

# Two types of somatic recombination are necessary for the generation of complete immunoglobulin heavy-chain genes

Hitoshi Sakano, Richard Maki, Yoshikazu Kurosawa, William Roeder & Susumu Tonegawa

Basel Institute for Immunology, 487 Grenzacherstrasse, Postfach CH 4005, Basel 5, Switzerland

*At least two types of somatic recombination are necessary for the generation of a complete immunoglobulin  $\gamma 2b$  gene from germ-line DNA sequences. The first type of recombination consists of the assembly of three separate DNA segments, each encoding a different part of the variable region. The second type of recombination replaces the exons coding for the constant region of the  $\mu$  chain with those coding for the same region of the  $\gamma 2b$  chain. The DNA sequencing studies suggest that the two types of recombination operate by different mechanisms.*

COMPLETE, active immunoglobulin genes are created by somatic recombination that occurs during the differentiation of lymphocyte precursor cells<sup>1-3</sup>. The organization of the gene sequences before and after somatic recombination has been studied extensively for both  $\lambda$ - and  $\kappa$ -type light chains in the mouse<sup>3-10</sup>. These studies established that in the embryonic genome, the major portion of the conventionally defined variable (V) region is encoded in a DNA segment (V DNA) that is located some distance away from a DNA segment (J DNA) coding for the rest of the V region. The J DNA segment has been mapped a few kilobases upstream of (that is, 5' side, with respect to the direction of transcription) the single-copy DNA segment (C DNA) coding for the constant (C) region. In the myeloma cells in which the light-chain gene in question is active, the V DNA segment and the J DNA segment are contiguous as a consequence of a recombination that apparently accompanies deletion of the DNA sequence occurring between the V and J DNA segments in the germ-line genome<sup>7</sup>. The entire region of the chromosome containing the V DNA segment, J DNA segment, J-C intron and C DNA segment is transcribed into a single RNA molecule from which a mature messenger RNA is generated by RNA splicing<sup>11,12</sup>. The V-J joining step is considered to be a key event in the activation of the immunoglobulin gene. It may also contribute to the increase of the organism's capacity to synthesize a large number of different antibodies<sup>7,13-15</sup>.

As the heavy-chain polypeptide is also composed of a V and C region, one might assume that what has been established at the DNA level for the light chains would also be true for heavy chains. However, the expression of heavy chains has some features which make it considerably more complex than that of light chains. For instance, unlike the light chains, there are in the mouse at least eight different C regions, each of which seems to share the same set of V regions, that is, a given V gene can be expressed with more than one heavy-chain class or subclass. In addition, it seems that, as the lymphocytes differentiate, the C region of the heavy chain synthesized switches from one class to another without alteration of the V region<sup>16-18</sup>.

Recently, we reported on the structure of the complete  $\gamma 2b$  gene isolated from a myeloma, MOPC141, an IgG2b secretor<sup>19</sup>. These studies showed that, unlike the light-chain genes, at least two recombinations are necessary to generate the heavy-chain gene. One of them occurs at or near the 3' end of the V DNA segment and in the vicinity of the J DNA segment. In contrast to the light-chain genes, where the J DNA segments are located in the 5'-flanking regions of their respective C genes, in the heavy-chain gene system the J DNA used for the MOPC141  $\gamma 2b$  gene does not lie in the 5'-flanking region of the germ-line  $C\gamma 2b$  gene,

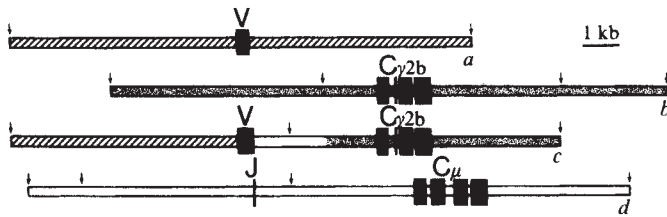
but is located in the 5'-flanking region of the germ-line  $C\mu$  gene. This recombination between V DNA and J DNA near the  $C\mu$  gene probably allowed the  $\mu$  gene to be activated in the precursor B cells to MOPC141 myeloma. The second recombination occurs between a pair of sites, one located between the J DNA segment and the  $C\mu$  gene and the other in the 5'-flanking sequence of the  $C\gamma 2b$  gene, and replaces the  $C\mu$ -coding exons with the  $C\gamma 2b$  exons. This recombination is referred to as 'switch recombination', for it seems to be a key event preceding the heavy-chain switch. A similar gene structure was also observed for a complete  $\alpha$ -chain gene<sup>20</sup> and for a  $\gamma 1$ -chain gene<sup>21</sup>.

We report here on the two types of recombination, studied by DNA sequencing. The sequences around the V-J joining and the switch recombination sites have different features, suggesting that at least two distinct enzyme systems are involved in the generation of the active immunoglobulin  $\gamma 2b$  gene. In addition, sequence analysis of the V DNA of MOPC141 and its germ-line components, the embryonic V DNA and J DNA, revealed that the third hypervariable region is encoded in a separate DNA segment(s) in the germ-line genome.

Mouse DNA inserts in the DNA clones used in this study are listed in Fig. 1. The homology among these DNA fragments has been determined by heteroduplex analysis and restriction enzyme mapping. The coding regions were mapped by R-loop analysis using mRNAs from MOPC141 ( $\gamma 2b$ ) or MOPC104E ( $\mu$ ). Isolation and electron microscopic characterization of clones M141-P21 (complete  $\gamma 2b$ ), MEP203 (embryonic J plus  $C\mu$ ) and MEP3 (embryonic  $C\gamma 2b$ ) have been described elsewhere<sup>19</sup>. The isolation of the embryonic V-gene clone, PJ14, is described below.

## Identification of J DNA for MOPC141

We have previously identified the putative V-J joining site on the complete  $\gamma 2b$  gene clone (M141-P21) and on an embryonic  $C\mu$  gene clone (MEP203) by heteroduplex analysis and R-loop mapping<sup>19</sup>. The electron microscopic studies indicated that the putative J DNA for the MOPC141 heavy-chain gene lies 3.6 kilobases 5' to the  $C\gamma 2b$  gene on M141-P21 and 5 kilobases 5' to the  $C\mu$  gene on MEP203. To analyse the structural features of heavy-chain V-J joining sites, we determined the nucleotide sequences of the appropriate regions of M141-P21 and MEP203. Although the amino acid sequence of the MOPC141 heavy chain is unknown, identification of the J DNA segment coding for this chain was possible because J-peptide sequences are highly conserved. As shown in Fig. 2, M141-P21 carries a DNA sequence in the vicinity of the putative V-J joining sites which can encode a peptide that is identical to the J peptides of



**Fig. 1** Mouse DNA inserts in the clones used in the study. Embryonic  $V_{M141}$  gene clone PJ14 (a), embryonic  $C\gamma 2b$  gene clone MEP3 (b), complete  $\gamma 2b$  gene clone M141-P21 (abbreviated to M141) (c) and embryonic  $C\mu$  gene clone MEP203 (d) are shown. Arrows indicate *EcoRI* cleavage sites. Filled boxes are exons identified by R-loop mapping. The regions of homology among the four DNA clones as determined by heteroduplex analysis are indicated by bars of different types of shading.

