## Exon-intron organization and complete nucleotide sequence of a human major histocompatibility antigen $DC\beta$ gene

(major histocompatibility complex/class II antigen)

Dan Larhammar, Jens Jørgen Hyldig-Nielsen, Bo Servenius, Göran Andersson, Lars Rask, and Per A. Peterson

Department of Cell Research, The Wallenberg Laboratory, University of Uppsala, Box 562, S-751 22 Uppsala, Sweden

Communicated by Jean Dausset, August 15, 1983

ABSTRACT We have determined the complete nucleotide sequence of a human class II histocompatibility antigen  $DC\beta$  gene. The gene spans more than 7 kilobases and contains five exons corresponding to the different domains of the DC $\beta$  polypeptide. The exon-intron organization is thus analogous to that of class II antigen  $\alpha$ -chain genes, class I antigen heavy chain genes, and the constant parts of immunoglobulin genes, emphasizing further the evolutionary relationship among these molecules. The mature polypeptide deduced from the DC $\beta$  gene shows 93% and 88% homology, respectively, to sequences derived from two DC $\beta$  cDNA clones of other haplotypes. The allelic polymorphism of  $DC\beta$  chains resides predominantly in the first extracellular domain, whereas the rest of the polypeptide is virtually constant. The exons of the DC $\beta$  gene display high homology to the corresponding exons of a murine I-A $\beta$  gene. Also, the introns show significant homology. The DC $\beta$  chains lack eight amino acids in the cytoplasmic tail, as compared to DR and I-A  $\beta$  chains. This is probably due to a nonfunctional splice junction of  $DC\beta$  genes, causing a separate cytoplasmic exon to be nonexpressed.

The D region of the human major histocompatibility complex (MHC) harbors genes involved in the regulation of the immune response (1). The known products of this region, the class II antigens, are polymorphic glycoproteins composed of two dissimilar subunits: an  $\alpha$  chain of approximately 34,000 daltons and a  $\beta$  chain of approximately 28,000 daltons (2). The polymorphism is primarily carried by the  $\beta$  chains (3). Class II antigens are expressed on the surface of cells involved in various aspects of an immune response such as B lymphocytes, subpopulations of T lymphocytes, and macrophages (4).

Three different class II heterodimers, denoted DR (2), DC (5), and SB (6), have been identified by immunochemical techniques and primed lymphocyte typing. Also a fourth antigen called BR has been proposed (7). Southern blot analyses with fragments derived from  $\beta$ -chain cDNA clones provide evidence that the class II family contains more than three  $\beta$ -chain genes (8). The exact function of the different class II antigens in cell interactions is as yet unknown.

Through cloning of cDNAs corresponding to class II antigen  $\alpha$  and  $\beta$  chains, we have recently shown that the two chains display structural homology to each other as well as to class I antigens and immunoglobulins (9–11), suggesting that these immunologically important proteins have evolved from a common ancestor by gene duplications.

As a step towards the understanding of the evolution and mechanisms for generation and maintenance of the polymorphism of class II antigens, we present here the nucleotide sequence of a human class II antigen  $\beta$ -chain gene.

## **MATERIALS AND METHODS**

Isolation of Genomic Clone. A genomic library was constructed from DNA donated by an HLA-homozygous human individual typed to be Dw4/DR4 (unpublished data) in the cosmid vector pHEP (12). The library was screened with a 627-basepair (bp) Ava I fragment containing almost the entire coding part of the DC $\beta$  cDNA clone pII- $\beta$ -1 [previously denoted pDR- $\beta$ -1 (9)]. Screening of the cosmid library and growth and analysis of cosmid clones were performed as described (13), as were Southern blot analyses (14). DNA probes were labeled by nicktranslation (15).

DNA Sequence Determination. Nucleotide sequences were determined with the chemical degradation procedure (16) and a modification of the dideoxy chain termination method (17) using exonuclease III and synthetic oligonucleotide primers (unpublished results).

## **RESULTS AND DISCUSSION**

Isolation of a Cosmid Containing a DC $\beta$  Gene. Genomic DNA from a panel of DR-homozygous human individuals has been analyzed in Southern blotting experiments using a fragment from a DC $\beta$  cDNA clone as probe (8). All blots with DNA digested with different enzymes show one strong constant band and one or two strong polymorphic bands. The  $\beta$ -chain gene of cosmid clone cosII-102 accounts for the strong polymorphic DC $\beta$ bands of the genomic donor DNA. Thus, this clone was chosen for further characterization.

Exon-Intron Organization of the DCB Gene. Two overlapping fragments containing the DC $\beta$  gene [a 7.7-kilobase (kb) EcoRI fragment and an 11-kb BamHI fragment] were subcloned in pUC9 to facilitate sequence determination (see Fig. 1). A restriction map of the DC $\hat{\beta}$  gene and the sequencing strategy are shown in Fig. 1. Exons were localized by comparison with the DC $\beta$  cDNA clone pII- $\beta$ -2 (unpublished results). The nucleotide sequence and the translated amino acid sequence are shown in Fig. 2. The DC $\beta$  gene encompasses more than 7 kb and contains five exons correlating with the different domains of the DC $\beta$  polypeptide (see Fig. 3). The first exon corresponds to the 5' untranslated region, the signal sequence, and four amino acids of the first domain. The remainder of the first domain and the second domain are encoded by exons 2 and 3, respectively. Exon 4 encodes the connecting peptide, the membrane-spanning segment, and six amino acids of the cytoplasmic tail. The last four amino acids of the cytoplasmic tail are encoded by a separate exon also containing the 3' untranslated region. All splice junctions conform to the G-T-A-G rule (18). Overall, the exon-intron organization is analogous to that of genes for class II antigen  $\alpha$  chains (refs. 19 and 20; unpub-

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "advertisement" in accordance with 18 U.S.C. §1734 solely to indicate this fact.

