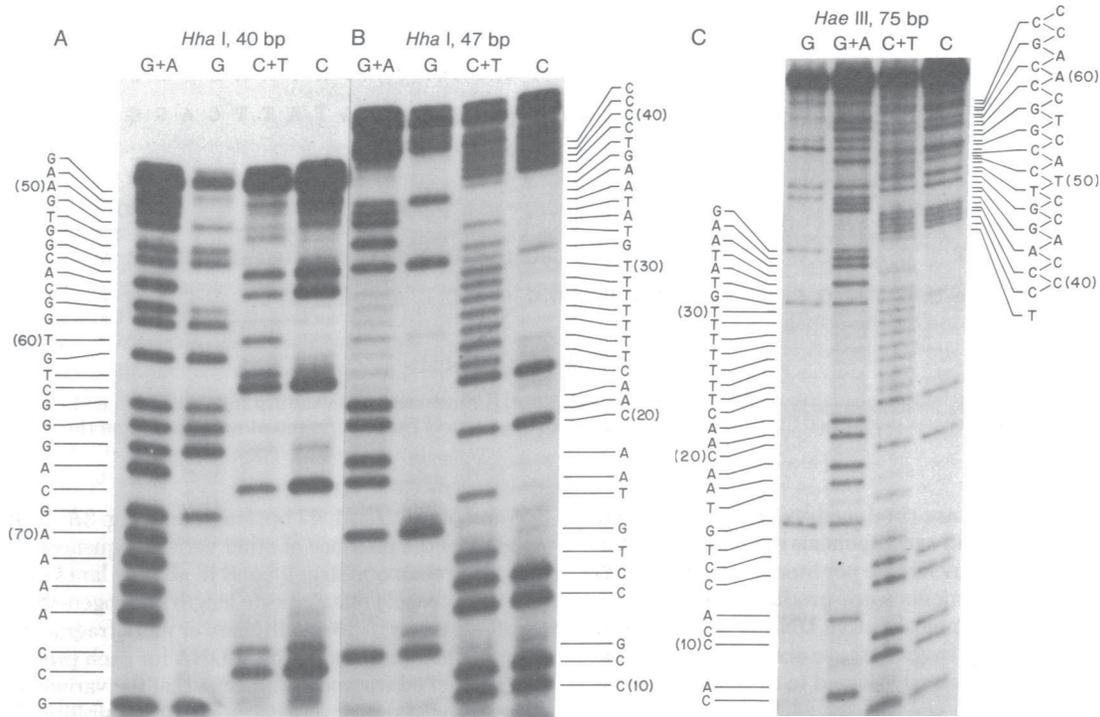
3′ end of the *Hha* I, 47 bp, fragment (Fig. 3B) from the dimer indicates the presence of other variant sequences. The presence of more than one strong band in all four lanes near the 3′ end of the molecule may indicate length heterogeneity; the observed pattern is compatible with three or more fragments of different sizes. The relative amount of DNA for each presumptive band of uncleaved fragments indicates that the variant sequences are not trivial fractions. Because of the possibility of artifacts in sequencing methodologies (21), our assessment of the sequence variants in the dimer is based on repeated observations in numerous sequencing gels utilizing separately prepared reaction sets.

**Sequence Variants Consistent with Random Mutation.** Several minor fragments produced by cleavage at sites not found in the basic repeat can be explained by random mutation. Completely random alterations at every position in the entire array of repeating units would appear only as background in the sequence ladders of autoradiograms. Thus, although random mutational alterations cannot be detected by sequence analyses they are detectable by the appearance of "new" restriction sites within the basic repeat unit. As expected, these arise most often by single nucleotide changes. As long as the proportion of fragments containing these "new" sites is only a few percent of the total sequences, their presence can be readily explained by random mutation.

The site for *Alu* I is the four central nucleotides of that for *Hin*dIII, and there are no additional *Alu* I sites located within the basic repeat. Nevertheless, treatment of the dimer with *Alu* I produced several fragments (Fig. 4, lane 1), which made up 5% of the dimer DNA. The sequences of the *Alu* I, 45 bp, fragment appear identical to those of the *Hha* I, 47 bp, fragment. The change of a single base (G to T) at position 47 would produce the *Alu* I recognition site of (5′)A-G↓C-T(3′) and yield a 45-bp fragment. Sequence analyses of the *Alu* I, 83 bp, fragment from the dimer indicated that this class is a mixture of at least two fragments. Cleavage at the *Alu* I site in the center of the dimer would produce two labeled fragments of 83 bp, one from each 5′ terminus. Random mutation might easily have produced a fraction of dimer fragments that are missing the middle *Hin*dIII site by alteration of either the first or last nucleotide in the site.