MOPC21 and MOPC173 heavy chains (see Table 1). We also analysed the MEP203 insert to determine whether the same J DNA segment is present in the expected position. Our previous study on clones MEP203 and M141-P21 suggested that the  $J_{M141}$  DNA segment and the sequence that follows this segment are derived from a region 5' to the  $C\mu$  gene. In fact, the sequence of MEP203 showed that the J DNA segment and its 3' non-coding sequence found on M141-P21 are also present in the region about 5 kilobases 5' to the  $C\mu$  gene on MEP203 (Fig. 2). Comparison of these two  $J_{M141}$  DNA sequences, one on M141-P21 and the other on MEP203, indicates that the germ-line  $J_{M141}$  DNA can encode the M141 J peptide starting with the second nucleotide of the Leu codon at position 112.

In  $\gamma 2b$  heavy-chain peptides, the boundary of V and C regions has been thought to be around Ala 126 (ref. 22). The  $J_{M141}$  DNA segment can fully encode Ser 125, but only partially Ala 126. We found a triplet, GGT, instead of the Ala codon GCN. Because, in almost all cases studied thus far, an intron starts with the doublet GT (ref. 23), we tentatively conclude that the coding of the  $J_{M141}$  DNA ends with the first letter G in the triplet GGT.

#### Four heavy-chain J DNA segments around 8 kilobases 5' to the germ-line $C\mu$ gene

Previous studies showed that most, if not all,  $\kappa$ -light-chain J DNAs are tightly linked<sup>7,13</sup>. To examine whether the same applies to the heavy-chain J DNA segments, we isolated DNA clones containing the rearranged V gene from myelomas MOPC603, MOPC315, HOPC8 and MOPC173 and analysed the heteroduplexes formed between each of these DNA clones and the embryonic  $\mu$  DNA clone MEP203. It is expected that the divergency point observed on a Y-shaped heteroduplex molecule corresponds to the J DNA segment used for expression in the respective myeloma cells. In this way, we mapped four different J DNA segments in the region 5' to the  $C\mu$  gene, at 1.1 (MOPC173-II J), 1.6 (HOPC8-II J), 1.9 (MOPC315 J) and 2.2 (MOPC603 J) kilobases, 5' to the *EcoRI* site located 3.5 kilobases from  $C\mu$  DNA on MEP203. As MEP203 carries a 3-kilobase deletion at 1.5 kilobases 5' to the  $C\mu$  gene (see below), the *EcoRI* site is 6.5 kilobases from the  $C\mu$  gene on the mouse chromosome. Using  $\mu$ -chain mRNA from MOPC104E in R-loop analysis, we also identified the  $J_{M104E}$  DNA segment at 2.3 kilobases 5' to the *EcoRI* site (data not shown). As the MOPC104E  $\mu$  chain and the MOPC603  $\alpha$  chain seem to have the same J peptide (Table 1), this J DNA segment is probably the same as that detected by heteroduplex analysis using the MOPC603-derived V-gene clone.

Figures 3Aa and 4 show the restriction map and the DNA sequence of the 1.5-kilobase J-rich region, respectively. So far, 22 heavy chains have been sequenced in the J-peptide regions.

They can be classified into four different groups according to the amino acid sequences of the J peptides (Table 1). Comparison of these J-peptide sequences with the DNA sequence allowed identification of the corresponding four J DNA segments. These four segments can account for all heavy-chain J peptides sequenced so far (see Table 1) and, therefore, seem to encode the J peptides of all different classes or subclasses. We designate these J DNAs as  $J_{H1}$ ,  $J_{H2}$ ,  $J_{H3}$  and  $J_{H4}$ , from 5' towards 3'. No additional J-like DNA sequence has been observed in the 1.5-kilobase region sequenced thus far. There is one unexpressed J DNA,  $J_{\kappa 3}$  (ref. 7), in a  $\kappa$ -light-chain gene, and there may also be some in the heavy-chain gene.

#### Isolation of the germ-line V-gene clone for the MOPC141 heavy chain

To determine the recombination site on the germ-line V gene, we attempted to identify and clone the *EcoRI* DNA fragment carrying the germ-line V gene for the MOPC141 heavy chain. First, we analysed embryonic DNA by Southern hybridization using the V DNA probe purified from M141-P21. This probe is a 257-base pair *AvaII-HgaI* fragment encoding residues 9–94 of the M141 heavy chain (Fig. 3Ab). In the *EcoRI* digest of mouse embryo DNA, about 10 hybridizable bands were detected (data not shown). Detection of 5–10 bands by a single V-gene probe has previously been reported for both  $\kappa$  light chains<sup>5,9</sup> and heavy chains<sup>24–26</sup>. This is due to cross-hybridization among a subset of germ-line V genes that have presumably diverged more recently. As the length of the 5'-noncoding region of the  $V_H$  gene on M141-P21 is 6.5 kilobases, the size of the *EcoRI* fragment carrying the germ-line  $V_H$  DNA for the M141 heavy chain should be larger than 6.5 kilobases. We separated the *EcoRI* digest of embryonic DNA on a preparative agarose gel, eluted the DNA from the 6.5–23-kilobase region and used it for DNA cloning. Using the V-region probe, we isolated some 30 clones from 2 separate experiments and classified them into 5 different types according to the size of the *EcoRI* insert. Heteroduplex molecules formed between the clones of each type and M141-P21, a complete  $\gamma 2b$  gene clone, were examined by electron microscopy.

Although the *EcoRI* inserts of these clones were all hybridizable with the V-gene probe, only the clones containing a 14-kilobase insert formed Y-shaped heteroduplex molecules with M141-P21, in which the homology extended from the 5' *EcoRI* end to a site 6.9 kilobases away from it. R-loop analysis using the MOPC141  $\gamma 2b$  mRNA indicated that the V-coding DNA segments map at or around the heteroduplex fork point. We further analysed one of these 14-kilobase clones, PJ14, by restriction mapping and DNA sequencing. Evidence that the 14-kilobase clone contains the germ-line V gene for the MOPC141  $\gamma 2b$  chain is obtained from restriction enzyme maps. The cleavage sites for *BamHI*, *KpnI*, *XbaI*, *SacI* and *HindIII* in the 5'-flanking region of the V gene on M141-P21 were all found at the corresponding positions on the embryonic V-gene clone, PJ14 (data not shown).

