# Human β satellite DNA: Genomic organization and sequence definition of a class of highly repetitive tandem DNA

(human genome/concerted evolution/acrocentric chromosomes)

JOHN S. WAYE AND HUNTINGTON F. WILLARD*

Department of Medical Genetics, University of Toronto, Toronto, ON M5S 1A8 Canada

Communicated by Charles R. Cantor, May 15, 1989 (received for review August 23, 1988)

ABSTRACT    We describe a class of human repetitive DNA, called β satellite, that, at a most fundamental level, exists as tandem arrays of diverged ≈68-base-pair monomer repeat units. The monomer units are organized as distinct subsets, each characterized by a multimeric higher-order repeat unit that is tandemly reiterated and represents a recent unit of amplification. We have cloned, characterized, and determined the sequence of two β satellite higher-order repeat units: one located on chromosome 9, the other on the acrocentric chromosomes (13, 14, 15, 21, and 22) and perhaps other sites in the genome. Analysis by pulsed-field gel electrophoresis reveals that these tandem arrays are localized in large domains (50–300 kilobase pairs) that are marked by restriction fragment length polymorphisms. In total, β satellite sequences comprise several million base pairs of DNA in the human genome. Analysis of this DNA family should permit insights into the nature of chromosome-specific and nonspecific modes of satellite DNA evolution and provide useful tools for probing the molecular organization and concerted evolution of the acrocentric chromosomes.

It has been recognized for nearly 20 years that a considerable proportion of the human genome consists of constitutive heterochromatin, regions known to be rich in highly repetitive "satellite" DNA (1, 2). Recently, considerable progress has been made toward the systematic classification of satellite DNAs based on the size and composition, as well as the genomic distribution and subchromosomal localization, of the tandem repeat unit. Studies employing cloned satellite probes have revealed the existence of two broad classes of human satellite DNA. One class of sequences, comprising the major components of several classical satellite DNAs originally isolated by altered buoyant density (base composition), consists of relatively short oligonucleotide tandem repeat units organized in long chromosome-specific arrays. These sequences are located primarily in the heterochromatic long arm regions of human chromosomes 1, 9, 16, and Y (3–7). A second class of sequences, α satellite, consists of ≈171-base-pair (bp) tandem repeat units localized to the centromeric region of each human chromosome (8). Although both of these major classes of human satellite DNA are organized in a largely chromosome-specific manner, a notable exception involves sequences at the centromere and on the short arms of the five pairs of acrocentric chromosomes (9–12). Repetitive DNA sequences and the ribosomal RNA (rRNA) genes localized to these chromosomes do not exhibit chromosome specificity but rather demonstrate complex patterns of interchromosomal homology, reflecting partial or complete sequence homogenization and concerted evolution (13–18).

In this report, we describe the isolation and characterization of a human satellite DNA family, β satellite,† which is unrelated in structure or sequence to any previously described human satellite DNAs. These sequences comprise a minimum of several million base pairs of DNA in the human genome and are organized, at a fundamental level, as tandem arrays of diverged ≈68-bp monomer repeat units. The β satellite DNA family is subdivided into two or more distinct subsets, at least one of which is shared by the five pairs of acrocentric chromosomes.

## MATERIALS AND METHODS

**Isolation of Cloned Repetitive DNAs.** Human genomic Sau3A fragments less than ≈200 bp in length were isolated from a 2.0% low-melting-point agarose gel and ligated into BamHI-digested pUC9. One recombinant plasmid (pB70-2) contained a 70-bp Sau3A insert that hybridized to highly repeated sequences in genomic DNA and was used to screen a plasmid library constructed from 1.5- to 2.5-kilobase (kb) fragments of Bgl II-digested human genomic DNA cloned in the BamHI site of pUC9.

**DNA Analysis.** Methods for DNA isolation, restriction endonuclease digestion, electrophoresis, transfer to nitrocellulose or Hybond membranes, prehybridization, and hybridization have been described (19). The human/rodent somatic cell hybrid mapping panel has been described in detail elsewhere (20, 21). Conditions of reduced stringency consisted of overnight hybridization at 42°C [3× SSC (1× SSC = 0.15 M NaCl/15 mM sodium citrate)/50% formamide] and washing at 65°C in a final solution containing 0.5 M NaCl. Conditions of high stringency consisted of overnight hybridization at 52–53°C (3× SSC/50% formamide) and washing at 65°C in a final solution of 0.1× SSC/0.1% SDS. Methods for separation of high molecular weight genomic DNAs by pulsed-field gel electrophoresis (PFGE) have been described (22). PFGE filters were provided by S. Kenwrick, M. Patterson, and K. Davies (University of Oxford).

