

## Sequence of an HLA-DR $\alpha$ -chain cDNA clone and intron-exon organization of the corresponding gene

Janet S. Lee, John Trowsdale, Paul J. Travers, Janet Carey, Frank Grosveld\*, John Jenkins & Walter F. Bodmer

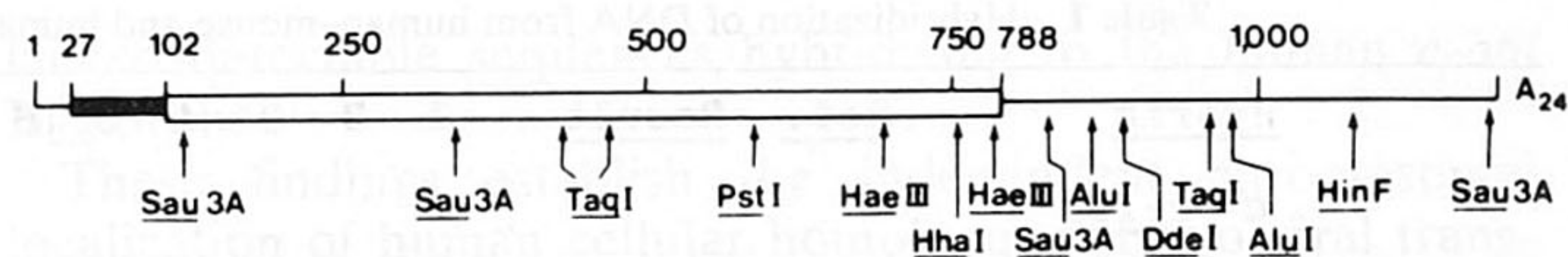
Imperial Cancer Research Fund, 44 Lincoln's Inn Fields, London WC2A 3PX, UK

\* National Institute for Medical Research, Mill Hill, London NW7 1AA, UK

The HLA-DR and related antigens are the major cell-surface glycoproteins identified so far as involved in regulation of the immune response in humans<sup>1</sup>. Here we present and discuss the sequence of cloned cDNA copies of the mRNA for the HLA-DR  $\alpha$ -chain<sup>2</sup>. The sequence is 1,195 nucleotides long, with one open reading frame encoding a 254 amino acid primary translation product. After cleavage of the signal sequence this yields a mature polypeptide of 229 residues. Four introns have been located within the corresponding genomic sequence revealing a correlation between the protein domain structure and the genomic exon organization. Analysis of the HLA-DR  $\alpha$ -chain amino acid sequence shows a clear homology with HLA-ABC,  $\beta_2$ -microglobulin and immunoglobulin domains.

The structure of the cDNA insert of clone pDRH2 is shown in Fig. 1, and the corresponding nucleotide sequence is shown in Fig. 2. Starting from position 27, with an ATG initiation codon, there is a continuous stretch of 762 nucleotides in an open reading frame, encoding 254 amino acids. The known N-terminal sequence begins at position 102 with the isoleucine codon ATC<sup>2,3</sup>, suggesting that there is a signal sequence<sup>4,5</sup> of 25 predominantly non-polar amino acids, and giving a mature product of 229 amino acids. The remaining 5'-untranslated region in pDRH2 is incomplete.

Beginning with amino acid residue 192 there is a strongly hydrophobic sequence of 23 amino acids, which has all the expected features of a transmembrane region. The resulting 191 amino acid extracellular sequence corresponds reasonably well with the size estimated from papain water-soluble fragments of the HLA-DR antigens<sup>6</sup>. This leaves a cytoplasmic portion of 15 amino acids, from position 214 to the C-terminus at 229, which contains several hydrophilic residues, as is characteristic of the cytoplasmic segments of other transmembrane glycoproteins, such as HLA-ABC, H-2K and IgM<sup>7-9</sup>. The sequence shown in Fig. 2 has two potential glycosylation sites conforming to the canonical asparagine-X-(threonine or serine) sequence, at amino acid residues 78 and 118 as expected from previous biochemical analysis. There are two cysteine residues at positions 107 and 163 in the extracellular part of the HLA-DR  $\alpha$ -chain, which are presumably available for disulphide bonding and correspond quite closely in their relative positions to the cysteine pairs found, for example, in the membrane proximal domain of HLA-ABC and other related glyco-