#### The third hypervariable region is encoded separately

The DNA sequences around the V-coding regions of PJ14 and M141-P21 are shown in Fig. 2. The sequences include a 46-base pair exon presumably coding for the hydrophobic signal peptide, an 81-base pair intron, the entire V-gene exon and the 3'-flanking regions. As in all  $\kappa$ - and  $\lambda$ -chain genes studied so far (refs 4, 6, 8 and G. Heinrich, unpublished results), the V coding begins with the residue at position -4 within the signal peptide. Also, as in the light-chain genes, the coding by the embryonic clone PJ14 ends prematurely with Ser 97, whereas that of the myeloma clone, as mentioned above, continues up to Ser 125. In the 0.6-kilobase stretch of DNA starting with the *HindIII* site in the 5'-noncoding region and ending with the Ser 97 codon, the

**Table 1** Amino acid sequences of mouse immunoglobulin heavy chains around the carboxyl ends of the V regions

Myeloma	End of V	HV3	J	J DNA
M104E	... YYCAR	DY	d W Y F D V W G A G T T V T V S S	
J558	... YYCAR	D	r Y - - - - -	
T15	... YYCAR	DYYGS	s Y - - - - -	
M603	... YYCAR	NYYGST	- - - - -	
M167	... YYCTR	DADYGDSYF	g - - - - -	J <sub>H1</sub>
T601	... YYCAR	LGYY	g - - - - -	
S107	... YYCAR	DYYGS	s Y - - - - -	
H8	... YYCAR	DYYGN	s Y - - - - -	
M511	... YYCAR	DGDYGS	s Y - - - - -	
M315	... YYCAG	DNDHL	Y F D Y W G Q G T T L T V S S	
X24	... YYCAR	LGYYG	- - - - -	J <sub>H2</sub>
Hdex4	... YYCAR	DK	n - - - - -	
Hdex5	... YYCAR	DS	n - - - - -	
A4	... YYCTT		s W F A Y W G W G T L V T V S A	
E109	... HYCTT		g - - - - -	
U61	... YYCTT		g - - - - -	J <sub>H3</sub>
A47N	... YYCST		g - - - - - P - - -	
X44	... YYCAR	LHYYGYA	- - - - -	
J539	... YYCAR	LHYYGYN	- - - - -	
M141	... YYCAS	VSIYYYGRSDKYFT	I D Y W G W G T S V T V S S	
M21	... YYCAR	HGNYPW	Y A M - - - - -	J <sub>H4</sub>
M173	... YYCAR	SP	Y Y A M - - - - -	

Amino acid sequences of 22 mouse heavy chains<sup>27,32,34-40</sup> for residues 90-113 (numbering is after Kabat *et al.*<sup>22</sup>) are shown by one-letter codes<sup>41</sup>. According to the J-region peptides, the chains are classified into four groups. The amino acid sequence of the MOPC141  $\gamma$ 2b chain was deduced from the nucleotide sequences determined in this study. The first amino acid residues in the J peptides are shown in small letters whenever the first letters of their codons are not included in the corresponding J DNA segments.

PJ14	AAGCTTGCCTGTGCTTCTTTATCCTCTCAGGAACCTCCCCAATGCAAATCAGCCCTCAGGCAGAGGATAAAAAGCTCACACAAGATGAGAAGCCCC
M141	AAGCTTGCCTGTGCTTCTTTATCCTCTCAGGAACCTCCCCAATGCAAAGCAGCCCTCAGGCAGAGGATAAAAAGCTCACACAAGATGAGAAGCCCC
	-19 L -5 MetAlaValLeuAlaLeuLeuPheCysLeuValThrPheProSerC
PJ14	ATCATCTTCTCATAGAGCCTCCATCAGAGCATGGCTGTCTGGCATTACTTCTGCTGGTAAACATCCCAAAGCTGTAAGTGTGTCCAGGGTTTCAAGAGGGACTAAAGACATGTCAG
M141	ATCATCTTCTCATAGAGCCTCCATCAGAGCATGGCTGTCTGGCATTACTTCTGCTGGTAAACATCCCAAAGCTGTAAGTGTGTCCAGGGTTTCAAGAGGGACTAAAGACATGTCAG
	-4 1 ysIleLeuSerGlnValGlnLeuLysGluSerGlyProGlyLeuValAlaProSerGlnSerLeuSerIleThrCysThrVal
PJ14	CTAATGTGTGACTAATGGTAATGTCACCTGTGCTAGGTATCCTTTCCAGGTACAGCTGAAGGAGTCAGGACCTGGCTGGTGGCCCTCAGAGCCTGTCCATCAGATGCACCGTC
M141	CTAATGTGTGACTAATGGTAATGTCACCTGTGCTAGGTATCCTTTCCAGGTACAGCTGAAGGAGTCAGGACCTGGCTGGTGGCCCTCAGAGCCTGTCCATCAGATGCACCGTC
	SerGlyPheSerLeuThrGlyThrGlyValAsnTrpValArgGlnProProGlyLysGlyLeuGluTrpLeuGlyMetIleTrpGlyAspGlySerThrAspTyrAsnSerAlaLeuLys
PJ14	TCAGGGTCTCATTAACCGCTATGGTCTAACTGGTTCGCCAGCCTCCAGGAAGGGTCTGGAGTGGCTGGGAATGATATGGGGTATGGAAGCACAGACTATAATTCAGCTCTCAAA
M141	TCAGGATCTCATTAACCGCTATGGTCTAACTGGTTCGCCAGCCTCCAGGAAGGGTCTGGAGTGGCTGGGAATGATATGGGGTATGGAAGCACAGACTATAATTCAGCTCTCAAA
	* HV1 *
	* HV2 *
	* HV3 *
PJ14	SerArgLeuSerIleSerLysAspAsnSerLysSerGlnValPheLysLysMetAsnSerLeuGlnThrAspAspThrAlaArgTyrTyrCysAlaSer TCCAGACTGAGCATCAGCAAGGACAACTCCAAGACCAAGTTTCTTAAAAATGAACAGTCTGCAAACCTGATGACACAGCCAGGTACTACTGTGCCAGAGACACAGTGAGGGAAGTCCAA
M141	SerArgLeuThrIleThrLysAspAsnSerLysSerGlnValPheLysLysMetAsnSerLeuGlnThrAspAspThrAlaArgTyrTyrCysAlaSerValSerIleTyrTyrTyrGly TCCAGACTGAGCATCAGCAAGGACAACTCCAAGACCAAGTTTCTTAAAAATGAACAGTCTGCAAACCTGATGACACAGCCAGGTACTACTGTGCCAGCGTCTCAATTTATTAATCACTACGGT
	97
	* HV3 *
PJ14	TGTGAGCCTGCACAAATACCTCTCTGCAGGGATGATCACACCAGCAGGGGGCGCTGAGGACCAAGGACTT
M141	ArgSerAspLysTyrPheThrLeuAspTyrTrpGlyGlnGlyThrSerValThrValSerSer CGTAGCCACAATACTTCACTTTGGACTACTGGGGTCAAGGAACCTCAGTCACTGTCTCTCAGGTAAGAATGGCCCTCCAGGCTTTATTTTAACTTTGTTATGGAGTTTCTC
	112 125
	AspTyrTrpGlyGlnGlyThrSerValThrValSerSer
MEP203	CACTATTGTGATTACTATGCTATGGACTACTGGGGTCAAGGAACCTCAGTCACTCTCTCAGGTAAGAATGGCCCTCCAGGCTTTATTTTATTCCTTTGTTATGGAGTTTGTG
	J