**Nucleotide Sequencing.** The nucleotide sequence of ≈1.5 kb from each of clones pB3 and pB4 was determined by the dideoxy termination method using double-stranded plasmid templates as described (21).

## RESULTS

**Isolation of a Human Repetitive DNA Sequence.** Restriction endonucleases that cleave at regular intervals within a tandem DNA array generate discrete fragments, which, depending on their abundance relative to the bulk of genomic DNA

Genetics: Waye and Willard

Proc. Natl. Acad. Sci. USA 86 (1989)    6251

fragments, may be visualized by ethidium bromide staining (5, 23). Digestion of human genomic DNA with Sau3A generates several discrete, highly repeated fragments <200 bp in length. To determine the nature of these sequences, these fragments were isolated and cloned. One such clone (pB70-2), containing a 70-bp Sau3A insert, detected repetitive sequences in genomic DNA. Two longer clones (designated pB3 and pB4), consisting largely of ≈68-bp Sau3A repeat units and corresponding in size to repetitive Bgl II fragments detected in genomic DNA, were subsequently isolated from a genomic library screened with pB70-2 under conditions of reduced stringency (see Materials and Methods).

**Genomic Organization of Cloned Repeat Units.** As an initial step in the characterization of these sequences, pB3 and pB4 were used as hybridization probes in Southern blot analysis of human DNA digested with various restriction endonucleases. Under conditions of high stringency, pB3 and pB4 detect genomic sequences that have dramatically different hybridization patterns (Fig. 1). For tandemly repeated DNAs, identical restriction periodicities generated by different enzymes can provide an indication of the minimum length of a higher-order tandem repeat unit. Accordingly, the hybridization patterns revealed by pB3 and pB4 (Fig. 1) indicate that the cognate genomic sequences may be organized as tandem arrays of ≈2.5-kb and ≈2.0-kb higher-order repeat units, respectively. Partial digestions of genomic DNA with either EcoRI or Acc I have demonstrated the tandem organization of the genomic sequences detected by these probes (data not shown). Clone pB3 hybridizes to tandemly arranged genomic sequences with EcoRI sites located at ≈2.5-kb intervals. Clone pB4, on the other hand, hybridizes to genomic sequences that are organized as tandem ≈2.0-kb units marked by Acc I sites. The majority of pB3- or pB4-homologous genomic sequences were not cleaved by the enzyme characteristic of the other probe (e.g., Acc I or EcoRI, Fig. 1). Thus, the tandem arrays detected by pB3 and pB4 are not overlapping; rather, these probes appear to recognize independent domains of tandemly repeated sequences.

**Nucleotide Sequence Organization of Higher-Order Repeat Units.** The nucleotide sequence of ≈1.5 kb from each clone was determined to define the internal suborganization of these higher-order repeat units. Examination of the sequence of pB3 and pB4 indicates that each is comprised of diverged, tandemly reiterated units of monomer length ≈68 bp. In total, 33 monomer repeats were sequenced; they are generally punctuated by Sau3A restriction sites, have an average

length of 68 bp (range = 66–70 bp) and an average G + C content of 46% (range = 39–51%), and have no readily apparent sequence redundancy smaller than the monomer. Based on sequence homology, we have aligned 33 monomer sequences (17 from pB3 and 16 from pB4) and derived an "average" monomer sequence for these higher-order repeat units (Fig. 2). A few positions appeared to show a consistent difference between clones (e.g., positions 23, 25, and 31), as might be expected for clones defining different subfamilies.

These data define a human tandem repeat DNA family based on a monomer repeat length of ≈68 bp. These repeat units are organized as longer higher-order repeat units that consist of a characteristic number of diverged monomer units. Since the ≈68-bp monomer sequence bears no resemblance to other reported human satellite DNAs, we suggest that these sequences represent a previously undefined class of repetitive DNA in the human genome that we refer to as "β satellite," in accordance with a proposed nomenclature for primate satellite DNAs (24).