**Fig. 1** Schematic diagram showing the structure of the cDNA coding for the HLA-DR  $\alpha$ -chain cloned in pDRH2. Nucleotide residues are numbered in the 5' to 3' direction in the same orientation as the mRNA, beginning at the first nucleotide following the 17 guanine residues that were used to insert the cDNA into the *Pst* site of pAT153<sup>2</sup>. The signal sequence (solid bar) begins at the ATG initiation codon in position 27-29, and the codon for the N-terminal amino acid, isoleucine, is at position 102-104. The coding region (open bar) extends to nucleotide 788, followed by the 3' untranslated region of 402 nucleotides before the poly(A) extension, which is 24 residues long in pDRH2. The terminal G and C extensions are not represented. The *Pst*I and *Sau*3AI cDNA-specific fragments were separated by electrophoresis in 9% polyacrylamide gels, cut out of the gels and eluted<sup>29</sup>. The isolated fragments were used in all further experiments, whether or not subjected to cleavage by other enzymes before end labelling. The only restriction enzyme sites shown are those that were labelled at the 5' or 3' ends. Unique end labelled fragments were then obtained by strand separation<sup>29</sup> or by cleavage with another restriction enzyme followed by electrophoresis in polyacrylamide gels and dilution as already described. Nucleotide sequence of an extensive series of overlapping fragments was determined by the chemical degradation procedure of Maxam and Gilbert<sup>29</sup> and confirmed on *Pst*I-*Tag*I fragments inserted into M13mp7 phage by the chain termination technique<sup>30</sup>.

proteins. This clearly suggests a two domain overall structure for the extracellular portion of the HLA-DR  $\alpha$ -chain in which the N-terminal domain has no cysteine pairs, while the membrane proximal region is perhaps an immunoglobulin-like domain with a characteristic disulphide bond.

The amino acid sequence given in Fig. 2 is in general agreement with the partial sequence given by Korman *et al.*<sup>10</sup>, with the exception of valine instead of leucine at position 217. This difference may be genuine, since a change from G to T in the first nucleotide of the valine codon GTG, to produce leucine codon TTG, would be difficult to account for by cloning or sequencing artefacts.

Following the C-terminal stop codon TAA at nucleotide positions 789-791 there is a 404 nucleotide 3'-untranslated sequence, before the poly(A) sequence. The 3'-untranslated region shown in Fig. 2 appears to be a complete non-coding region as there is a consensus polyadenylation signal AATAAA starting 28 nucleotides upstream from the poly(A) at position 1,168. There is an additional polyadenylation signal in the same orientation starting a further 102 nucleotides upstream (position 1,066), which may occasionally be used (unpublished observations), as has been found in other situations.

To obtain cloned HLA-DR  $\alpha$ -chain related genomic sequences we screened a cosmid library, made from human placental DNA (F.G., unpublished), with a labelled pDRH2 DNA probe. Out of several positive cosmids, one (10ii) contained a strongly hybridizing 3.4-kilobase (kb) *Eco*RI fragment which included most of the pDRH2 sequences, but did not contain the 5'-untranslated region or the signal sequence for the HLA-DR  $\alpha$ -chain (data not shown). Further studies, to be reported in detail elsewhere, have now shown that in addition to the gene corresponding to the pDRH2 cDNA insert on cosmid 10ii, there are other related (but significantly different) sequences on chromosome 6 (unpublished data). Figure 3 shows that the first identifiable intron whose 3' junction sequence conforms to the consensus for a splice site, occurs within the codon for amino acid 3 (before nucleotide 109). There may, of course, be introns further upstream within the signal sequence and 5'-untranslated region, but these can only be located by analysing the sequences upstream of the 5' terminus of the 3.4-kb fragment, which are not present in cosmid 10ii. The second and third introns fall within the codons for amino acids 85 and 179 respectively. These introns define two exons corresponding