**Fig. 2** Nucleotide sequences of MOPC141 V<sub>H</sub> DNA (M141), its embryonic V<sub>H</sub> DNA (PJ14) and the embryonic J DNA (J<sub>4</sub>) coding for the J peptide of the MOPC141  $\gamma$ 2b chain (MEP203). Sequencing strategy is shown in Fig. 3A. Sequence determination was carried out according to the method of Maxam and Gilbert<sup>43</sup>. The amino acid sequence of the complete MOPC141 V region predicted by the nucleotide sequence of clone M141 is shown in italics. The amino acid sequences predicted by the germ-line V DNA segment in clone PJ14 and the J<sub>4</sub> DNA segment in clone MEP203 are also shown. The nucleotide differences observed between the two clones, M141 and PJ14, are shown by asterisks. Vertical lines indicate the possible recombination sites. The two conserved sequences in the 3'-flanking region of the embryonic V DNA are underlined. HV1, HV2 and HV3 designate the three hypervariable regions.

sequences of the two DNA clones differ in 9 base pairs. Of these differences, 6 are within the first and the second hypervariable (HV) regions. Similar, limited, single-base differences have previously been observed in the hypervariable regions of a germ-line V gene and its somatic variants for the light chains of both  $\lambda$  and  $\kappa$  types and were taken as direct evidence for somatic diversification of antibody genes<sup>6</sup>. Three other differences may be a consequence of the similar somatic mechanism or may have occurred during the propagation of the myeloma cells.

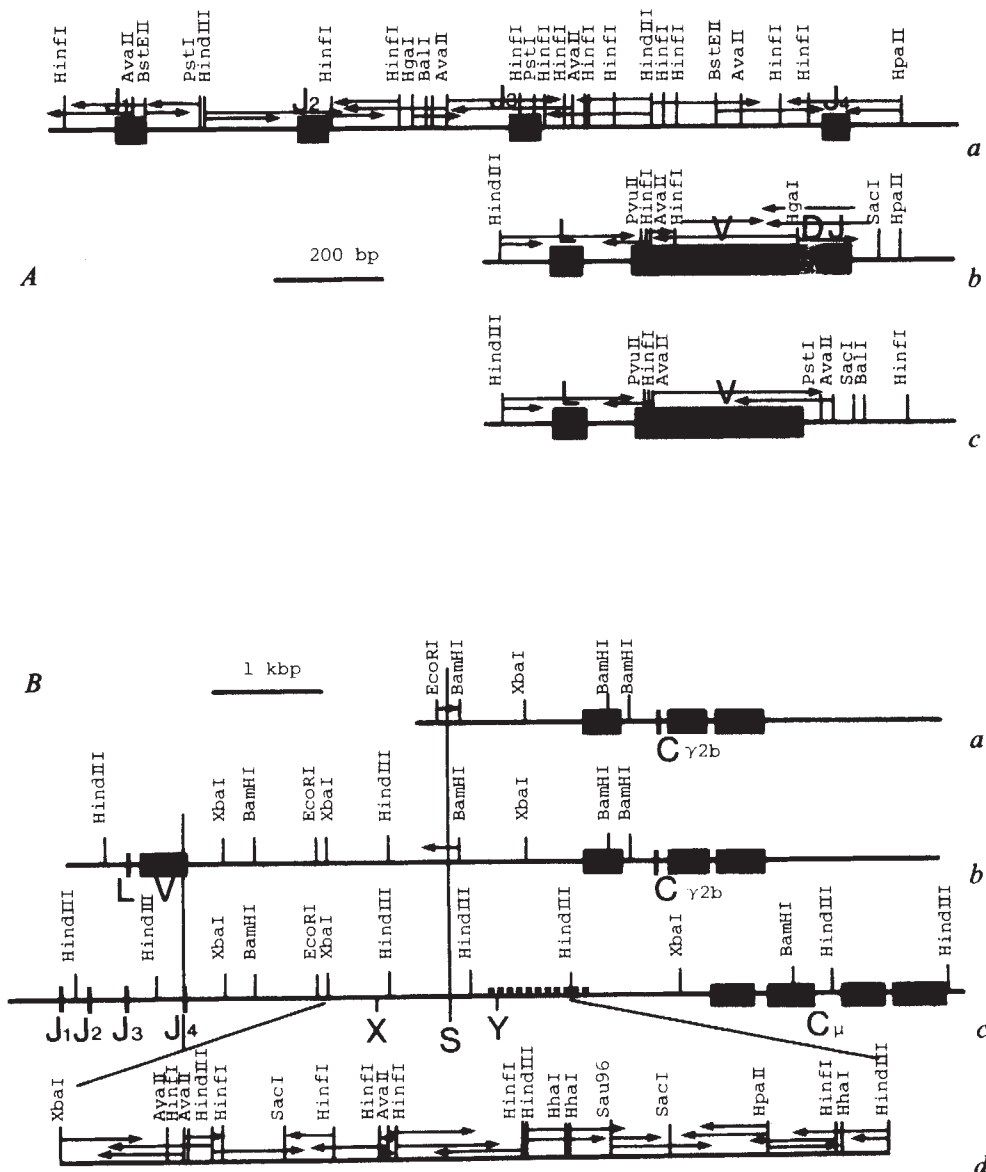
Figure 2 also shows the DNA sequence around the  $J_{M141}$  DNA segment of the embryonic clone MEP203. As described before, coding begins within the Leu 112 triplet and ends with the Ser 125 triplet. It thus seems that the 14-residue peptide comprising the third hypervariable region (HV3) beginning with Val 98 and ending with Thr 111 is encoded in neither the germ-line V gene nor the J DNA segment. This situation is in contrast to that of light-chain genes previously studied<sup>6-8</sup>, in which no such coding gap was observed. We suggest that the third hypervariable region of this heavy chain is encoded in one or more discrete DNA segments (D DNA segments for 'diversity'<sup>27</sup>) that lie elsewhere in the germ-line genome. The V, D and J DNA segments must be assembled by recombination to generate the complete V gene present in the myeloma cells.