**Long-Range Organization of β Satellite Domains.** The technique of PFGE has been used to investigate the long-range organization of genomic β satellite DNA sequences. In theory, restriction enzymes that do not cleave within the higher-order repeat unit should cleave only flanking sequences, releasing arrays of tandemly repeated units intact.

We have analyzed the β satellite hybridization patterns of genomic DNAs from a small two-generation family, which have been digested with Xba I or EcoRI, and separated by PFGE. Xba I does not cleave within the pB3 or pB4 higher order repeat units (e.g., Fig. 1). PFGE analysis of Xba I-digested genomic DNAs reveals multiple, high molecular weight fragments detected by pB3 (Fig. 3A). A different set of high molecular weight Xba I fragments was detected by using pB4 (Fig. 3C). EcoRI, a noncutter of the pB4 higher-order repeat unit, generates several high molecular weight fragments homologous to pB4 (Fig. 3D). In contrast, most of the genomic EcoRI fragments homologous to pB3 are localized to a low molecular weight band (≈2.5 kb) (Fig. 3B), consistent with the presence of a single EcoRI site within most copies of the pB3 higher-order repeat unit (Fig. 1).

With both of the β satellite probes, multiple high molecular weight fragments were generated (Fig. 3 A, C, and D). This may indicate that there are multiple arrays of a given higher-order repeat unit per haploid genome. Alternatively, there may be a single array that is interrupted by infrequent restriction sites. Interestingly, the pattern of high molecular weight bands is different in different individuals, and the long-range restriction fragment length polymorphisms appear to be inherited in Mendelian fashion (i.e., there are no bands in the offspring that are not represented in one or both of the parents, Fig. 3). The polymorphic Xba I patterns seen with pB3 (Fig. 3A) are consistent with segregation of one major "locus" in the genome. The patterns observed with pB4, however (Fig. 3 C and D), are more complex and may reflect independent segregation of several different loci on different chromosomes.

**Estimated Prevalence of β Satellite Sequences in the Genome.** The copy numbers of the pB3 and pB4 higher-order repeat units were estimated by quantitative dot blot hybridization (data not shown). By using high stringency conditions, we estimate that there are ≈250–500 copies of the ≈2.0-kb higher-order repeat unit (pB4) and significantly lesser amounts (≈50–100 copies) of the ≈2.5-kb higher-order repeat unit (pB3) per haploid genome. These estimates are consistent with independent estimates derived from summing up fragment lengths in PFGE experiments (e.g., Fig. 3).

To estimate the total percentage of human DNA represented by β satellite sequences (pB3- and pB4-homologous sequences, and other less closely related sequences), the above experiment was repeated under less stringent condi-
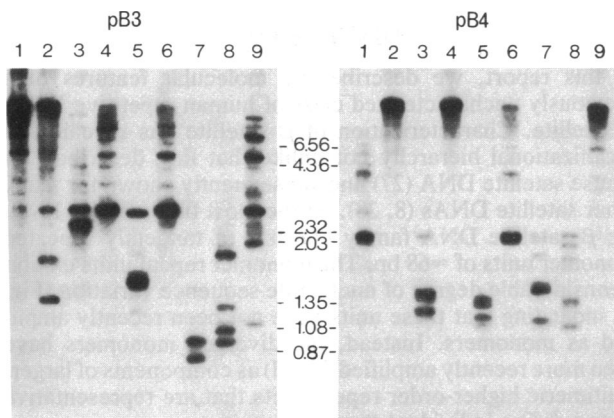