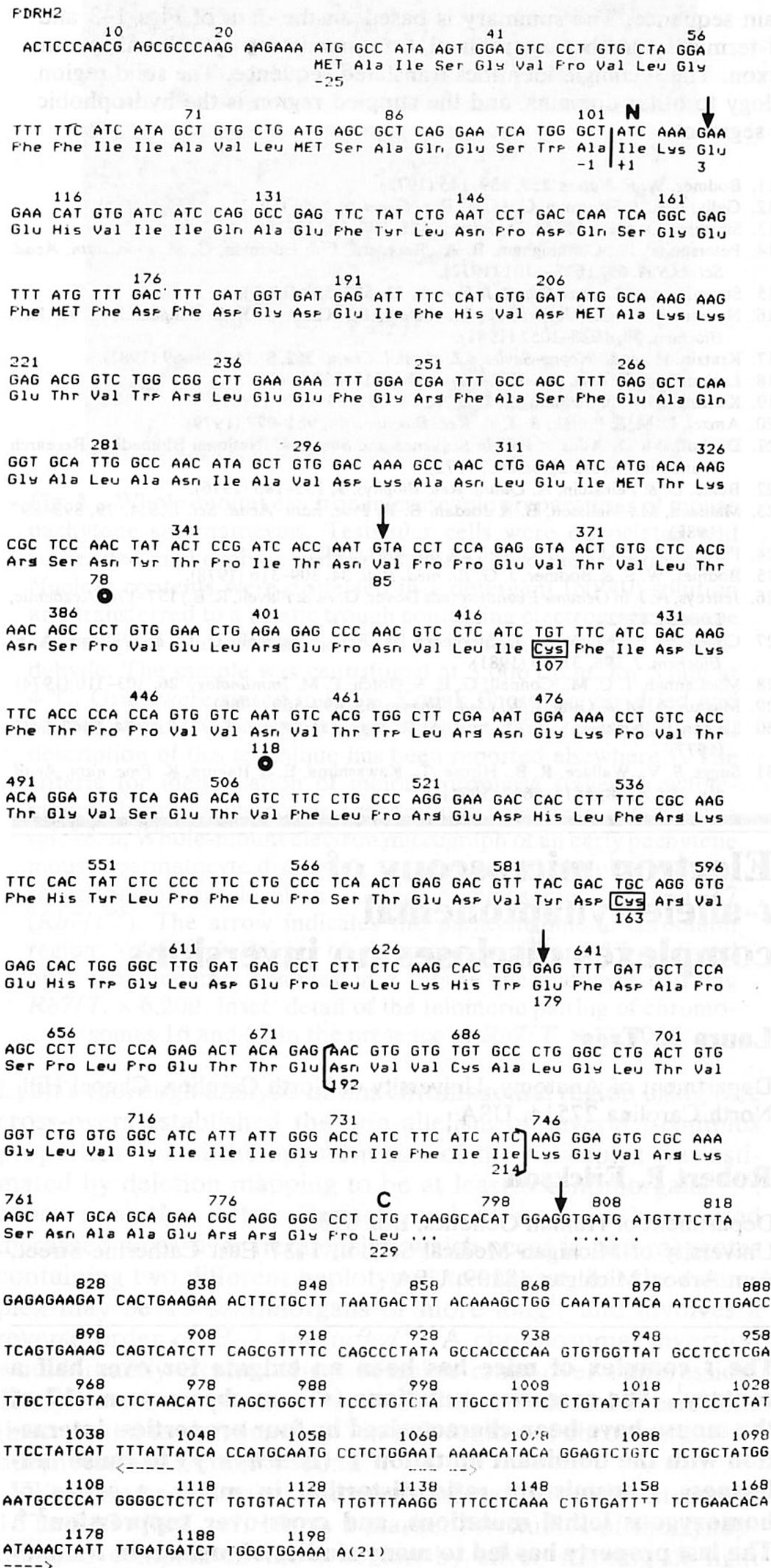
**Table 1** Per cent homologies between various HLA-associated and immunoglobulin domain sequences