### The exact joining ends of V and J DNA segments are unfixed

Close examination of the amino acid sequences around the third hypervariable region and the nucleotide sequences around the 3' end of the germ-line  $V_H$  gene and around the 5' end of the  $J_H$  DNA segments suggest several interesting features of the DNA joining events. As previously suggested for the light-chain J DNA segments<sup>7,13-15</sup>, the exact 5' boundary of  $J_H$  DNA segments seems to be unfixed (Table 1). For instance, in the case of  $J_{H4}$ , the coding of the MOPC141  $\gamma 2b$  chain begins with the second base of the Leu 112 triplet UUG as indicated in Fig. 2. In contrast, coding of the MOPC21  $\gamma 1$  chain and the MOPC173  $\gamma 2a$  chain by the same  $J_{H4}$  DNA segment seems to start with a Tyr triplet UAU and another Tyr triplet UAC, respectively. The first letters of the two Tyr triplets are 7 and 10 bases ahead (5' side) of the first base used for the coding of the MOPC141 J region. Similar examples for  $J_{H1}$ ,  $J_{H2}$  and  $J_{H3}$  are indicated in Fig. 4 and Table 1.

The DNA sequence of the clone PJ14 demonstrated that coding by the germ-line  $V_{M141}$  gene ends with the second letter of the Ser 97 codon AGC (Fig. 2). As shown in Table 1, the amino acid sequences immediately preceding Ser 97 are highly conserved among various V regions. Particularly noteworthy is

**Fig. 3** Restriction enzyme cleavage maps of part of the mouse DNA inserts in clones MEP 203 (Aa, Bc, Bd), MEP3 (Ba), M141-P21 (Ab, Bb) and PJ14 (Ac). Cleavage sites were determined by the end-labelling method of Smith and Birnstiel<sup>42</sup>, by the separation of single and double digests of DNA fragments by various restriction enzymes, and/or by DNA sequence determination. Horizontal arrows indicate the extent and direction of sequence determination. The filled boxes are exons identified by R-loop mapping and/or DNA sequencing. Vertical line S in B indicates the switch recombination sites for MOPC141 heavy-chain genes. X and Y in Bc are sequences similar to the  $S_{\mu}(\gamma 2b)$  sequence. Four J DNA segments,  $J_1$ ,  $J_2$ ,  $J_3$  and  $J_4$ , were identified by heteroduplex analysis of rearranged V-gene clones isolated from MOPC603, MOPC315, HOPC8 and MOPC173, respectively. The J DNA segment for MOPC104E was identified as  $J_1$  by R-loop mapping with  $\mu$ mRNA of MOPC104E. The region of clone MEP203, composed of repetitions of CTGAG or sequences closely related to it, is indicated by a broken line in Bc. bp, base pairs; kbp, kilobase pairs.





**Fig. 4** Nucleotide sequence of the 1.4-kilobase region of the embryonic clone MEP203 containing four J<sub>H</sub> DNA segments. The sequencing strategy is shown in Fig. 3Aa. The J<sub>H</sub> DNA segments are referred to as J<sub>1</sub>, J<sub>2</sub>, J<sub>3</sub> and J<sub>4</sub> from 5' to 3'. The amino acid sequences encoded by these J DNA segments are shown in italics. The sequences closely related to the palindromic heptamer, CACTGTG, and present in front of every J DNA segment, are underlined. The vertical lines indicate determined (solid lines) or predicted (broken lines) sites where coding by the J DNA segments begin in the various heavy-chain genes present in the indicated myelomas. The sites were determined by comparing the J DNA sequences with the sequences of the rearranged V DNA cloned from the respective myeloma. (See Fig. 2 for the M141  $\gamma$ 2b chain gene. The others are based on unpublished sequence data.) We have identified, isolated and sequenced two rearranged V genes, one from HOPC8 and the other from MOPC173. The V regions encoded by these rearranged V genes are different from the published ones, and are therefore referred to as H8-II and M173-II. It is not clear whether these V genes are part of the active heavy-chain genes of the respective myelomas. The putative sites for the heavy-chain genes of S107, M167, M315, A4, J539, M173 and M21 were deduced by comparing the J-peptide amino acid sequence of each heavy chain with the DNA sequence of the corresponding J DNA segment (Table 1).

the universal Cys at position 95. Tyr 93, Tyr 94, Ala 96 and Arg 97 are also well conserved. We therefore believe that the coding by most, if not all, germ-line V genes ends in the immediate vicinity of the 97th triplet. If so, coding gaps of various lengths exist in the majority of the cases listed in Table 1. As suggested above, these residues must be encoded elsewhere in the genome. The gaps in the heavy chains of four myelomas, A4, E109, U61 and A47N, are the shortest: only the first letter of Gly 98 is unaccounted for. Similarly, only several more bases are necessary to code for the HV3 of the J558  $\alpha$  chain. Previous studies on  $\kappa$ -chain genes indicated that the 3' end of a given germ-line V gene may not be precisely fixed and could shift within a range of several nucleotides<sup>7</sup>. If the same applies to the heavy-chain genes, the heavy chains with shorter HV3, such as those of A4, E109, U61 and J558, could be fully encoded by the respective germ-line V gene and the J DNA segment. It thus seems that in some cases the germ-line V DNA segment directly joins with the J DNA segment, as is the case in the light-chain genes. However, it is also possible that the assembly of the V<sub>H</sub> always occurs in two steps, but that in some cases the D segments are very short or do not even contribute to protein encoding at all.

In either case, by modulation of the joining sites and by various combinations of these DNA segments, diversity could be generated in the third hypervariable region of heavy-chain molecules. However, the location and number of the D DNA segments in the germ-line genome are unknown. They may be at the 3' end of each of the multiple germ-line V genes or clustered between the V DNA segments and the J DNA segments. In either case we assume that the order of the V, D and J DNA segments in the germ-line genome is the same as in the completely assembled V genes, so that successive joining of the three types of DNA segments is possible by the looping-out of the intervening spacers.

## Sequence characteristics around the V and J joining sites

We have previously reported that the five J<sub>K</sub> DNA segments and the one J<sub>L</sub> DNA segment are preceded by two blocks of short, conserved sequences: a palindromic hexamer interrupted by an AT base pair at the centre of symmetry, CACTGTG, and a nanomer, GGTTTTTGT (ref. 7). Similar sequences are also present 5' to all four J<sub>H</sub> DNA segments, although deviation from the basic sequences seems to be slightly higher for the J<sub>H</sub> segments than for the J<sub>L</sub> segments (Table 2). We also reported previously that germ-line V<sub>K</sub> and V<sub>L</sub> DNA clones contain inverted complements of the nanomer and heptamer in the 3'-noncoding regions. As shown in Fig. 2 and summarized in Table 2, the same or similar sequences are also present in the equivalent positions of the germ-line V<sub>H141</sub> gene clone. Based on the fact that these sequences are highly conserved, we propose that they constitute recognition signals for the putative recombinase. The finding of similar sequences near the heavy-chain gene DNA segments strengthens the validity of this hypothesis.