Abbreviation: bp, base pair(s).

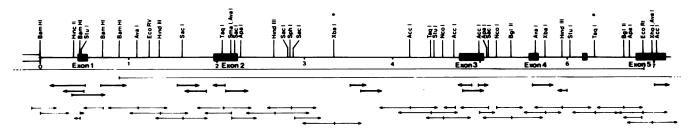


FIG. 1. Restriction map of the  $DC\beta$  gene of cosmid clone cosII-102. Exons are shown as filled boxes. The hatched box shows the nonfunctional cytoplasmic exon, whereas the open box corresponds to the cosmid vector. Subcloned fragments are indicated by bars below the map. Sequencing strategy is shown by bold arrows for the dideoxy chain termination method with synthetic primers and fine arrows for the chemical degradation procedure. At sites indicated by an asterisk, no overlapping sequence was obtained. Scale shows distances in kb.

lished data) and class I antigens (21, 22), as well as to the constant parts of immunoglobulin genes (23). This finding, together with the previously described structural homologies between the class II, class I, and immunoglobulin polypeptides (9-11), suggests that these molecules, all of which are of central importance in immune responses, have evolved from a com-

Operational Control Contre Control Control Control Control Control Control Cont	• • • • • • • • • • • • • • • • • • • •	
MARTINE TO COMPARISATION TO CONTROL TRADITION TO CONTROL AND CONTRA	GGATCCCCACTTAATTTGCCCTACTGAAAGAATC <u>CCAAGTAT</u> AAAAAACAACCAGTTTTTAATC <u>AATATTA</u> CAAAGATGTTTACTGTTGAATCG <u>CATTTTT</u> CTTTGGCTTCTTAAAATCCC	120
The Ant CALARANCE AND ADDRESS	TROGCATTCAATCTTCAGCTCTTCTATAATTGAGAGGAAGTTTTCACCTCAAATGTTCATCCAGTGCAATTTGAAGACGTCACAGTGCCAGGGACTGGATTGAGAACCTTCACAAAAAA	240
OTHER MERICAN CONSTRUCTION OF CONCENTRATION TO THE TABLE AND AND CONTROL AND CONTRO		360
GEA ACT OF ACT THE ATH CHE GO AND CHE GOA ACT CHE GAE ACT CHE GAE GAE ACT THE CHE GYMANDEAGACHTECHTECHANDER       562         CECANTINGCHANDEXCHTTERCHANDETUCICUMATICEANTETUCIANACHTECHTENTACTACHTECHTENTERCANDEANDEAUNGUTURANTECHANDEANDEA       692         CATTINETTITIKARGANDECHOOTDAGECCENTECCEANDECTIMECCEGAAATTITAGEGABAGABAAANGTUTAÉNTITICEANDEANDEAUNGUTURANTECHANDEANDEAUNGUTURANTECHANDEANDEAUNGUTURANTECHANDEANDEAUNGUTURANTECHANDEANDEAUNGUTURANTECHANDEANDEAUNGUTURANTECHANDEANDEAUNGUTURANTECHANDTECHANDEAUNGUTURANTECHANDTECHANDEAUNGUTURANTECHANDTECHANDEAUNGUTURANTECHANDTECHANDEAUNGUTURANTECHANDTECHANDTECHANDEAUNGUTURANTECHANDTEUNGUTURA	GTACATCAGATCCATCAGGTCCAAGCTGTGTGACTACCACTACTTTTCCCTTCGTCTCAATT ATG TCT TGG AAG AAG GCT TTG CGG ATC CCT GGA GGC CTT CGG	
CATTINETTITION GRAND AGA CANTOGUNE CCCARTCE CANAGE CENTRE CENTRE CONSTRAINED AND AGAIN	Val Ala Thr Val Thr Leu Met Leu Ala Met Leu Ser Thr Pro Val Ala Glu Gly Arg Asp Ser Pro G GTA GCA ACT GTG ACC TTG ATG CTG GCG ATG CTG AGC ACC CCG GTG GCT GAG GGC AGA GAC TCT CCC G GTAAGTGCAGGGCCACTGCTCTCCAGAGCC	•
MOTEGEOCRATICUTATIVITAGENERALCIAETICIC CEGEGEOCOCACEGEOCOCACIÓN COCTOCICCICACIÓN COCTOCICACIÓN COCTOCICAC	GCCACTCTCGGGAACAGGCTCTCCTTGGGCTGGGGTAGGGGGATGGTGATCTCCATGATCTCGGACACAATCTTTCATCAACATTTCCTCTCTTTGGGGAAAGAGAACGATGTTGCATTCC	682
TOTALGOCINGCINT GUARTETIC CONTACT TATTETIC CONTACT CONT	CATTTATCTTTTAGTGATGAGGTGAGGCACAGTCGGATCCCCATCCTACAGGCTTAAGCCTGGAATATAAGGAGAGAGA	802
ACAGACAGGAGTATTA CAGATACAGTATTATOGGAATTGATGACTTCAGTOGGÀTTCAGAGACCCGAGTTGATGACACAATTATACAGAGGAAGAATTAA 1152 CATATATTGTGAAACCACTCATTTCAGAGACAGCACAATTATTTAGATAAATTGTCTCTCATGTTGTGATGTGGACTGAGTATGATGATGTTTTACACAATTATATGTT CATATATTGTGAAACCACTCATTTCAGAGACAGCACAAGCTTATTGTCAAATTGTTCTCATGATTGAGAGTAGGAGTATTTACCAGGAGCAATTAATT	AAGTGAGCATGGGGTTATTTTTGAAGATACGAATATCCCAGAGACACAGCAGGATTTGTCATTTAGGCGTGCCCCCAAGACTTTGCCTGACTAAATATTATGGGATCCTGCATTGGGAAA	922