	DR $\alpha$	DR $\beta$	HLA-ABC	$\beta_2$ m	C $\lambda$	C $\kappa$
DR $\alpha$	—	35	31	29	22	19
DR $\beta$		—	30	33	25	21
HLA-ABC			—	21	21	18
$\beta_2$ m				—	17	16
C $\lambda$					—	35
C $\kappa$						—

The sources of the sequence data were DR $\beta$ , ref. 17; HLA-ABC, the B7 sequence (S. Weissman, personal communication),  $\beta_2$  microglobulin, ref. 31, C $\lambda$  and C $\kappa$ , ref. 21.

to the proposed N-terminal domain and the immunoglobulin-like membrane proximal domain respectively.

The last intron to be identified occurs 10 nucleotides downstream from the first stop codon, defining an exon that includes mainly the transmembrane and cytoplasmic portion of the HLA-DR  $\alpha$ -chain. It is interesting that this last splice moves

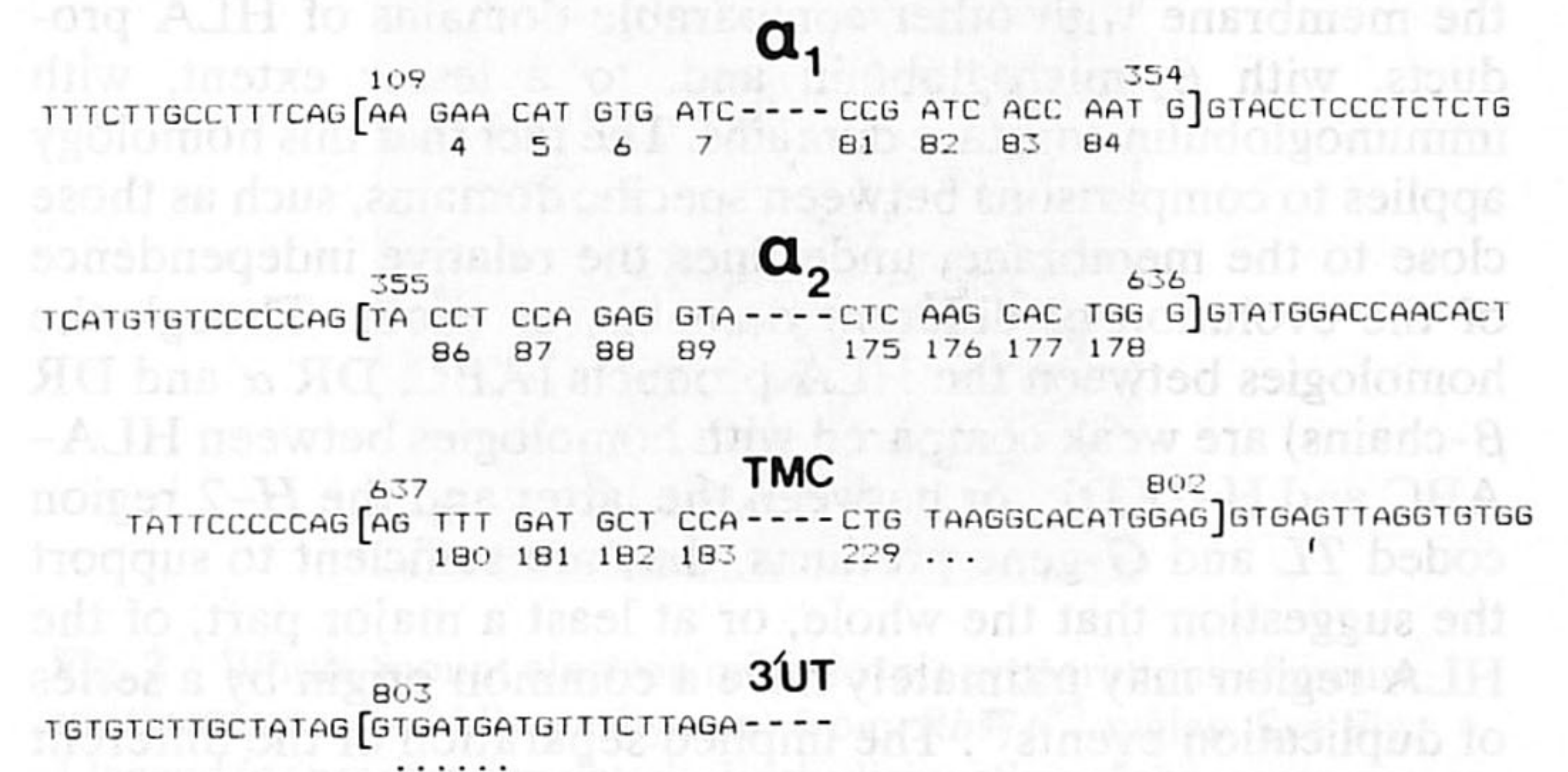


**Fig. 2** The complete nucleotide sequence of cDNA insert of pDRH2 and the predicted amino acid sequence of the HLA-DR  $\alpha$ -chain. The nucleotide sequence was determined as explained in Fig. 1 legend and follows the same numbering scheme. The derived amino acid sequence is shown below the nucleotide sequence, and numbered below from -25 to -1 for the signal sequence and from 1-229 for the mature HLA-DR  $\alpha$ -chain. The two cysteine residues that are presumed to form the interchain S-S bond are indicated by boxes; N and C indicate the amino, and carboxy termini respectively; ● indicate the glycosylation sites; arrows above the nucleotide sequence designate the known intron-exon boundaries (see text and Fig. 4); and the vertical lines after amino acids 191 and 214 delineate the transmembrane (TM) hydrophobic region. The in-phase termination codons are indicated by dots, and the poly(A) addition sites by arrows indicating direction.