In the case of J<sub>H3</sub>, two sequences, CACAATG and CAATGTG, are closely related to the basic heptamer sequence, CACTGTG. Similarly, two sequences, TATTGTG and TACTATG, are present near the 5' end of J<sub>H4</sub> (Fig. 4). In these cases it is not clear which of the two sequences constitutes part of the proposed recognition signal.

The presence of the closely related conserved sequences in the corresponding positions of both light- and heavy-chain germ-line V genes and their J DNA segments suggests that all recombination events necessary for the generation of complete, somatic V genes are carried out by the same or a similar mechanism. We predict that the similar recognition sequences are also present on both sides of the D DNA segments.

An additional feature of the conserved signal sequences is the striking regularity in the length of the spacer between the heptamer and the nanomer. As summarized in Table 2, in almost all cases the spacers are 12 or 23 base pairs long. One exception is one of the two alternative cases of  $J_{H4}$ , where the spacer is 31 base pairs long. As a DNA double helix completes a turn every 10.4 base pairs<sup>28</sup>, the above indicates that the internal boundaries of the heptamer and the nanomer of a given DNA segment are orientated in nearly the same direction relative to the axis of the helix, regardless of the length of the spacer. The spacers of all  $V_{\kappa}$  segments correspond to one turn, and those of  $V_{\lambda I}$  and  $V_{\lambda II}$  to two turns. The spacer of  $V_{M141}$  is two turns long. All  $J_{\kappa}$  spacers are two turns whereas the  $J_{\lambda I}$  spacer is one turn. All spacers of the four  $J_H$  segments are two turns. The 31-base pair spacer, one of the two alternative spacers of  $J_{H4}$ , could be considered as an exceptional three turns. It should be emphasized that although the lengths of the spacers are conserved in a given type of gene segment, their sequences have diverged widely.

### A model of the joining enzymes

The above mentioned features around the joining sites suggest some predictions about the nature of the recombinases. In both  $\lambda$  and  $\kappa$  light-chain V-J joinings, one partner in the recombination displays the heptamer and nanomer one turn of the helix apart from the recombinase, and the other displays these sequences two turns apart. We predict that the recombinase contains two DNA-binding proteins: one recognizing the heptamer and nanomer one turn apart, and the other recognizing them two turns apart. As the two partners carry essentially identical recognition sequences, the two DNA-binding proteins would be structurally closely related.

The germ-line  $V_{M141}$  gene and all  $J_H$  segments show the recognition sequences separated by two-turn spacers. If the D segment contains recognition sequences with one-turn spacers on each side, then V-D and D-J joinings would be mediated by recombinases having components equivalent to those acting on the light-chain genes. In this way, all recombinations leading to the assembly of immunoglobulin V genes would follow a one-turn/two-turn spacer rule.

The two recombinase components hold the two recombining partners together and cut and rejoin the strands in the vicinity of the heptamers. As the recognition sequences of one partner are complementary to those of the other, it is possible that intra-strand base pairing occurs, however transiently, in the enzyme-DNA complex and facilitates the ligation reaction.

This rule prohibits some unwanted recombinations, such as those between  $V_{\lambda}$  and  $J_{\kappa}$ , and  $V_{\kappa}$  and  $J_{\lambda}$ . However, some other unwanted recombinations, such as those between  $V_H$  and  $J_{\lambda}$ , D and  $J_{\kappa}$ , and  $V_{\lambda}$  and  $V_{\kappa}$ , could occur. This may be the reason that the three gene families are distributed in three different chromosomes.

### The switch recombination sites for the complete $\gamma 2b$ gene of MOPC141

Our previous electron microscopic studies<sup>19</sup> demonstrated that the complete  $\gamma 2b$  gene clone M141-P21 shows some homology to the germ-line  $C_{\mu}$  gene clone MEP203 and some to the germ-line  $C_{\gamma 2b}$  gene clone MEP3. As shown in Fig. 1, the homology between clone M141-P21 and MEP203 is 2.4 kilobases long and extends, on M141-P21, from the 3' end of the V DNA segment to a site 1.2 kilobases 5' to the  $C_{\gamma 2b}$  gene, and, on MEP203, from the  $J_4$  DNA segment to a site 2.5 kilobases 5' to the  $C_{\mu}$  gene. The remaining part of the 3.6-kilobase V- $C_{\gamma 2b}$  intron of the complete  $\gamma 2b$  gene and the sequence that follows this part are entirely homologous to the germ-line  $C_{\gamma 2b}$  gene clone. The restriction enzyme maps of the three DNA clones shown in Fig. 3B confirm the results of the electron microscopic studies. Thus, the cleavage sites present in the first 2.2-kilobase portion of the V- $C_{\gamma 2b}$  intron are also present at the cor-

**Table 2** The two conserved blocks of sequences near the V-J or V-D-J joining sites

J DNA segments	Nanomers	Heptamers	Ref(s)
$J_{\kappa 1}$	GGTTTTGT 23	CACTGTG	7, 13
$J_{\kappa 2}$	<u>AGTTTTGT</u> 23	CAGTGTG	7, 13
$J_{\kappa 3}$	GGGTTTTGT 21	CACTGTA	7, 13
$J_{\kappa 4}$	GGTTTTGT 24	CACTGTG	7, 13
$J_{\kappa 5}$	GGTTTTGT 23	CACTGTG	7, 13
$J_{\lambda I}$	GGTTTTG <u>C</u> 12	CACAGTG	6
$J_{H1}$	<u>AGTTTTAGT</u> 22	GACTGTG	*
$J_{H2}$	GGTTTTGT 23	TAGTGTG	*
$J_{H3}$	<u>ATTTATTGT</u> 21	<u>CACAATG</u>	*
	23	CAATGTG	*
$J_{H4}$	GGTTTTGT 22	<u>TATTGTG</u>	*
	31	<u>TACTATG</u>	*
Basic sequence	GGTTTTGT	CACTGTG	
V DNA segments	Heptamers	Nanomers	Ref.
$V_{\kappa 21C}$	CACAGTG 11	ACAAAAACC	7
$V_{\kappa 21B}$	CACAGTG 12	ACAAAAACC	†
VK41	CACAGTG 12	ACATAAACC	8
VK2	CACAGTG 12	ACATAAACC	5
$V_{\lambda I}$	<u>CACAATG</u> 22	<u>TCAAGAACA</u>	6
$V_{\lambda II}$	CACAATG 23	ACAAGAACA	4
$V_{H141}$	CACAGTG 23	ACAAATACC	*
Basic sequence	CACAGTG	ACAAAAACC	

Two blocks of conserved sequences found in the 5'-flanking region of J DNA segments and in the 3'-flanking region of embryonic V DNA segments are compared. The numbers between the two types of sequences indicate the distance between them in base pairs. The bases different from those of the basic sequences in the corresponding positions are underlined. The refs from which the sequences were taken are listed (\* this paper; † G. Heinrich, personal communication).

responding positions of the germ-line  $C_{\mu}$  gene clone, whereas the cleavage sites located in the rest of the intron and the DNA stretch that follows are also found in the corresponding regions of the germ-line  $C_{\gamma 2b}$  gene clone. The maps show that switch recombination occurs between a pair of sites, one located within the 750-base pair *HindIII-HindIII* segment of the germ-line  $C_{\mu}$  gene clone and the other within the 200-base pair *EcoRI-BamHI* segment of the germ-line  $C_{\gamma 2b}$  gene clone, and generates the 650-base pair *HindIII-BamHI* segment present in the complete  $\gamma 2b$  gene clone (see Fig. 3B).