CATATATIONAMICALTERTICAGEMENGENERGE       1282         ANATOGRATHETETTICAGEMENGENERGE       1282         ANATOGRATHETETTICAGEMENGENERGENERGENERGENERGENERGENERGENER	TOTANGCCAGCAATGGTGTCTGTAGTCTCCGTATTTGAGGAAAAGTTGTCTGTATTCCTGACTGA	1042
ATATOGATATICTOTICATACATACTOCCCCTAGCTGGGCCGCCACAGCTTAATTGGAATCTAGTTATCAAAATTGAAATCGAGCATGGCAAGCGGGGGGGCCACAGAAATGGCCATATTCTTCTTCTTGGGGCGGCGGGGGGGG	ACAGACAGGGAGTCTTCAGGTTTCACTGATTTATGGGCAAATGGTGACTTCAGTGGGATTCAGAGACCCGAGTTGGTGGACTGAATTTAGCAGAAAGGAGGATGTAAAGAAGGGAAATAA	1162
CGAAGGUCTAAGTUCTTICTAAAAAAGTUAACUCCAACUCCAACUCCAACUCCAACUCGAACUCCAACUCAGGUCAAGUCAAGUCAAAAAAAA	CATATATTGTGAAACCACTCATTTCAGACACAGCACAGC	1282
AGGAGACAGGCTTANTTACTGGACCGTCTCATCATCACCTGAACTGACAGGTTATGGGATAATTTATCTCGGAGGCTGGAGGGACAGGAAAATGGAATTCATGGAGTTTATGGGAGAGCGAGGAGAGAGA	ATATGGATATTCTGTTACATACCTGCCCTAGCTGGTGACTGCCACAGCTTAATTGGAATCTAGTTTATCAAAATCAAAAGCTTGTGCTCTTTCCATGAATAAATGTTTCTTCTAGGACT	1402
ACCCCCCCANTCCCCTANGTOCAGAGGTCTATOTAÁAATCAGCCCGACTGCCCCCCCCCCCCCCCCCCCCCCC	CGGAGGTGTAGGTCCTTTCTAACATAAAAGTGAGTGAACCTCACATGGCATCGGAAGGGTAAATCCAGGCATGGGAAGGGAGGTATTTTACCGAGGGACCAAGAGAATACGCATATCAG	1522
GRITICUCGAMCCCCCAGAGAGGGGGGGAGGGCAGGGCGGGGCGG	AACCAGGACAGGCTTAATTTCTGGACCCGTCTCATCATTCCCTTGAACTCACAGGTTTATGTGGATAATTTTATCTCTGAGGTTTCCAGGAGCTCAATGGAAAATGGGATTTCATGCGAG	1642
CTCCGGGCCGGGGCGGCGGGGGGGGGGGGGGGGGGGGG	ACCCCCCTGATTCCCTCTAAGTGCAGAGGTCTATGTAAAATCAGCCCGACTGCCTCTCCCTCGGTTCACAGGCTCCGGCAGGGACAGGGCCTTTCCCGCCCTTTCCTGCCTG	1762
CTCC0000C000FCL0000C000CC000C000C000C000C0	GATTCCCGAAGCCCCCAGAGAGGGGGGGGGGGGGGGGGG	1882
Asp Val Gly Val Tyr Arg Ala Val Thr Pro Leu Gly Pro Pro Ala Ala Glu Tyr Trp Asn Ser Gln Lys Glu Val Leu Glu Arg Thr Arg72 2173Ala Glu Leu Asp Thr Val Gys Arg His Asn Tyr Gln Leu Glu Leu Gu Cre Gee Gee Gee Cre Gee Cae Cre Gee Cae Cae Cre Gee Gee Gee Cre Gee Cae Cae Cre Cre Gee Cae Cae Cae Cre Gee Cae Cae Cre Gee Cae Cae Cre Gee Cae Cae Cre Gee Cae Cae Cae Cre Gee Cae Cae Cre Gee Cae Cae Cre Gee Cae Cae Cre Gee Cae Cae Cae Cae Cae Cae Cae Cae Cae C	lu Asp Phe Val Tyr Gin Phe Lys CTCCGGGCCGGGTCAGGGCGGGCGGGCGGGGCGGGGCCGGGGCCGGGGCCGGGGCCGGGCGGCGGGCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCCGGGCCGGGCCGGGCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCCGGGCGGGCCGGGCCGGGCCGGGCGGGCGGGCGGGCGGGCGGGCGGGCGGGCCGGGCGGGCGGGCGGGCGGGG	
GAC GTC GGC GTC TAT CGC GCC GTC ACC CCC CTC GCC GCC GCC GCC GCC GCC GCC G	Gly Met Cys Tyr Phe Thr Asn Gly Thr Glu Arg Val Arg Leu Val Thr Arg Tyr Ile Tyr Asn Arg Glu Glu Tyr Ala Arg Phe Asp Ser GGC ATG TCC TAC TTC ACC AAC GOG ACG GAG CGC GTG CCT CTT GTG ACC AGA TAC ATC TAT AAC CGA GAG GAG TAC GCA CGC TTC GAC AGC	
GCG GAG TTU GAC ACG GTG[ <u>ITUC</u> ]AGA CAC AAC TAC CAG TTA GAG TTA GAG TTA GAG CTA CGC ACG ACC TTA CAG CGA G GTGAGCGTCATCGCCCCCCCCCC	Asp Val Gly Val Tyr Arg Ala Val Thr Pro Leu Gly Pro Pro Ala Ala Glu Tyr Trp Asn Ser Gln Lys Glu Val Leu Glu Arg Thr Arg GAC GTG GGG GTG TAT CGG GCG GTG ACG CCG CTG GGG CCG CCT GCC GAC GAC TAC TGG AAC AGC CAG AAG GAA GTC CTG GAG AGG ACC CGG	
TGGTGCATCGGAGGGCAGGGCACCTAGGCAGGGGGACAGGCAGAGGTTGGTCAAGCTGCCTAGTTTCGCCCCATCCTCCCCGTCGGCCTGGCCTCGCCCCCCCC	Ala Glu Leu Asp Thr Val Cys Arg His Asn Tyr Gln Leu Glu Leu Arg Thr Thr Leu Gln Arg Arg V GCG GAG TTG GAC ACG GTG TGC AGA CAC AAC TAC CAG TTG GAG CTC CGC ACG ACC TTG CAG CGG CGA G GTGAGCGTCGTCGCCGTCGCCGAGGCCCA	