two more nonsense codons at nucleotide positions 804-809 into phase. The remainder of the 3'-untranslated region is completely contained in the last exon (data not shown), but no sequence in the 3' flanking region in the genomic DNA has yet been determined.

A number of protein sequence studies have revealed homologies between various immunoglobulin domains,  $\beta_2$ -microglobulin and HLA-ABC and H-2K D domains<sup>13-16</sup>, as had been predicted<sup>11,12</sup>. More recently the membrane proximal external domain of HLA-DR  $\beta$ -chains has been added to this list<sup>17-19</sup> and from the sequence data given in Fig. 2 the membrane proximal external domain,  $\alpha_2$ , of HLA-DR  $\alpha$ -chains can also be included in this set of homologies.

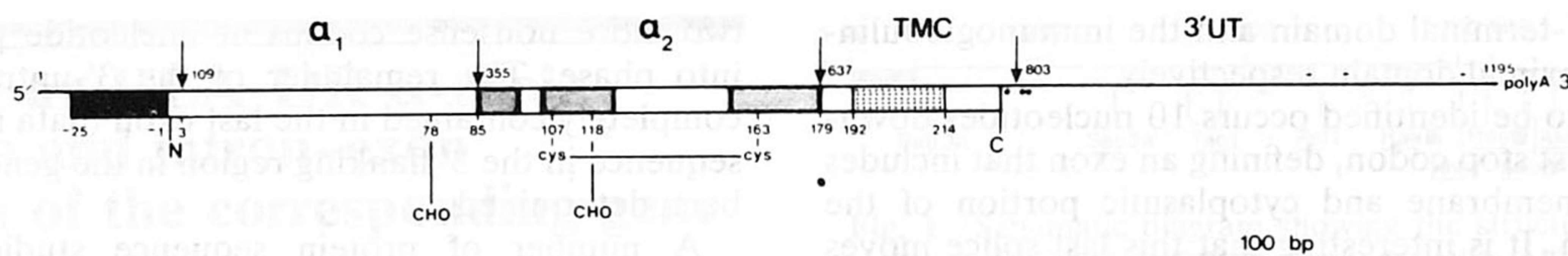
The  $\alpha_2$  domain of the HLA-DR  $\alpha$ -chain, defined by the splice points within the codons for amino acids 85 and 179 (see Fig. 2), was compared with the corresponding membrane proximal external domains of HLA-DR  $\beta$  and HLA-ABC chains as well as with  $\beta_2$  microglobulin and the constant regions of human  $\kappa$  and  $\lambda$  immunoglobulin light chains. In addition to the pair of cysteines, proline residues at positions 87 and 114, phenylalanine at 112, tyrosine at 161, valine at 165, and his-