FIG. 1. Organization of genomic DNA sequences homologous to pB3 and pB4. Genomic DNA was digested with different restriction enzymes and hybridized under high stringency conditions to pB3 or pB4 as indicated. Lanes: 1, Acc I; 2, Bgl I; 3, Bgl II; 4, EcoRI; 5, Hae III; 6, Hpa I; 7, Hinf I; 8, Rsa I; 9, Pst I. Sizes are indicated in kb.

```
        1        10        20        30         40        50        60       68
        .         .         .         .    A     C   .         .         .         .
        GATCAGTGCAGAGATATGTCACAA-GCCCC-TGTAGGCAGAGCCTAGACAAGAGTTACATCACCTGGGT
                                    T

pB3     ...T.....................AT....C...........TA...GG......CC............
        ....G......T........T.AC.T.....-...........G.....G.................
        .......T.......G.....A..T.T...A..........TT.TA.......T.-.C.......A....
        ...................T.........T...T.CAAA.A................AG..............A
        ...........A.....C.....T......-A..G......TA-.......TA............T..
        .T........A.TG.T..A.....T......-CA...A....A.......G..C..C...........
        T.........T.....T..T.C.AAG...C..........--.G.T...C........A...AA...
        ...................T.C.AA....C..AC...GA...-..T....T........A.TA.C...
        ......G.....A.......T.C.AAG...-.......C.....A..........T....T..A....
        ................A..T...ATAA..C....A.......A.................C.A
        .................A...TTC.....-..A...C.T.A...GC.............A...
        ........A..G......ACT..T.T...A.....T..CT..............G.C....G.CA.A
        ............GG......A....A-...........-..............TA.....T...T.
        T.........CACT...AA...T..A.AAT.C.......C...A.........AC.C..C..ACA..
        ......G.....A............A....--......C.....A.............TA.T..A-...
        ...G...A....A.........C.AAG...C.C.....CA...........GA..G....T..A-...
        ................C.......TAA...-......T.AGA....................A....

pB4     ......C...TG...........TATG..G..A.A......G.....................-...
        .T.................A.TATTCC..TT-A....-.C.......C.....T.........TAA...
        ............T.......CT.TGT....-.........A...-...G....G...........TT..
        C....C.T..CG.......G.A..C....--...........T.C..T...GA.....C.....A.T.
        ...T....................T...AT-C...........-......G.A.........G....A...
        ................T.....A..TT...--.........T..T.A...T...C...G....AA..
        .....................A..T..-.T.....C..A.T....T--.T...............
        ....T.....A.TG...........A..A..C....C....A.C........TG..........T.....
        .............T..............T...G--....C.....T....G.T.AT.....G........A
        ...........C.....T...TC....-...G.-..C.................A........-...
        T.A.............G..A.GT...TAT...C.......C........C.....T..........TAC...
        ...T..........G..A.GT..CT...TA-C.............-....G.A........G....A...
        .................CA.TG...T.....-...........-.T.G....T.G..T...........
        .............G.....G....G......-.T..A.........C.....T..............A..
        ...........G.C........T.....-.T..A.......T..A..........G.......AC..
        ...............T.C.G....T.....-CA........T..A...G.......CGC-.........
```
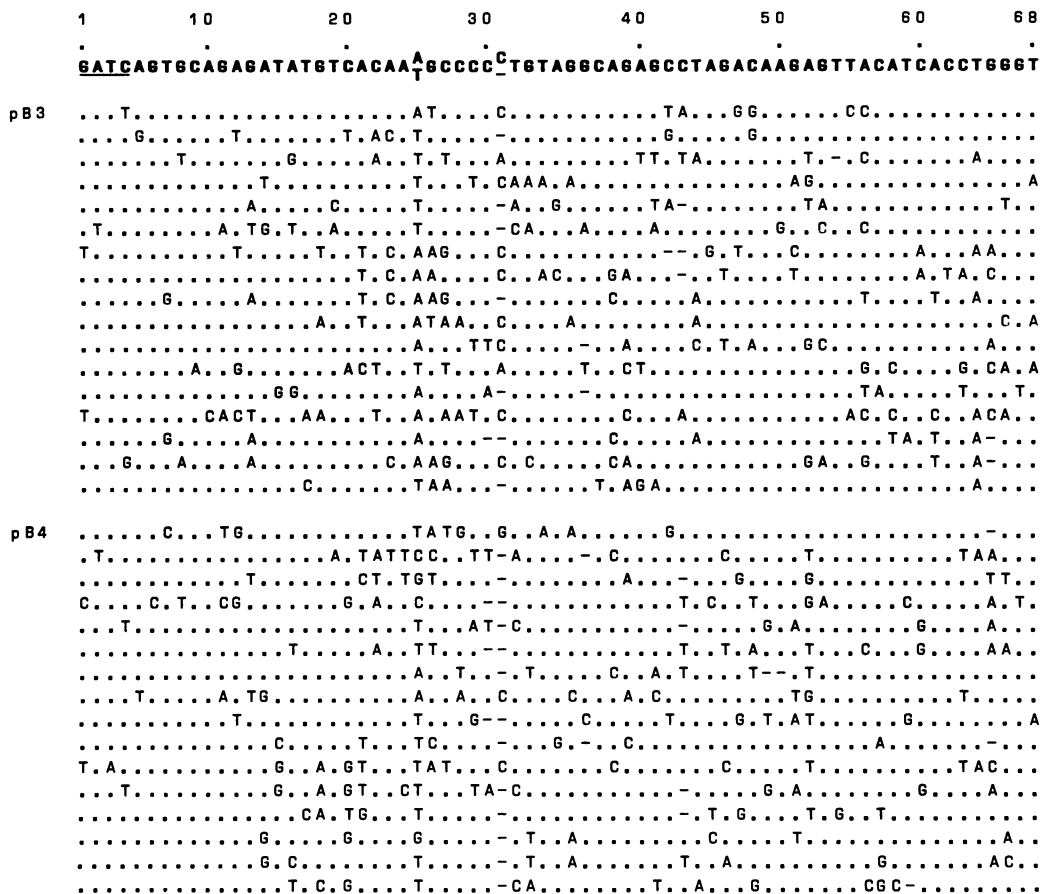
FIG. 2. Derivation of a β satellite average monomer sequence. The sequences of 33 β satellite monomer units (17 from pB3 and 16 from pB4) were aligned for maximum homology. The average sequence (bold) is based on the most abundant base at a given position. A dash (–) indicates a gap introduced to optimize alignment. A period indicates agreement with the average base.