TGCCTCGTGCCTTATGCGTCCCGGGCCTACCTTTACCTAAGCAGTTCTCTCGGCCCCCAATTTCCGCCCGC	TCCTTGGCAGGGGCCCAGAGTCTCTGCCGCAGGGGGGGGG	2390
AGGTCCACCTACACAGGTCAATGCGGAAGCTTCAGACTTGCGCCTGATGGAGTTAGGGCTGCCCCACACAGTGGCGAGCGCATCCAGCAATTACAGTGTTGTAAATAAGAATATTTGA2750CTITTGACTTCAAATTAATGATCGTAATTCATGGTTTTCTTAAATGGCTCCAGTGCGGAGCGCGCCCCACAGGGAATGCAGAGGAGAAAAAAGGGCTTAGAATGGACCAATGGA2870CTGGCATGTGGTATGAGCTCAATGATCTTTGTTAAATTAATGAATAAATGGCTCAGCTGGCGAGCGCCAATTAGGGGAGAGAGA	TOGTOCATCOCAGOGGCAGOGACCTAGOGCAGAGCAGGGGGACAGGCAGAGTTOGTCAAGCTGCCTAGTTTCGCCCCATCCTCCCCGTCCGTCC	2510
AGGTCCACCTACACAGGTCAATGCGGAAGCTTCAGACTTGCGCCTGATGGAGTTAGGGCTGCCCCACACAGTGGCGAGCGCATCCAGCAATTACAGTGTTGTAAATAAGAATATTTGA2750CTITTGACTTCAAATTAATGATCGTAATTCATGGTTTTCTTAAATGGCTCCAGTGCGGAGCGCGCCCCACAGGGAATGCAGAGGAGAAAAAAGGGCTTAGAATGGACCAATGGA2870CTGGCATGTGGTATGAGCTCAATGATCTTTGTTAAATTAATGAATAAATGGCTCAGCTGGCGAGCGCCAATTAGGGGAGAGAGA	TGCCTCGTGCCTTATGCGTTTGCCTCCTCGTGCCTTACCTTTACTAAGCAGTTCTCTCTC	2630
CTTITGACTICAAATTAITAITCATCGTAATTCAGTITTCTTAAATGGCCCCCATTCATGGCGGGGCCCCTTTGAGATGAGGGGGGGG	· · · · · · · ·	
CTGGCATGIGGTATGAGCTCAATGATCTTCTGTAAAATAAATGAATAAATGTGCTCAGCTGCCAATCCACTTAGGGCTCAAGGGAAAGCAGAGGATAAATAGAGCCTTAAAAAAAGGGCTT2990TATTAATTATTTTCTGTCATTTTGCTTAAATTCTTAAAGTAAAACTCTTATTGACTGGGATCTAATAGAGGTTGGGAATACAAAGTCGAGGAAAAAAAGCCTTAAAAAAAA		2870
TATTANTATTTTCTGGCATTTTGGCTTANTTCTTANAGTANACTCTTATTGACTGGATCTTANTAGAGTTTGGGATACAAAGTCGAGGAAGAAAAAGTGTTGGGAAAAAAGTGTTGGGAAAAAA		
ACACCTGAATGATGTTTGTAAGGCAGGTTTTAAATTTCTTAGAAAAGCTGAACAAATGGCACAATGCAAAGGCAGAAGTTTTGGAATAAATA		3110
ATAGTITCAGGITGCITTIGGCTTAGGTTCCCCCCCCCCC		
CTTATCAAAATACCGTCTATGTTACGTCATTTAATCTCACAGTTGCCTGTCCATTAGAGATTAGCATCACCACTTTATATATCGTAATATTAGTACATGACAAAACACTTTAAGTAATCAGC       3470         CCACAAGTACTCACCAAGACCTTAAGCCTCCCAAAGTACACCAATATTCTTTAAGTCTTCACTACACACTTTAAAGGGCAAAGGGACATAAAGCCTTGTTAAAGCCAAGTAATGCAAAAGCCTAGTTAAGCCAAGTACACCAAGGAGAACCAAAGGAGAAACCAATGCAAAGGGACATCAAAGGGACATCAAAGGGACATAAAGCCTTGTTAAAGCCTGGTCGACTCCAGGGAGACCAAAGGAGAACCAAAGA       3590         GCAGCAAAGAACTTTTAGTTTACTTTAAAACTCCCTTGCCTTGCCTGCC		
CCACAMGTACTCACCAMGACCTTANGCCTCCCAMAGTACACAMTATTCTTTATGTTCTTCACTACACACATCTATAGAGTCAMAGGGACATAAAGCCTTGTTAAAGCCAGTTTTGACTAGAA       3590         GCAGCAATGAGTCTCTTCCTGTTGATCCCATGTTAATGGGACAAAATGATACTTTCAAGGCATTGAAAATTTATGGATTAATCAATC		
GCAGCANTGAGTCTCTTCCTGTTGATCTCCATGTTAATGGGACAAAATGATACTTTCAAGGCATTGAAAATTTATGGATTAATCAATC		3590
ATTCACAAAACTTITAGTTTACTITTAAACTCCCTTGCCTTCTTTTGACTCACATCGTAGTGCCAGCAAGTACTTACATTTTGCTTAITTTGGGTCTATTCCATAAAATTTATTTTAT		
ТСТТТСТСАТАЛОТТГОТОСССТСТАТТТТАСТСССЛОГТСТОТОТАЛАЛАТСТАЛАЛАТСТТАТАЛОСОСАСАТАССТОЛОТТСАСОТОСАСТССАХОЛЛОСОАСААЛАСС АЛАОТТСАЛОТССАЛОСАСАЛАССОТОАТТССТТСССОАТОЛТСАТОЛОТСАЛОДОТОССТТТАЛАТОСОСТСАЛСТСАССТСАССАЛАТСТСАССТАХОЛЛОСОАЛАСС СЛОССТСАТТАЛССТЛОССТОЛСТСССАЛОДОТСТСАСАЛАТАТТТТСТАСЛОДАТАЛСАТАЛСТАЛАТАТСТОЛТТСАСОЛСАЛАЛАТСТСАСАТАТАТСТОЛОСОАЛОС		
AAAGTTCAAGTCCAAGCACAAAACCGTGATTCCTTCCGGATGATGGCTCAAGAGTGCCTTTTAATTGGGGTGCAACCTGCTGACCTCAGCAAATCTCAGGCTATATTTATATGTTCACAATTA CAGGCTCATTAACCTAGGCTGATCTCTGCAAGGATCTCAGAAATATTTTCTACAGATAACATACAT		
CAGGCTCATTAACCTAGGCTGATCTCTGCAAGGATCTCAGAATATTTTCTACAGATAACATACAT	• • • • • • •	
	······································	
	AGGGTAGACAGCTAGTAATTAAACTCACTTGTATGTAAAAAAAA	