**Fig. 3** Nucleotide sequences at the identified intron-exon junctions. The 3.4-kb *EcoRI* fragment from cosmid 10ii was first subcloned in pAT153 using standard techniques, and then *Sau3AI* and *EcoRI-PstI* fragments of the insert were further subcloned into the single-stranded DNA phage M13 to locate the intron-exon junctions within the genomic DNA by sequence determination. Templates with pDRH2 cDNA related sequences were identified by hybridization of dot-blot with the labelled cDNA insert. Several of these phages were then sequenced by the chain termination technique<sup>30</sup>.  $\alpha_1$  and  $\alpha_2$  refer respectively to the exons encoding the N-terminal and membrane proximal external domains. TMC refers to the transmembrane and cytoplasmic exon. Numbers above the sequence indicate nucleotide position and below indicate amino acid positions corresponding to those in Fig. 2. The three in-phase codons are underlined with dots.

tidine at 167 are common to all the domains compared in Table 1, and in fact to nearly all immunoglobulin constant region domains<sup>15,16,20-22</sup>. These conserved residues may have some important structural role in maintaining the characteristic immunoglobulin fold<sup>21</sup>.

A number of homologies are found among the HLA associated domains themselves, which are more closely related to each other than they are to immunoglobulins. For example, all of the HLA related sequences have tryptophan at position 178. Moreover, in the HLA-DR  $\alpha_2$  domain,  $\beta_2$  microglobulin and the HLA-ABC  $\alpha_3$  domain, this codon occurs one base before the intron between this exon and the exon encoding the transmembrane region. The valine at position 85 in the HLA-DR  $\alpha$ -chain contains the splice site at the start of the  $\alpha_2$  exon and is also present in the HLA-DR  $\beta$ -chain sequence. These splice sites may lie in the same positions in the  $\beta$ -chain, a possibility further emphasised by the fact that the aspartate in the homologous position in the HLA-ABC  $\alpha_3$  domain<sup>23</sup> and the arginine in  $\beta_2$  microglobulin<sup>24</sup> also contain the upstream splice sites for the respective domains. Other conserved sites of interest for the HLA associated set of domains are the proline at position 155 and the phenylalanine at position 145.



**Fig. 4** Schematic summary of the organization of the HLA-DR  $\alpha$ -chain sequence. The summary is based on the data of Figs 1–3 and follows their numbering and conventions.  $\alpha_1$ ,  $\alpha_2$ , and TMC refer to the N-terminal, membrane-proximal, transmembrane-cytoplasmic exons or domains respectively, and 3'-UT refers to the 3' untranslated region exon. The rectangle identifies translated sequence. The solid region is the signal sequence, the shaded regions are those with greatest homology to other domains, and the stippled region is the hydrophobic transmembrane segment.