tions that allow substantial sequence mismatch. Under these conditions, the results given by both probes indicate that the related sequences are present in amounts equivalent to 1000–2000 copies of the respective higher-order repeat units (or 30,000–60,000 copies of the 68-bp monomer repeat unit) per haploid genome. Thus, the β satellite DNA family comprises at least several million base pairs of DNA in the human genome, of which only a relatively small proportion corresponds to the specific subsets represented by pB3 and pB4.

**Distribution of β Satellite Sequences in the Human Genome.** To determine the chromosome distribution of the pB3 and pB4 β satellite subsets, we have conducted Southern blot experiments using DNAs of an extensive panel of human/ rodent somatic cell hybrids (20, 21). The genomic sequences corresponding to pB3 (≈2.5-kb EcoRI fragment characteristic of the higher-order repeat unit) were only detected in hybrids that contain human chromosome 9 (Fig. 4A). In a series of 32 hybrids, all other human chromosomes (including the acrocentrics) were discordant in at least five hybrids. Fluorescence in situ hybridization experiments (e.g., ref. 26), using biotinylated pB3 as probe, have localized pB3-homologous sequences to the heterochromatic region of the long arm of chromosome 9 (9qh) (M. Bedford and H.F.W., unpublished data).

The pB4 β satellite subset appears to be more widely distributed and could not be uniquely assigned to a single chromosome (Fig. 4B). Considering only those hybrids that are negative for pB4 at high stringency, it appears that this subset is not present on human chromosome 3, 4, 6, 7, 8, 12, 17, 19, or X. In conjunction with this conclusion, analysis of hybrids containing single (or only a few) human chromosomes indicates that the pB4 subset is located on at least five

human chromosomes, specifically, 13, 14, 15, 21, and 22. Fluorescence in situ hybridization experiments have localized pB4-homologous sequences to the short arms of the acrocentric chromosomes, both proximal and distal to the rRNA gene clusters (M. Bedford and H.F.W., unpublished data). At least in the case of chromosome 21, this has been verified by Southern blot analysis of human/rodent somatic cell hybrids containing only the proximal or distal portion of a chromosome 21 translocation, with a breakpoint through the rRNA gene cluster (18, 25) (Fig. 4B and data not shown).