The presence of variants of each multimer was further indicated by digestion with *Hae* III, whose site is at positions 74–77 of the basic repeat. The dimer presumably contained two *Hae* III sites and should, when digested by *Hae* III, give fragments corresponding to the presence of either one or both sites.

FIG. 3.   Representative autoradiograms of DNA sequencing gels. Maxam–Gilbert analyses for fragments *Hha* I, 40 bp (*A*), *Hha* I, 47 bp, from the *Hin*dIII dimer (*B*), and *Hae* III, 75 bp, from the monomer (*C*) were applied to 12% polyacrylamide/urea gels.

23, and 28 bp; Fig. 4, lane 2) that represented only a few percent of the digest. These would occur if there had been a change in only a single base at three positions of fragment *b* of the basic repeating unit. Changes to G at positions 64 and 56 would give 21- and 28-bp fragments, respectively. A similar change at position 62 would give a 22-bp fragment and could account for the minor fragment of 23 bp. The proposed changes in the basic repeat unit in the three positions mentioned would give rise to the fragments observed, and these changes appear to be the most likely interpretation of the data at this time.

**Estimate of Divergence in the Multimers.** The extent of the divergence in classes of repetitive DNAs can be determined by examination of the organization of restriction sites (1, 23). This approach indicates the fraction of altered sites in an array of sites originally spaced uniformly. From the fraction of sites that has been altered, the divergence in the entire sequence can

be calculated. Assessment of a set of restriction sites in a repetitive DNA is equivalent to selecting a small population of nucleotides from which to predict the amount of change that has occurred in the entire population. As divergence within the sequences increases, the number of sites altered at consecutive positions will increase. As a result, a multimeric series of fragments will be produced from the original population of monomers. Regardless of the overall extent of divergence, the relative amount of each successively larger multimer should be less than the preceding multimer. Other multimeric series digested from either satellites (1, 24, 25) or total DNA (26, 27) were compatible with formation by random mutation alone. For the *Hin*dIII multimeric series of *Geryon* total DNA, the plot of the logarithm of the relative amount of each multimer versus $(n - 1)$, in which $n = 1$ for a monomer, deviates significantly from the straight line typically observed (1, 11). In particular, there are enhanced amounts of dimer, tetramer, and octamer. These results, coupled with specific indications of amplified divergent sequences in the dimer and in the octamer (see below), suggest that random mutation alone cannot explain the observed distribution of restriction sites.

After the observation of a nonrandom spacing of *Hin*dIII sites in the *Geryon* genome, the distribution of *Hha* I sites within specific *Hin*dIII fragments was investigated. *Hin*dIII fragments were digested with *Hha* I. Fragments were divided into classes[§] according to the number of altered *Hha* I sites. For each multimer, class 0 fragments are those with no altered *Hha* I sites (always the 47- and 40-bp fragments), while classes 1–4 have between one and four consecutive altered sites. In *Hha* I digests of the monomer, class 1 will be the enzyme-resistant fraction, whereas for the dimer it will be fragments of 121 and 128 bp (Fig. 2). The relative amounts of each class are plotted in Fig. 5 for the first four multimers. On each graph there are two broken lines representing two different levels of divergence,
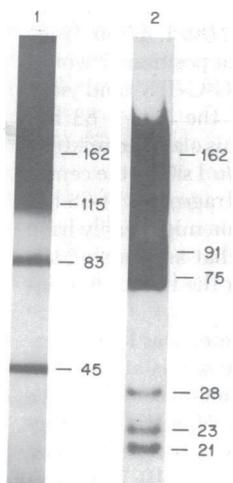


FIG. 4.   Autoradiogram of digestion products of the *Hin*dIII dimer. The dimer was labeled at the 5′ end, digested with *Alu* I (lane 1) or *Hae* III (lane 2), and electrophoresed on a 7% polyacrylamide gel. In