A schematic summary of the organization of the HLA-DR  $\alpha$ -chain sequence based on the data of Figs 1, 2 and 3 is given in Fig. 4. The most obvious features revealed by the sequence analysis are the organization of the HLA-DR  $\alpha$ -chain into exons that correspond surprisingly well to those defined for HLA-ABC chains,  $\beta_2$  microglobulin and, of course, immunoglobulins, and the demonstration of a remarkable correspondence between exon organization and the presumptive protein domains. Another striking feature of the HLA-DR  $\alpha$ -chain sequence data is the homology of the external domain next to the membrane with other comparable domains of HLA products, with  $\beta_2$ -microglobulin and, to a lesser extent, with immunoglobulin constant domains. The fact that this homology applies to comparisons between specific domains, such as those close to the membrane, underlines the relative independence of the evolution of different domains, or exons. Though the homologies between the HLA products (ABC, DR  $\alpha$  and DR  $\beta$ -chains) are weak compared with homologies between HLA-ABC and H-2KDL, or between the latter and the H-2 region coded *TL* and *Q* gene products, they are sufficient to support the suggestion that the whole, or at least a major part, of the HLA region may ultimately have a common origin by a series of duplication events<sup>25</sup>. The implied separation of the different classes of HLA products by duplication must, however, have occurred up to 500 million years ago, and certainly well before the evolutionary divergence of the mammals<sup>26</sup>.

The suggestion that there might be homology between HLA products and immunoglobulins was made on the basis of a functional analogy, namely the possibility that HLA products might have recognition functions analogous to those of the variable regions of immunoglobulins<sup>11,12</sup>. Subsequently, however, it became clear that while HLA products do play a part in cellular interactions and recognition, at least in the immune system, they do not have an analogue of the immunoglobulin variable region. Thus, it is not surprising that homologies between HLA and immunoglobulins revealed by sequence analysis involve particularly the constant region domains, as might be expected if HLA products are involved in interactions similar to those that occur between immunoglobulins and complement<sup>25</sup>, Fc receptors, monocyte receptors and secretory components<sup>27</sup> as well as between complement and cellular component receptors. Another relevant analogy might be that between the C<sub>H</sub>3 domain of IgG and the appropriate receptor on specific antibody dependent killer cells<sup>28</sup>.

We thank P. Friedland, R. Staden, A. Williams, J. Rogers, A. Mellor, T. Lund, R. Flavell, D. Larhammar, A. Korman and C. Furse for their help.

Received 16 July; accepted 19 August 1982.

- Katz, D. H. & B. Benacerraf (eds) *The Role of the Histocompatibility Gene Complex in Immune Response* (Academic, London, 1976).
- Lee, J. S., Trowsdale, J. & Bodmer, W. F. *Proc. natn. Acad. Sci. U.S.A.* **79**, 545–549 (1982).
- Springer, T. A., Kaufman, J. F., Terhorst, C. & Strominger, J. L. *Nature* **268**, 213–218 (1977).
- Korman, A. J., Ploegh, H. L., Kaufman, J. F., Owen, M. J. & Strominger, J. L. *J. exp. Med.* **152**, 655–82s (1980).
- Owen, M. J., Kissonerghis, A. M., Lodish, H. F. & Crumpton, M. J. *J. biol. Chem.* **256**, 8987–8993 (1981).
- Kaufman, J. F. and Strominger, J. L. *Proc. natn. Acad. Sci. U.S.A.* **76**, 6304–6308 (1979).
- Robb, R. J., Terhorst, C. & Strominger, J. L. *J. biol. Chem.* **253**, 5319–5324 (1978).
- Coligan, J. E., Kindt, T. J., Uehara, H., Martinko, J. & Nathanson, S. G. *Nature* **291**, 35–39 (1981).
- Rogers, J. *et al. Cell* **20**, 303–312 (1980).
- Korman, A. J., Knudsen, P. J., Kaufman, J. F. & Strominger, J. L. *Proc. natn. Acad. Sci. U.S.A.* **79**, 1844–1848 (1982).