To understand this novel somatic event better, we determined the nucleotide sequences around the switch recombination sites (Fig. 5a). As expected, the 5' portion of the M141-P21 sequence is the same as the sequence of the germ-line  $C_{\mu}$  gene clone, MEP203, whereas the 3' portion is identical to the sequence of the germ-line  $C_{\gamma 2b}$  gene clone, MEP3. The exact  $\mu$ - $\gamma 2b$  switch recombination sites are indicated by a vertical line in Fig. 5a. We propose to refer to the  $\mu$ - $\gamma 2b$  switch sites on the germ-line  $C_{\mu}$  gene and the germ-line  $C_{\gamma 2b}$  gene as  $S_{\mu}(\gamma 2b)$  and  $S_{\gamma 2b}(\mu)$ , respectively.

### Sequences around the $\mu$ - $\gamma 2b$ switch sites

Are there any characteristic sequences around the  $\mu$ - $\gamma 2b$  switch sites that might be considered as part of the recognition signal for the recombinase? As shown in Fig. 5a, two short blocks of sequences, TCCTGG and AGA, present in front of (towards the 5' side of)  $S_{\mu}(\gamma 2b)$  are also present near the  $S_{\gamma 2b}(\mu)$  site in the equivalent positions. The two sets of sequences are in the same orientation relative to the direction of transcription, and thus, differ from those sequences conserved near the V-J joining sites not only in the sequence content but also in the relative



the matter. (Recently, we learned that the Y sequence is indeed used for the recombination of the  $\gamma 1$  gene sequenced by Kataoka *et al.*<sup>21</sup> (N. Takahashi and T. Honjo, personal communication).)

In the 0.8-kilobase region extending from the 5' side of the Y sequence to the *Hind*III site, a pentameric sequence, CTGAG, or its close variants are repeated in tandem. As the  $\Sigma\mu(\alpha)$  site reported by Davis *et al.*<sup>20</sup> falls within this region, involvement of the repeated sequence in the  $\mu$ - $\alpha$  switch recombination can be invoked. Indeed, our recent DNA sequencing study on the germ-line  $C\alpha$  gene clone indicates that the same pentameric sequences are prominent in the region where the  $C\alpha(\mu)$  site had been mapped by Davis *et al.* (H.S., unpublished observations).

Clone MEP203 carries a deletion of about 3 kilobases which seems to have occurred during the cloning or subsequent propagation of the clone in *Escherichia coli*. Restriction enzyme maps of several independent  $C\mu$  gene clones isolated by our laboratory and another laboratory<sup>20</sup> suggest that most clones carry deletions of various lengths that map between the two *Hind*III sites surrounding the Y sequence. As shown in Fig. 5b, this is the region where the pentameric sequences are repeated. Our recent DNA sequencing studies on two independent  $C\mu$  gene clones carrying different deletions confirmed that these deletions occur within the pentamer-rich regions.

## The two types of somatic recombination

The present studies deal with two types of somatic recombination required for the generation of complete heavy-chain genes, namely V-D-J or V-J joining and switch recombination. The two types of recombination are clearly different in several aspects, although they both occur somatically during lymphocyte differentiation. Although V-D-J or V-J joining should precede the synthesis of the free  $\mu$  chain detected in pre-B cells<sup>31</sup>, switch recombination almost certainly occurs after these cells mature to small B cells bearing complete IgM molecules on the cell surface. V-D-J or V-J joining takes place at the margins of two protein-encoding DNA segments and generates complete V genes. In contrast, switch recombination occurs in the sequences located away from the protein-encoding DNA segments and in effect replaces the  $C\mu$  exons of the complete  $\mu$ -chain gene with exons coding for the C regions of other heavy-chain classes or sub-classes, except, perhaps, for the  $\delta$  class.

In V-D-J or V-J joining, the recombinant DNAs are to be translated in the recombined regions. This imposes certain restrictions on the precise cutting and joining sites because the encoded polypeptide chains would have to fold properly to be able to function as immunoglobulin subunits. Indeed, this type

of recombination seems to occur within a short (several base pairs) region adjacent to the conserved palindromic heptamer present at the margin of every V and J DNA segment studied to date. However, the recombination site does not seem to be unique in a given gene segment and this flexibility, however limited, seems to have been exploited during evolution to increase the genetic capacity of the organism's antibody repertoire. Switch recombination can *a priori* occur anywhere within the 7.5-kilobase intron separating the J DNA segments and the  $C\mu$  exons, and anywhere within a several-kilobase region located 5' to the respective C exon of other classes, as long as the transcription unit on the recombined DNA can encompass the V and C exons and the transcript can be processed properly to the mature mRNA. This suggests that switch recombination sites may be scattered widely in these regions. Indeed, the  $\Sigma\mu(\gamma 2b)$  site for the MOPC141  $\gamma 2b$  chain and the  $\Sigma\mu(\alpha)$  site for the MOPC603  $\alpha$  chain are at least 500 base pairs apart. Furthermore, our recent studies demonstrated that in the *Eco*RI digests of various myeloma DNA, the sizes of the DNA fragments detected by the 3-kilobase *Xba*I-*Xba*I fragment carrying the switch region (see Fig. 3Bc) are diverse, even among those myelomas secreting the heavy chains of the same class or subclass (R.M., unpublished observations). This suggests that one or both switch recombination sites are not class or subclass specific.

We previously proposed that the V-J joining can be considered as a reversal of an ancient, accidental insertion of an IS-like DNA element carrying invertedly repeated sequences at the margins<sup>7</sup>. In contrast, the evolution of immunoglobulin chains suggests that the C exons of various heavy-chain classes or subclasses arose by duplication of ancient  $C\mu$  exons and their flanking sequences<sup>44</sup>. The switch recombination sites have probably arisen from these duplicated flanking sequences by making use of the sequence homology among them. Whether the switch to various classes or subclasses is actively controlled at the level of the DNA recombinations remains to be determined.