(Fig. 2 continues on the next page.)

ACACTAAGAGAATAAAGGAAATGCAATAAAGTGGCCTGAAAGATAAAGGATGAGACGTGTAAAGAGACAGGGAAAGATGTGTCATTTTTTACTATGAGCAGCAGCTGAGAAGATAAAG	4430
GAATCGAGTTATCGGCAAACATGATGTTTGATCAGTGTTATTTGTTTTGAAGGCCTGCCT	4550
TTACTAGATCGCACATTCTGTAAAGGCAGGGACCATGGTATGTTGTTTATCTTTGGATTCTCAGTGATTGTCATATTTATATTTGTTGAATGAA	4670
al Glu Pro Thr Val TATCTTTCCACTCTGGTTCCAAGGAGGAGCATTCCCCAATGGTAGACGTGCTGTGTGTG	99 4784
Thr Ile Ser Pro Ser Arg Thr Glu Ala Leu Asn His His Asn Leu Leu Val Cys Ser Val Thr Asn Phe Tyr Pro Ala Gln Ile Lys Val ACC ATC TCC CCA TCC AGG ACA GAG GCC CTC AAC CAC CAC AAC CTG CTG GTC TCA GTG ACA GAT TTC FAT CCA GCC CAG ATC AAA GTC	129 4874
Arg Trp Phe Arg Asn Asp Gin Giu Giu Thr Thr Giy Val Val Ser Thr Pro Leu Ile Arg Asn Giy Asp Trp Thr Phe Gin Ile Leu Val CGG TGG TTT CGG AAT GAC CAG GAG GAG ACA ACT GGC GTT GTG TCC ACC CCC CTT ATT AGG AAC GGT GAC TGG ACC TTC CAG ATC CTG GTG	159 4964
Met Leu Glu Met Thr Pro Gln Arg Gly Asp Val Tyr Thr Cys His Val Glu His Pro Ser Leu Gln Asn Pro Ile Ile Val Glu Trp A ATG CTG GAA ATG ACT CCC CAG CGT GGA GAC GTC TAC ACC TGC CAG GAG CAC CCC AGC CTC CAG AAC CCC ATC ATC GTG GAG TGG C GT	188 5054
AAGGGGATATTGAGTTTCTGTTACTATGGGCCCCACAAGACAAAGGGCAGAGCTCCTTCTGACCCATTCCTTCC	5174
CTAGAGCACCTCTTGCTCCATGGCAAGTGCATCAGAAGAATCCTGATCTCATCACCTTTCCAGATGCTAGGGAAATTATTCTACGTACTGTTTCTCCAGATCCCAGTCCTGATAGCTCCG	5294
AGGGACTTATTATTAGGGCTGGTGACTGGGATCTTAGGGTTTAAGGTATGGATGAGTTCCTGAGGAGTG JAGATCTGCTTCCCCGCTCTCTCACCTACTCACTATACCGAAGGACCTATT	5414
GGCTGGCTTTCCCCTCCCTTAGGGGTGGTCTGAATGGAGGACXAGGTTCCTTTGACACTTTCACCTCCTGGACTGGA	5534
rg Ala Gin Ser Giu Ser Ala Gin Ser Lys Met Leu Ser Giy Ile Giy Giy Phe Val Leu Giy Leu GACAAACGCTGACACTCAGGCTCTGCTTCTTAG GG GCT CAG TCT GAA TCT GCC CAG AGC AAG ATG CTG AGT GGC ATT GGA GGC TTC GTG CTG GGG CTG	210 5632
Ile Phe Leu Gly Leu Gly Leu Ile Ile His His Arg Ser Gln Lys G ATC TTC CTC GGG CTG GGC CTT ATT ATC CAT CAC AGG AGT CAG AAA G GTGAGGAACCCCAAGGGGAAGAAGGGGAAGATGGGCTGTGACCCAGACCCTCTGTTCAGG	225 5736
GAGGTCCTGTCTCTAGATGTGGCTCTTTCCTCCTGACCCTGAGAGGAAGAAAACTGAGCTGGAGGTGGGAGGAGACAGGACAAGATTGGAGGAGGAGGAGATTGGAATCTGATTTTACTAGTTG	5856
AAAGGTAGCCCTGTCACACAGGGTGACTGATAGAGCTTATTCCAGGATATACTTACCATTCATCATCATCATTGGCTCCTTTCCAAAAGCTTCCTCCATTAAGAGGGTCAGAGCCTTGGCCT	5976
CCTTGCCTTCTAGTGACAATTTTCTTTGTTTTAGGGGATTTTTAAATTAGGGTACTTAAGGCCTTGAAGAACATGAGTGGTAAGAGAATATAACTCTAATTAAGTCACATGTGTCATTTTC	6096
CTTTGGGGTGAAAGAGTGGCTGTTTGTGTAATGAGACCTTTCTCTGCATAACTTCCTTTTGTAA <u>GACCTCAAGGGCCTCCACCAGG</u> TGATATTTCAGCCATGAGCCAGTGTGGGGGGG	6216
GCACAGGTGTAAGAGGGAAGAGCATGAGCTGAATGCACCTGACCACAATGGTCTCTGTTCATGGTATATTTGCTGCTATGAGGATCAAGACTTAGGGTCGAAGTTTGCCAGTTTCTAGGA	6336
ATCTCCAGAGGTTGTTCCCCAGAACCAAGCCTTAACTTTGGTGGTATCTTCTTGTGAAATGTGAAGCCAGAACCACAGCTTAAATGTTAGACAAGAGGATGATGCCCCACTTTGTGCCACA	6456
TOTTOGTOGCTACTGCCTGTAGGCATTTTCCAGTGACTGAAAGAGGCTGCTAGTGGTAGGGATGAGGTATCATCCCAATTTTCTAAAAAGATTGAACCCTTCATATTCCCCAGAAGAGTAA	6576
CAGCTGTTCCGCCACTTCCCACATATCTGCATCAAGCTGAAGTTCTGTGTCCTCACGAGCTGATTTCACCTTTGCACAGATCTTGCGGGAGGTGACAATAATACATTCTGGACCTCAGCT	6696
ly Leu H15 *** TTCTCTGTCTGAAGCTGCAGGGGGGCCCCTGAGGGGGGGG	229 6808
CTCCTCAGACTATTTTAACTCGGATTCGTTATCACTTTTCTGTAACCCCTCCTTGTCCCTCCC	6928
CTGTCACGCAGCCACCAGGTCATCTCCTTTCATCCCCACCTCGAGGCTGATGGCTGTGACCCTGCTTCCTGCACTTACCCAGAGCCTCTGCCTGTGCACGGCCAGCTQCGTCTACTGAGG	7048
CCCCAACCCCTTCTCTTTCTATTCTCTCCTCACACTCCTC	7168
TCAACTTCCTTAATTGAGCAGAGGCAGGAAATCACTGCAGAATGAAGGAACATACCTGAGGTGACCCAGCCAACCTGTGCCCAGAAGGAGGGTTGTACCTGAAA	7272