- Bodmer, W. F. *Nature* **237**, 139–145 (1972).
- Gally, J. A. & Ederman, G. M. *A. Rev. Genet.* **6**, 1–46 (1972).
- Smithies, O. & Poulik, M. D. *Science* **175**, 187–189 (1972).
- Peterson, P. A., Cunningham, B. A., Berggard, I. & Edelman, G. M. *Proc. natn. Acad. Sci. U.S.A.* **69**, 1697–1701 (1972).
- Strominger, J. L. *et al. Scand. J. Immun.* **11**, 573–592 (1980).
- Nathanson, S. G., Uehara, H., Ewenstein, B., Kindt, T. J. & Coligan, J. E. *A. Rev. Biochem.* **50**, 1025–1052 (1981).
- Kratz, H., *et al. Hoppe-Seyler's Z. physiol. Chem.* **362**, S. 1665–1669 (1981).
- Larhammar, D. *et al. Scand. J. Immun.* **14**, 617–622 (1981).
- Kaufman, J. F. & Strominger, J. L. *Nature* **297**, 694–697 (1982).
- Amzel, L. M. & Poljak, R. J. *A. Rev. Biochem.* **48**, 961–977 (1979).
- Dayhoff, M. O. *Atlas of Protein Sequence and Structure*. (National Biomedical Research Foundation, Washington DC, 1972).
- Beale, D. & Feinstein, A. *Quant. Rev. Biophys.* **9**, 135–180 (1976).
- Malissen, M., Malissen, B. & Jordan, B. R. *Proc. natn. Acad. Sci. U.S.A.* **79**, 893–897 (1982).
- Parnes, J. & Seidman, J. G. *Cell* **29**, 661–669 (1982).
- Bodmer, W. F. & Bodmer, J. G. *Br. med. Bull.* **34**, 309–316 (1978).
- Jeffreys, A. J. in *Genome Evolution* (eds Dover, G. A. & Flavell, R. B.) 157–176 (Academic, London, 1982).
- Cohen, F. E., Novotney, J., Sternberg, M. J. E., Campbell, D. G. & Williams, A. F. *Biochem. J.* **195**, 31–40 (1981).
- MacLennan, I. C. M., Connell, G. E. & Gotch, F. M. *Immunology* **26**, 303–310 (1974).
- Maxam, A. M. & Gilbert, W. *Meth. Enzym.* **65**, 499–560 (1980).
- Sanger, F., Nicklen, S. & Coulson, A. R. *Proc. natn. Acad. Sci. U.S.A.* **74**, 5463–5467 (1977).
- Suggs, S. V., Wallace, R. B., Hirose, T., Kawashima, E. & Itakura, K. *Proc. natn. Acad. Sci. U.S.A.* **78**, 6613–6617 (1981).

## Electron microscopy of *t*-allele synaptonemal complexes discloses no inversions

Laura L. Tres

Department of Anatomy, University of North Carolina, Chapel Hill, North Carolina 27514, USA

Robert P. Erickson

Department of Human Genetics, Box 015, University of Michigan Medical School, 1137 East Catherine Street, Ann Arbor, Michigan 48109, USA

The *t* complex of mice has been an enigma for over half a century<sup>1</sup>. The recessive mutations (*t*<sup>n</sup>) on chromosome 17 of the mouse have been characterized by four properties: interaction with the dominant mutation *T* (*Brachyury*) to cause taillessness, transmission ratio distortion in males, a series of homozygous lethal mutations, and cross-over suppression<sup>2–4</sup>. The last property has led to many studies of meiosis in *t*-allele-bearing males. Forejt found fewer chiasmata, the cytological correlate of crossing-over, in the region of *t* haplotypes<sup>5,6</sup>. Lyon *et al.*<sup>7</sup> confirmed these results and suggested that chiasma suppression determined by *t* alleles is probably desynaptic but that electron microscopy would probably be needed to settle this point. We have now extended these observations by using a centromeric translocation and analysed synaptonemal complexes by electron microscopy to study *t*-haplotype pairing. We have found normal synaptonemal complexes without any evidence of inversions but detected early disjunction of the chromosome 17 homologues, supporting the idea of desynaptic chiasma suppression.

The cross-over suppression of *t* haplotypes masks the true genetic nature of the region. It has previously been convenient to consider the *t* alleles as point mutations mapping near *T*.