Evolution: Christie and Skinner

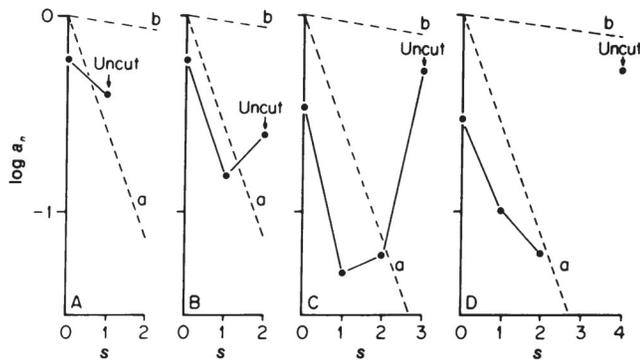Proc. Natl. Acad. Sci. USA 77 (1980)    2789



FIG. 5.  Distribution of *Hha* I sites in *Hin*dIII multimers. The relative amounts of each *Hha* I digestion product of the monomer (*A*), dimer (*B*), trimer (*C*), and tetramer (*D*) were quantified and plotted (solid line) according to the arrangement of *Hha* I sites (see Fig. 2). Classes of fragments are produced in which cleavage occurred at the first through the fourth site from each labeled 5' terminus, depending on the number of consecutive altered sites. The plot of log $a_n$, the fractional amount of each *Hha* I digestion product, versus $s$, the number of consecutive altered sites per fragment, should yield a straight line if the distribution of *Hha* I sites were due entirely to random mutation of a previously uniform array of sites (1, 11). The broken lines are plots of log $c_n$ versus $s$ for divergence values of 3.7% (*a*) and 15.7% (*b*), in which $c_n$ is the number of copies of each multimer relative to the monomer; i.e., $c_n = 1$ when $n = 1$ (11). *Uncut* refers to the enzyme-resistant fraction. For the tetramer, class 3 fragments were not detected.

(*a*) 3.7% and (*b*) 15.7%. These two lines are included to allow a comparison to the distribution of *Hin*dIII sites in the total genome (11). For the *Hin*dIII multimeric series, the relative amounts of trimer, pentamer, hexamer, and heptamer best fit line *a*, whereas the dimer and octamer best fit line *b*. The tetramer value is intermediate. The dimer, tetramer, and octamer compose 0.25, 0.5, and 3.0% of the genome, respectively (11).

Examination of the distribution of *Hha* I sites reveals two points. First, there is always a large amount of enzyme-resistant fragments (indicated as *uncut* in the figure). Although a fraction of this DNA may be unrelated or "background fragments," a similar result was observed for the *Drosophila melanogaster* 1.688 g/cm³ satellite (28). Random inactivation of restriction sites is insufficient to explain completely this organization of restriction sites. Such resistant fractions may well be the result of amplifications of DNA segments in which the original spacing of restriction sites was lost.

The second and more easily interpretable result obtained from these analyses is the indication of the extent of divergence of *Hha* I sites within the *Hin*dIII sequences. A precise determination of divergence of the *Hha* I sites could not be made because of the small number of sites within the first four multimers for which quantitation was performed. Multimers of five, six, and seven monomer lengths were present in too few copies to be used for quantitation. The octamer was not included in this analysis because it contained enhanced amounts of class 3 fragments, which will be discussed later. Only the dimer, trimer, and tetramer will be considered for an estimate of divergence in *Hha* I sites, because an estimate from the monomer would depend on the enzyme-resistant fraction. The relative amounts of each class of fragments for the dimer, trimer, and tetramer are most compatible with the lower divergence of line *a*. These results support the idea that the divergence in at least one subset of the *Hin*dIII fragments is about 4% and that the large amounts of dimer, tetramer, and octamer may be due to selective amplification.
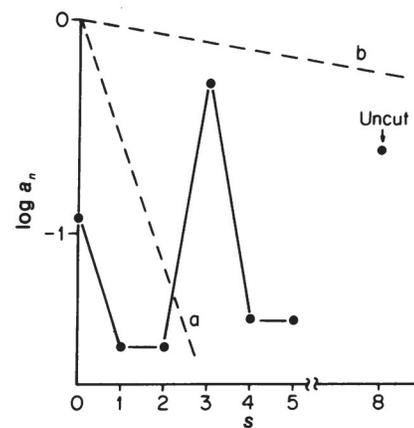


FIG. 6.  Distribution of *Hha* I sites in the *Hin*dIII octamer. The relative amount $a_n$ of each class of *Hha* I digestion products from the octamer is plotted versus $s$, the number of consecutive altered sites, as in Fig. 5. Fragments missing six or seven sites were not detected.

from a quantification of the *Hha* I products of the octamer (Fig. 6). Over 50% of the digest is composed of class 3 fragments, which have presumably lost three consecutive *Hha* I sites. Amplification of a DNA segment containing a divergent octamer would yield this result. Clearly there are also some less divergent octamers, because classes 0, 1, and 2 are also present. From the relative amounts of these latter classes we conclude that the divergence due to random mutation in a fraction of this DNA is only 4%. In addition, the presence of unusually large amounts of particular *Hin*dIII multimers suggests amplification of a DNA segment containing divergent multimers. It is also possible that the large amounts of *Hha* I-resistant fractions in each multimer resulted from amplification of divergent sequences.