FIG. 2. Nucleotide sequence and translated amino acid sequence of the DC $\beta$  gene. Underlined in the 5' part of the gene are the putative "CAT," "TATA," and "cap" elements, with an arrow marking the putative cap site. Cysteine residues and the glycosylation site are within boxes. The non-functional cytoplasmic exon is underlined, as is the 3' untranslated region. The polyadenylylation signal is within a box.

mon ancestor by processes involving gene duplications.

Polymorphism of DC $\beta$  Chains Resides in the First Extracellular Domain. The  $\beta$  chain encoded by cosII-102 has a signal peptide of 32 amino acids displaying the characteristic features of archetypal signal sequences. It contains a core of hydrophobic amino acid residues, and the amino acid preceding the mature protein has no side chain (24). The mature DC $\beta$  chain contains 229 amino acid residues (9). Four cysteines, most probably forming two disulfide loops, are present as in all class II  $\beta$ -chain sequences known to date. Likewise, the single glycosylation site at asparagine-19 is common to all known  $\beta$  chains. A hydrophobic membrane-spanning segment of 21 residues is followed by a short cytoplasmic tail of 10 residues.

Comparisons between the DC $\beta$ -chain amino acid sequences deduced from cosII-102 and two cDNA clones of other haplotypes (9, 11) reveal a conspicuously uneven distribution of amino acid replacements along the polypeptides (see Fig. 3). The first extracellular domain encoded by cosII-102 displays 84% and 80% homology to the pII- $\beta$ -1 and pII- $\beta$ -2 DC $\beta$  chains, respectively, whereas the second domain shows 97% and 95%

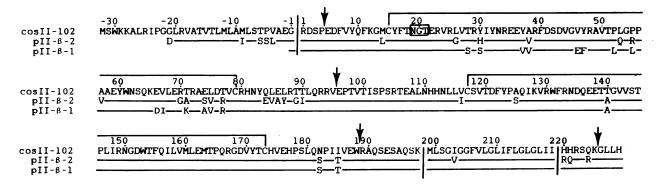


FIG. 3. Comparison of amino acid sequences deduced from the DC $\beta$  gene and two DC $\beta$  cDNA clones, pII- $\beta$ -1 (9) and pII- $\beta$ -2 (11). The standard one-letter code is used; a horizontal line indicates identity with the amino acid in the sequence above it. Arrows show exon boundaries. The boundary between the signal peptide and the mature protein and the boundaries of the membrane-spanning segment are marked with vertical bars. The glycosylation site is within the box.

homology, respectively. Thus, within one subset of class II  $\beta$ chains, the first extracellular domain carries most of the allelic polymorphisms, whereas the second domain appears virtually constant (25). In the nucleotide sequences encoding the first domains of the three DC $\beta$  chains, mutations in the first and second positions of codons are at least as prevalent as third-base mutations. This pattern of nucleotide substitutions is also found in the variable domain exons of immunoglobulins (26-28), particularly in the codons of the hypervariable regions (29). Models involving either positive selection favoring substitutions leading to amino acid replacements (26) or suppression of silent mutations (28) have been proposed. Whereas the class I antigen polymorphism may be due partly to gene conversion events (30, 31), the DC $\beta$  pattern of nucleotide substitutions is consistent with multiple independent mutational events, as has been suggested for the immunoglobulins (27). In the nucleotide stretches encoding the DC $\beta$  second domain, the pattern of substitutions is as expected for a conserved region-i.e., silent third-base substitutions in the codons are more frequent than first- and second-base substitutions. At only 4 of 34 variant positions do all three DC $\beta$  chains have unique amino acids. This low number may be due to structural constraints, allowing only a few different amino acids at a certain position.