Early *et al.*<sup>33</sup> have independently reported the analysis of M603  $\alpha$  heavy-chain V DNA and its germ-line V DNA and J DNA. They found that the 18-base pair D DNA segment in M603 V DNA is not accounted for by either germ-line V DNA or J DNA. They also identified two heavy-chain J DNAs by nucleotide sequencing that correspond to our  $J_{H1}$  and  $J_{H2}$ . Our determined or predicted joining sites on  $J_{H1}$  DNA for M603, M167 and S107 heavy-chain genes are in agreement with their direct sequencing of cDNAs for heavy-chain mRNAs prepared from respective myelomas.

We thank Dr Michael Potter for myelomas and Ms L. Sultzman, Mr G. Dastoornikoo, Ms L. Gibson, Mrs P. Riegert and Mr A. Traunecker for technical assistance.

Received 5 May; accepted 23 June 1980.

- Hozumi, N. & Tonegawa, S. *Proc. natn. Acad. Sci. U.S.A.* **73**, 3628-3632 (1976).
- Tonegawa, S., Hozumi, N., Matthysens, G. & Schuller, R. *Cold Spring Harb. Symp. quant. Biol.* **41**, 877-889 (1976).
- Brack, C., Hiram, M., Schuller, R. & Tonegawa, S. *Cell* **15**, 1-14 (1978).
- Tonegawa, S. *et al. Proc. natn. Acad. Sci. U.S.A.* **74**, 3171-3175 (1978).
- Seidman, J. G. *et al. Proc. natn. Acad. Sci. U.S.A.* **75**, 3881-3885 (1978).
- Bernard, O., Hozumi, N. & Tonegawa, S. *Cell* **15**, 1133-1144 (1978).
- Sakano, H., Hüppi, K., Heinrich, G. & Tonegawa, S. *Nature* **280**, 288-294 (1979).
- Seidman, J. G., Max, E. E. & Leder, P. *Nature* **280**, 370-375 (1979).
- Lenhard-Schuller, R., Hohn, B., Brack, C., Hiram, M. & Tonegawa, S. *Proc. natn. Acad. Sci. U.S.A.* **74**, 4709-4713 (1978).
- Seidman, J. G. & Leder, P. *Nature* **276**, 790-795 (1978).
- Gilmore-Herbert, M. & Wall, R. *Proc. natn. Acad. Sci. U.S.A.* **75**, 342-345 (1978).
- Schibler, U., Marcu, K. B. & Perry, R. P. *Cell* **15**, 1495-1509 (1978).
- Max, E. E., Seidman, J. G. & Leder, P. *Proc. natn. Acad. Sci. U.S.A.* **76**, 3450-3454 (1979).
- Weigert, M., Gatmaitan, L., Loh, E., Schilling, J. & Hood, L. *Nature* **276**, 785-790 (1978).
- Weigert, M. *et al. Nature* **283**, 497-499 (1980).
- Pernis, B., Forni, L. & Amante, L. *Ann. N.Y. Acad. Sci.* **190**, 420-431 (1971).
- Pernis, B., Forni, L. & Luzzati, A. L. *Cold Spring Harb. Symp. quant. Biol.* **41**, 175-183 (1976).
- Cooper, M. D., Kearney, J. F., Lydyard, P. M., Grossi, C. E. & Lawton, A. R. *Cold Spring Harb. Symp. quant. Biol.* **41**, 139-145 (1976).
- Maki, R., Traunecker, A., Sakano, H., Roeder, W. & Tonegawa, S. *Proc. natn. Acad. Sci. U.S.A.* **77** (in the press).
- Davis, M. M. *et al. Nature* **283**, 733-739 (1980).
- Kataoka, T., Kawakami, T., Takahashi, N. & Honjo, T. *Proc. natn. Acad. Sci. U.S.A.* **77**, 919-923 (1980).
- Kabat, E. A., Wu, T. T. & Bilofsky, H. in *Sequences of Immunoglobulin Chains* (NIH, Bethesda, 1979).
- Breathnach, R., Benoist, C., O'Hare, K., Cannon, F. & Chambon, P. *Proc. natn. Acad. Sci. U.S.A.* **75**, 4853-4857 (1978).
- Rabbitts, T. H., Matthysens, G. & Hamlyn, P. H. *Nature* **284**, 238-243 (1980).
- Davis, M., Early, P., Calame, K., Livant, D. & Hood, L. in *Eucaryotic Gene Regulation* (eds Axel, R., Maniatis, T. & Fox, C. F.) 393-406 (ICN-UCLA Symp., Academic, New York, 1979).
- Cory, S. & Adams, J. M. *Cell*, **19**, 37-51 (1980).
- Schilling, J., Cleavinger, B., Davie, J. M. & Hood, L. *Nature* **283**, 35-40 (1980).
- Wang, J. C. *Proc. natn. Acad. Sci. U.S.A.* **76**, 200-203 (1979).
- Salsano, P., Froland, S. S., Natvig, J. B. & Michaelsen, T. E. *Scand. J. Immun.* **3**, 841-846 (1974).
- Fu, S. M., Winchester, R. J. & Kunkel, H. G. *J. Immun.* **114**, 250-261 (1975).
- Burrows, P., LeJeune, M. & Kearney, J. F. *Nature* **280**, 838-840 (1974).
- Rao, D. N., Rudikoff, S., Krutzsch, H. & Potter, M. *Proc. natn. Acad. Sci. U.S.A.* **76**, 2890-2894 (1979).
- Early, P. *et al. Cell* **19**, 981-992 (1980).
- Kehry, M., Sibley, C., Fuhrman, J., Schilling, J. & Hood, L. *Proc. natn. Acad. Sci. U.S.A.* **76**, 2932-2936 (1979).
- Vrana, M., Rukikoff, S. & Potter, M. *Proc. natn. Acad. Sci. U.S.A.* **75**, 1957-1961 (1978).
- Francis, S. H., Leslie, R. G. Q., Hood, L. & Eisen, H. N. *Proc. natn. Acad. Sci. U.S.A.* **71**, 1123-1127 (1974).
- Milstein, C., Adetugbo, K., Cowan, N. J. & Secker, D. S. *Prog. Immun.* **2**, 157-168 (1974).
- Bourgeois, A., Fougereau, M. & Depreval, C. *Eur. J. Biochem.* **24**, 446-455 (1972).
- Rudikoff, S. & Potter, M. *Proc. natn. Acad. Sci. U.S.A.* **73**, 2109-2112 (1976).
- Rudikoff, S. & Potter, M. *Biochemistry* **13**, 4033-4038 (1974).
- IUPAC-IUB Commission J. *Biochem.* **113**, 1 (1969).
- Smith, H. O. & Birnstiel, M. L. *Nucleic Acids Res.* **3**, 2387-2398 (1976).
- Maxam, A. & Gilbert, W. *Proc. natn. Acad. Sci. U.S.A.* **74**, 560-564 (1978).
- Hill, R. L., Delaney, R., Fellows, R. E. Jr & Lebovitz, H. E. *Proc. natn. Acad. Sci. U.S.A.* **56**, 1762-1769 (1966).