## DISCUSSION

The relative proportions of *Hin*dIII multimers in total *Geryon* DNA cannot be explained solely by the random accumulation of mutations in a homogeneous array of monomers (11). The amounts of the trimer, pentamer, hexamer, and possibly the heptamer could result from approximately 4% divergence, whereas the amounts of the dimer, tetramer, and octamer would require 2–4 times that divergence to have been produced by random mutation. The distribution of *Hha* I sites was examined to give an independent assessment of divergence. Their distribution best fits a divergence of 4%, easily distinguishable from a divergence of 16%, and much too low to explain the large amounts of particular multimers observed. Mutation may produce specific changes in the sequence if access to the DNA is limited by chromatin structure (29). In addition, random base changes may produce "hot spots" for further mutation in only a few multimers. Single-nucleotide changes have produced altered mutation rates in the rII locus of phage T4 (30). We conclude that while selective mutation may occur to some extent, selective amplification seems necessary to explain the observed distributions of restriction sites.

Nonrandom distribution of restriction sites has been observed for mouse satellite DNA, although only a small fraction, approximately 5%, of the basic repeat units contain *Hae* III sites (1). For *Geryon*, sequences characterized by nonrandom distribution of restriction sites represent at least 20% of the *Hin*dIII sequences, using the assumption that half of each of the current amounts of dimer, tetramer, and of one of the octamers was produced by amplification of divergent sequences. Although

taxonomic groups are likely. Highly repetitive DNA composes between 25 and 50% of the genomes of several crabs (10, 31, 32). Other arthropods, specifically certain insects, also have large amounts of highly repetitive DNA (7, 33). Although such examples are not considered indicative of a requirement for large amounts of highly repetitive DNAs, they might well suggest that there are active mechanisms for their formation. A rodent, the kangaroo rat, *Dipodomys ordii,* contains over 60% highly repetitive DNA (34). By contrast the genomes of six primates, including humans, contain less than 10% highly repetitive DNA (6).

The homogeneity of repeating units in certain arthropod satellite DNAs should be contrasted with the heterogeneity in rodent satellite DNAs. For guinea pig $\alpha$-satellite or mouse satellite the predominant sequence is representative of only 50% of each DNA (35, 36). For crab poly(dA-dT) (12, 37) and two satellites in the hermit crab (2, 13), the predominant sequence represents more than 90% of the total sequences. This feature may again imply that an unusually high rate of amplification occurs in some repetitive DNAs of crabs. Multiple rounds of amplification at frequent time intervals would have the effect of maintaining a closely related set of sequences. For example, crab poly(dA-dT) appears in widely divergent species (14, 37). Clearly, the inclusion of a divergent multimer in the segment for amplification would increase the magnitude of the divergent sequences. Additional support for this idea is derived from the low divergence observed in a fraction of three *Hin*dIII multimers of *Geryon.* Amplification of both divergent and nondivergent multimers would effectively maintain a portion of the original subset and simultaneously magnify a particular divergent subset. "Divergent" and "nondivergent" refer only to the set of restriction sites being tested; other modifications would be undetected. Because each of three multimers contains fractions showing 4% divergence, it is possible that a single amplification contained a dimer, trimer, and tetramer that had retained the original spacing of *Hha* I sites. This would indicate that the minimum amount of amplified DNA was the combined length of these three multimers (729 bp). The data do not require that these three multimers be contiguous for amplification. Similarly, from the presence of large amounts of a divergent octamer the amplification unit could be at least 648 bp. Whether this latter amplification occurred at the same time as that of the smaller multimers cannot be determined. These estimates of the length for an amplification unit are minimal, because a length of several thousand nucleotides would still be small enough so that the distribution of restriction sites in the amplified DNA might not reflect the distribution in the entire set of *Hin*dIII sequences, approximately $3 \times 10^8$ bp.