Exon 4, encoding the connecting peptide, the membranespanning segment, and the first six amino acids of the cytoplasmic tail, is highly conserved. Its translated amino acid sequence is identical to that of the cDNA clone pII- $\beta$ -1 and differs from that deduced from pII- $\beta$ -2 at only four positions.

The 3' untranslated region of the DC $\beta$  gene in cosII-102 shows 98% and 92% homology to pII- $\beta$ -1 and pII- $\beta$ -2, respectively. The polyadenylylation signal A-T-T-A-A-A (ref. 9; see also ref. 32 and references therein) is present 330 bp downstream of the termination codon. Although the sequence determination was extended more than 100 bp 3' of the region, no additional polyadenylylation signal was found.

In all four exons available for comparison, the DC $\beta$  gene shows greater homology to pII- $\beta$ -1 than to pII- $\beta$ -2. The mature DC $\beta$ chain deduced from the gene displays an overall homology of 93% and 88%, respectively, to the amino acid sequences deduced from pII- $\beta$ -1 and pII- $\beta$ -2. The DC $\beta$  chain encoded by cosII-102 thus appears to be more closely related to that deduced from pII- $\beta$ -1 than to that of pII- $\beta$ -2.

DCB Genes Lack a Cytoplasmic Exon. The cytoplasmic tail of DC $\beta$  chains is eight amino acids shorter than that of DR $\beta$ chains (11). The nucleotide sequence of the murine I-A $\beta$  gene (33), the structural homologue of  $DC\beta$ , has provided evidence that a separate exon of 24 nucleotides accounts for this difference. A nucleotide stretch of 24 bp is indeed present in the DC $\beta$  gene, displaying clear homology to the corresponding segments of the DR $\beta$  cDNA clones and the A $\beta$  gene. The absence of this segment in the two DC $\beta$  cDNA clones whose sequences have been determined to date is probably due to an altered 5' splice junction of this exon. Instead of the preferred sequence  $A-G^{\downarrow}G(18)$ ,  $A-A^{\downarrow}G$  is found. Whether the difference in cytoplasmic tails between  $DC\beta$  and  $DR\beta$  chains has any functional significance remains to be investigated. Although this exon is not expressed in the DC $\beta$  genes, it remains well conserved at the genomic level (21 of the 24 nucleotides are identical with the A $\beta$  cytoplasmic exon). However, the entire intron between exons 4 and 5 of the DC $\beta$  gene is highly homologous to the corresponding stretch of the  $A\beta$  gene, indicating that there could be a general conservative pressure acting on the nucleotides of this part of the DC $\beta$  gene. Also, the third DC $\beta$ intron shows clear homology to the corresponding  $A\beta$  intron. Significant homology can also be detected in parts of intron 1 and intron 2 by using a computer alignment program (34).

Promoter Region. The transcriptionally important DNA sequences G-G-T-C-A-A-T-C-T (CAT) and TATA are expected at about 70 bp and 30 bp upstream from the cap site of transcribed genes, respectively (18). Because no full-length  $DC\beta$  cDNA clone is available that can define the cap site, assignment of these elements can only be tentative. A computer comparison of the 5' part of the DC $\beta$  gene with the corresponding region of the  $A\beta$  gene aligns a TATA-like sequence and a possible cap sequence with the proposed  $A\beta$  TATA and cap sequences, respectively (33). These putative TATA and cap elements, as well as a CAT-like sequence, are underlined in Fig. 2. Provided the proposed cap site is correct, the DC $\beta$  mRNA will have a 5' untranslated region of 329 nucleotides, unless an intron occurs in this region. The homology between the DC $\beta$  gene and the A $\beta$ gene in the region extending from the putative cap sites to the initiation codons is about 60%, counting insertions as mismatches.

The nucleotide stretches flanking the first domain exon are extremely rich in G+C, a feature shared with the nucleotide stretches flanking the first domain exon of the A $\beta$  gene (33) as well as the first and second domain exons of class I heavy chain genes in man (21) and mouse (22), but not of the DR $\alpha$  (19), I- $E\alpha$  (ref. 20; unpublished data), and  $\beta_2$ -microglobulin (35) genes. The significance of the G+C-rich stretches is unclear. However, it is interesting to note that they flank exons encoding polymorphic domains.

**Concluding Remarks.** The DC $\beta$ -chain gene reported here corresponds to the highly polymorphic DC $\beta$  gene detected in Southern blotting analyses of genomic DNA. The gene displays all the characteristics of a functional gene. The polymorphic region of DC $\beta$  chains is predominantly located in the first extracellular domain. The exon-intron organization of the DC $\beta$ gene is similar to that of class II antigen  $\alpha$ -chain genes, class I antigen heavy chain genes, and the constant parts of immunoglobulin genes; the similarity supports previous evidence for an evolutionary relationship among these molecules.