## DISCUSSION

In this report, we describe the molecular features of a previously uncharacterized class of human repetitive DNA, β satellite. Characterization of β satellite has revealed an organizational hierarchy not unlike that first described for mouse satellite DNA (27) and subsequently shown for many other satellite DNAs (8, 28). At the most fundamental level, the β satellite DNA family consists of tandemly repeated monomer units of ≈68 bp. The monomer repeat units exhibit a considerable degree of nucleotide sequence variation (Fig. 2), indicating that these units have not been recently amplified as monomers. Instead, the diverged monomers have been more recently amplified (fixed) as components of larger, multimeric higher-order repeat units that are representative of and define individual β satellite arrays.

Our results indicate that the β satellite DNA family is composed of at least two, and probably more, distinct types of higher-order repeat unit. These units (pB3 and pB4) define nonoverlapping arrays of β satellite that have individually distinct genomic distributions (Figs. 3 and 4). The subsets
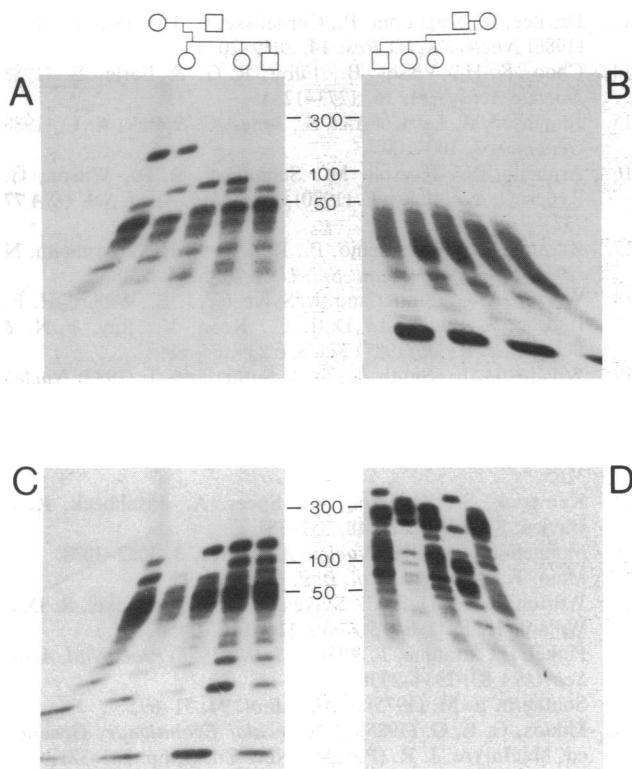
Genetics: Waye and Willard

*Proc. Natl. Acad. Sci. USA 86 (1989)* 6253



FIG. 3. Long-range organization of human β satellite subsets. Human genomic DNAs were digested to completion with *Eco*RI or *Xba* I, separated by PFGE, and probed under high stringency conditions with pB3 or pB4. (*A* and *B*) *Xba* I (*A*) and *Eco*RI (*B*) probed with pB3. (*C* and *D*) The same blots as in *A* and *B* (*C*, *Xba* I; *D*, *Eco*RI) hybridized with pB4. Size markers are indicated in kb as determined by a ladder of λ concatemers. The pedigrees indicate the relationships among the individuals tested.

defined by these higher-order repeat units represent only a fraction of the total β satellite sequences detected in the human genome.

We have demonstrated that one β satellite higher-order repeat unit (represented by pB3) is represented ≈50–100 times per haploid genome and can be localized specifically to chromosome 9 (Fig. 4*A*). This chromosome is marked by a large block of constitutive heterochromatin (29) and, by *in situ* hybridization, has been shown to contain various satellite DNAs (2, 7). One cloned sequence localized to 9qh is composed of tandem repeat units based on a satellite III-related sequence motif (7). Thus, chromosome 9 contains at least two distinct types of satellite (β satellite and satellite III), both of which are chromosome-specific. In addition, chromosome 9, like all human chromosomes, contains α satellite DNA (30), although such centromeric sequences have not as yet been cloned or characterized.