1. Southern, E. M. (1975) *J. Mol. Biol.* **94,** 51–69.
2. Chambers, C. A., Schell, M. P. & Skinner, D. M. (1978) *Cell* **13,** 97–110.
3. Smith, G. (1973) *Cold Spring Harbor Symp. Quant. Biol.* **38,** 507–513.
4. Britten, R. J. & Davidson, E. H. (1971) *Q. Rev. Biol.* **56,** 111–138.
5. Bozzoni, J. & Beccari, E. (1978) *Biochim. Biophys. Acta* **520,** 245–252.
6. Jones, K. W. (1977) in *Molecular Structure of Human Chromosomes,* ed. Yunis, J. J. (Academic, New York), pp. 295–326.
7. Barnes, S. R., Webb, D. A. & Dover, G. (1978) *Chromosoma* **67,** 341–363.
8. Graham, D. E. & Skinner, D. M. (1973) *Chromosoma* **40,** 135–152.
9. Rice, N. & Strauss, N. (1972) *Carnegie Inst. Washington Yearb.* **71,** 264–269.
10. Christie, N. T. & Skinner, D. M. (1979) *Nucleic Acids Res.* **6,** 781–796.
11. Christie, N. T. & Skinner, D. M. (1980) *Nucleic Acids Res.* **8,** 279–298.
12. Skinner, D. M. (1967) *Proc. Natl. Acad. Sci. USA* **58,** 103–110.
13. Skinner, D. M., Beattie, W. G., Blattner, F. R., Stark, B. P. & Dahlberg, J. E. (1974) *Biochemistry* **13,** 3930–3937.
14. Beattie, W. G. & Skinner, D. M. (1972) *Biochim. Biophys. Acta* **281,** 169–178.
15. Maio, J. J., Brown, F. L. & Musich, P. R. (1977) *J. Mol. Biol.* **117,** 637–655.
16. Rosenberg, H., Singer, M. & Rosenberg, M. (1978) *Science* **200,** 394–402.
17. Manuelidis, L. & Wu, J. (1978) *Nature (London)* **276,** 92–94.
18. Kopecka, H., Macaya, G., Cortadas, J., Thiéry, J.-P. & Bernardi, G. (1978) *Eur. J. Biochem.* **84,** 189–195.
19. Manuelidis, L. (1978) *Chromosoma* **66,** 23–32.
20. Maxam, A. & Gilbert, W. (1977) *Proc. Natl. Acad. Sci. USA* **74,** 560–564.
21. Maxam, A. & Gilbert, W. (1980) *Methods Enzymol.* **65,** in press.
22. Skinner, D. M. (1977) *BioScience* **27,** 790–796.
23. Slack, J. (1974) *Biopolymers* **13,** 2241–2264.
24. Altenberger, W., Horz, W. & Zachau, H. G. (1977) *Eur. J. Biochem.* **73,** 393–400.
25. Fittler, F. (1977) *Eur. J. Biochem.* **74,** 343–352.
26. Manteuil, S., Hamer, D. & Thomas, C. A., Jr. (1975) *Cell* **5,** 413–422.
27. Cooke, H. J. (1975) *J. Mol. Biol.* **94,** 87–99.
28. Carlson, M. & Brutlag, D. (1977) *Cell* **11,** 371–381.
29. Musich, P. R., Maio, J. J. & Brown, F. L. (1977) *J. Mol. Biol.* **117,** 657–677.
30. Koch, R. E. (1971) *Proc. Natl. Acad. Sci. USA* **68,** 773–776.
31. Vaughn, J. C. (1975) *Chromosoma* **50,** 243–257.
32. Holland, C. A. & Skinner, D. M. (1977) *Chromosoma* **63,** 223–240.
33. Brutlag, D., Appels, R., Dennis, E. S. & Peacock, W. J. (1977) *J. Mol. Biol.* **112,** 31–47.
34. Hatch, F. T. & Mazrimas, J. A. (1970) *Biochim. Biophys. Acta* **224,** 291–294.
35. Southern, E. M. (1970) *Nature (London)* **227,** 794–798.
36. Biro, P. A., Carr-Brown, A., Southern, E. M. & Walker, P. M. B. (1975) *J. Mol. Biol.* **94,** 71–86.
37. Skinner, D. M. & Kerr, M. S. (1971) *Biochemistry* **10,** 1864–1872.