We thank Ms. E. Rossi for help with computer analyses, Ms. P. Ågren for technical assistance, and Dr. P. Lind and Mr. P. Wenkler for preparation of oligonucleotide primers. This work was supported by grants from the Swedish Cancer Society, King Gustav V's 80-Years Fund, and Marcus Borgström's Fund.

- 1. Benacerraf, B. (1981) Science 212, 1229-1238.
- Klareskog, L., Sandberg-Trägårdh, L., Rask, L., Lindblom, J. B., Curman, B. & Peterson, P. A. (1977) Nature (London) 265, 248– 251.
- Charron, D. J. & McDevitt, H. O. (1979) Proc. Natl. Acad. Sci. USA 76, 6567–6571.
- 4. Hämmerling, G. J. (1976) Transplant. Rev. 40, 64-82.
- Tosi, R., Tanigaki, N., Centis, D., Ferrara, G. B. & Pressman, D. (1978) J. Exp. Med. 148, 1592-1611.
- Shaw, S., Johnson, A. H. & Shearer, G. M. (1980) J. Exp. Med. 152, 565-580.
- 7. Tanigaki, N. & Tosi, R. (1982) Immunol. Rev. 66, 5-37.
- Böhme, J., Owerbach, D., Denaro, M., Lernmark, Å., Peterson, P. A. & Rask, L. (1983) Nature (London) 301, 82-84.
- 9. Larhammar, D., Schenning, L., Gustafsson, K., Wiman, K., Claesson, L., Rask, L. & Peterson, P. A. (1982) Proc. Natl. Acad. Sci. USA 79, 3687-3691.
- Larhammar, D., Gustafsson, K., Claesson, L., Bill, P., Wiman, K., Schenning, L., Sundelin, J., Widmark, E., Peterson, P. A. & Rask, L. (1982) Cell 30, 153-161.
- Larhammar, D., Andersson, G., Andersson, M., Bill, P., Böhme, J., Claesson, L., Denaro, M., Emmoth, E., Gustafsson, K., Hammerling, U., Heldin, E., Hyldig-Nielsen, J. J., Lind, P., Schenning, L., Servenius, B., Widmark, E., Rask, L. & Peterson, P. A. (1983) Hum. Immunol., in press.
- 12. Grosveld, F. G., Lund, T., Murray, E. J., Mellor, A. L., Dahl, H. H. M. & Flavell, R. A. (1982) Nucleic Acids Res. 10, 6715–6732.

Immunology: Larhammar et al.

## Proc. Natl. Acad. Sci. USA 80 (1983) 7317

- 13. Grosveld, F. G., Dahl, H.-H. M., de Boer, E. & Flavell, R. A. (1981) Gene 13, 227-237.
- 14.
- Southern, E. M. (1975) J. Mol. Biol. 98, 503–517. Rigby, P. W. J., Dieckmann, M., Rhodes, C. & Berg, P. (1977) J. Mol. Biol. 113, 237–251. 15.
- 16. Maxam, A. M. & Gilbert, W. (1980) Methods Enzymol. 65, 499-560.
- Sanger, F., Nicklen, S. & Coulson, A. R. (1977) Proc. Natl. Acad. Sci. USA 74, 5463-5467. 17.
- 18. Breathnach, R. & Chambon, P. (1981) Annu. Rev. Biochem. 50, 349-383.
- 19. Das, H. K., Lawrence, S. K. & Weissman, S. M. (1983) Proc. Natl.
- Acad. Sci. USA 80, 3543–3547. Mathis, D. J., Benoist, C. O., Williams, V. E., II, Kanter, M. R. & McDevitt, H. O. (1983) Cell 32, 745–754. Malissen, M., Malissen, B. & Jordan, B. R. (1982) Proc. Natl. Acad. 20.
- 21. Sci. USA 79, 893-897.
- 22. Steinmetz, M., Moore, K. W., Frelinger, J. G., Sher, B. T., Shen, F.-W., Boyse, E. A. & Hood, L. (1981) Cell 25, 683-692.
- Gough, N. (1981) Trends Biochem. Sci. 6, 203-205. 23.
- Kreil, G. (1981) Annu. Rev. Biochem. 50, 317-348. 24.
- 25. Kaufman, J. F. & Strominger, J. L. (1982) Nature (London) 297, 694-697.

- 26. Baltimore, D. (1981) Cell 24, 592-594.
- 27. Bothwell, A. L. M., Paskind, M., Reth, M., Imanishi-Kari, T., Rajewsky, K. & Baltimore, D. (1981) Cell 24, 625-637.
- 28. Selsing, E., Miller, J., Wilson, R. & Storb, U. (1982) Proc. Natl. Acad. Sci. USA 79, 4681-4685.
- Rechavi, G., Bienz, B., Ram, D., Ben-Neriah, Y., Cohen, J. B., Zakut, R. & Givol, D. (1982) Proc. Natl. Acad. Sci. USA 79, 4405-29. 4409
- 30. Boss, J. M., Gillam, S., Zitomer, R. S. & Smith, M. (1981) J. Biol. Chem. 256, 12958-12961.
- Pease, L. R., Schulze, D. H., Pfaffenbach, G. M. & Nathenson, 31. S. G. (1983) Proc. Natl. Acad. Sci. USA 80, 242-246.
- 32. Weiss, E., Golden, L., Zakut, R., Mellor, A., Fahrner, K., Kvist, S. & Flavell, R. A. (1983) EMBO J. 2, 453-462.
  33. Larhammar, D., Hammerling, U., Denaro, M., Lund, T., Flav-
- ell, R. A., Rask, L. & Peterson, P. A. (1983) Cell 34, 179-188.
- 34. Orcutt, B. C., Dayhoff, M. O. & Barker, W. C. (1982) ALIGN: Alignment Score Program, NBR Report 820501-08710 (National Biomedical Research Foundation, Washington, DC).
- 35. Parnes, J. R. & Seidman, J. G. (1982) Cell 29, 661-669.