A second β satellite higher-order repeat unit (defined by pB4) is represented ≈250–500 times per haploid genome and is more widely distributed in the genome than pB3 (Figs. 3 and 4). Somatic cell hybrid mapping and *in situ* hybridization have identified pB4 β satellite domains on each of the acrocentric chromosomes. Several reports have previously identified a variety of repetitive sequences that are localized to the acrocentric short arms (9–12, 18) or centromeres (13–15). In one such study (published while this work was in progress), a cloned ≈67-bp *Sau*3A fragment was shown to hybridize *in situ* to the acrocentric short arms (11). The probe used in that study shares ≈90% sequence identity with the β satellite average sequence presented here (Fig. 2) and, as such, may be defined as β satellite.
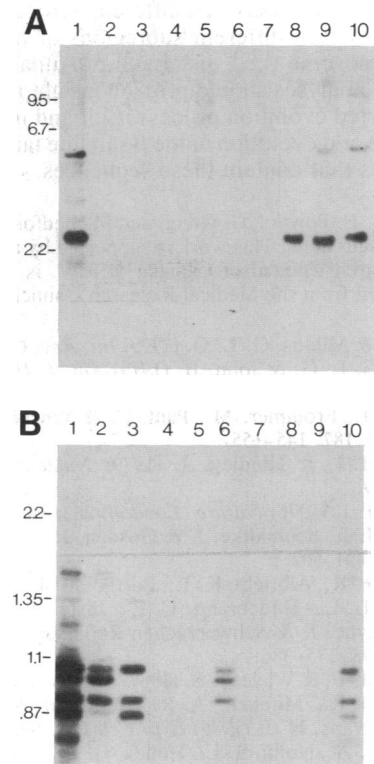
FIG. 4. Mapping of β satellite sequences using human/rodent somatic cell hybrid DNAs. (*A*) Human genomic DNA (lane 1), mouse genomic DNA (lane 2), and DNAs from eight different hybrid cell lines (lanes 3–10) were digested with *Eco*RI and hybridized under high stringency to pB3. Hybrids in lanes 8–10 contain human chromosome 9. The hybrid in lane 3 contains human chromosome 21. (*B*) Human genomic DNA (lane 1) and DNAs from different hybrid lines (lanes 2–10) were digested with *Hae* III and hybridized under high stringency to pB4. The hybrid in lane 2 contains human chromosome 21. The hybrid in lane 3 contains only the 21p12→21pter portion of chromosome 21, as part of a X;21 translocation chromosome (18, 25). Hybrids in lanes 6 and 10 contain human chromosomes 15 and 22, respectively. Hybrids in other lanes do not contain any human acrocentric chromosomes. Size standards are indicated to the left, in kb.

The short arms of the acrocentric chromosomes are rich in satellite DNAs and, with the exception of the rRNA gene clusters (31, 32), are thought to be devoid of functional genes. It has been proposed that the acrocentric short arms have evolved in a concerted manner, the presumed consequence of interchromosomal exchanges between the rRNA gene clusters or other sequences shared by these chromosomes (16, 17). If interchromosomal recombination between rRNA genes is responsible for concerted evolution, one might expect that the sequences distal to the rRNA clusters would be very similar on all of the acrocentric chromosomes (e.g., ref. 18) and that the proximal sequences might evolve in a more chromosome-specific manner.

It is therefore particularly noteworthy that β satellite sequences appear to map proximal and distal to the rRNA genes by *in situ* hybridization and, at least in the case of chromosome 21, by filter hybridization analysis (Fig. 4). Further, preliminary experiments indicate that, even within a single acrocentric chromosome, there exist multiple sub-families of β satellite that are distinguishable by nucleotide sequence divergence and by physical mapping (e.g., PFGE) (G. Greig, J.S.W., and H.F.W., unpublished data).

With respect to the possible mechanisms involved in concerted evolution of the acrocentrics, therefore, it will be of interest to determine the physical relationship between the rRNA clusters and various β satellite sequences and compare

β satellite sequences between different acrocentric chromosomes and between different subregions of the same acrocentric chromosome (i.e., distal and proximal to the rRNA genes). Such analyses should provide insight into the nature of the concerted evolution process itself and into the molecular structure and evolution of the β satellite family and of the chromosomes that contain these sequences.

1. John, B. & Miklos, G. L. G. (1979) *Int. Rev. Cytol.* **58**, 1–114.
2. Miklos, G. L. G. & John, B. (1979) *Am. J. Hum. Genet.* **31**, 264–280.
3. Prosser, J., Frommer, M., Paul, C. & Vincent, P. C. (1986) *Mol. Biol.* **187**, 145–155.
4. Cooke, H. J. & Hindley, J. (1979) *Nucleic Acids Res.* **6**, 3177–3179.
5. Cooke, H. J. (1976) *Nature (London)* **262**, 182–186.
6. Cooke, H. J., Schmidtke, J. & Gosden, J. R. (1982) *Chromosoma* **87**, 491–502.
7. Moyzis, R. K., Albright, K. L., Bartholdi, M. F., Cram, L. S., Deaven, L. L., Hildebrand, C. E., Joste, N. E., Longmire, J. L., Meyne, J. & Schwarzachen-Robinson, T. (1987) *Chromosoma* **95**, 375–386.
8. Willard, H. F. & Waye, J. S. (1987) *Trends Genet.* **3**, 192–198.
9. Gosden, J. R., Mitchell, A. R., Buckland, R. A., Clayton, R. P. & Evans, H. J. (1975) *Exp. Cell Res.* **92**, 148–158.
10. Devine, E. A., Nolin, S. L., Houck, G. E., Jr., Jenkins, E. C. & Brown, W. T. (1985) *Am. J. Hum. Genet.* **37**, 114–123.
11. Agresti, A., Rainaldi, G., Lobbiani, A., Magnani, I., DiLernia, R., Meneveri, R., Siccardi, A. G. & Ginelli, E. (1987) *Hum. Genet.* **75**, 326–332.
12. Kurnit, D. M., Roy, S., Stewart, G. D., Schwedock, J., Neve, R. L., Bruns, G. A. P., Van Keuren, M. L. & Patterson, D. (1986) *Cytogenet. Cell Genet.* **43**, 109–116.
13. Devilee, P., Slagboom, P., Cornelisse, C. J. & Pearson, P. L. (1986) *Nucleic Acids Res.* **14**, 2059–2073.
14. Choo, K. H., Vissel, B., Filby, R. G. & Earle, E. (1988) *Nucleic Acids Res.* **16**, 1273–1284.
15. Jorgensen, A. L., Kolvraa, S., Jones, C. & Bak, A. L. (1988) *Genomics* **3**, 100–109.
16. Arnheim, N., Krystal, M., Schmickel, R. D., Wilson, G., Ryder, O. & Zimmer, E. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 7323–7327.
17. Krystal, M., D'Eustachio, P., Ruddle, F. H. & Arnheim, N. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5744–5748.
18. Worton, R. G., Sutherland, J., Sylvester, J. E., Willard, H. F., Bodrug, S., Dubé, I., Duff, C., Kean, V., Ray, P. N. & Schmickel, R. D. (1987) *Science* **239**, 64–68.
19. Willard, H. F., Smith, K. D. & Sutherland, J. (1983) *Nucleic Acids Res.* **11**, 2017–2033.
20. Willard, H. F. & Riordan, J. R. (1985) *Science* **230**, 940–942.
21. Waye, J. S. & Willard, H. F. (1986) *Mol. Cell. Biol.* **6**, 3156–3165.
22. Kenwrick, S., Patterson, M., Speer, A., Fischbeck, K. & Davies, K. (1987) *Cell* **48**, 351–357.
23. Manuelidis, L. (1976) *Nucleic Acids Res.* **3**, 3063–3078.
24. Maio, J. J. (1971) *J. Mol. Biol.* **56**, 579–595.
25. Worton, R. G., Duff, C., Sylvester, J. E., Schmickel, R. D. & Willard, H. F. (1984) *Science* **224**, 1447–1450.
26. Pinkel, D., Straume, T. & Gray, J. W. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 2934–2938.
27. Southern, E. M. (1975) *J. Mol. Biol.* **94**, 51–69.
28. Miklos, G. L. G. (1985) in *Molecular Evolutionary Genetics*, ed. MacIntyre, J. R. (Plenum, New York), pp. 241–321.
29. Craig-Holmes, A. P. & Shaw, M. W. (1971) *Science* **174**, 702–704.
30. Manuelidis, L. (1978) *Chromosoma* **66**, 23–32.
31. Henderson, A. S., Warburton, D. & Atwood, K. C. (1972) *Proc. Natl. Acad. Sci. USA* **69**, 3394–3398.
32. Evans, H. J., Buckland, R. A. & Pardue, M. L. (1974) *Chromosoma* **48**, 405–426.