

Xiaoyi Liu · Zehuan Liu · Bin Lin · Yuanyuan Liu  
Zhuoxun Chen · Weicong He · Dong Zhong · Anlong Xu

## Catalog of 162 single nucleotide polymorphisms (SNPs) in a 4.7-kb region of the HLA-DP loci in southern Chinese ethnic groups

Received: 25 August 2003 / Accepted: 28 October 2003 / Published online: 15 January 2004  
© The Japan Society of Human Genetics and Springer-Verlag 2004

**Abstract** HLA class-II proteins are cell-surface molecules that present antigens to T cells, and their expressional regulation is crucial to the immune reaction. Sequence variation at the regulatory region can directly affect the gene expression level. We cloned and sequenced a 4.7-kb region containing the regulatory region, exon1, and partial intron1 of both *HLA-DPA1* and *DPB1* genes in 25 variable sequences from southern Chinese ethnic groups and got a high-density map of 162 single nucleotide polymorphisms (SNPs): seven in 5'-flanking regions, four in 5'-untranslated regions, and four in the coding regions. By comparing these data with SNPs in dbSNP database in the NCBI, 145 SNPs (89.5%) were novel. In addition, eight genetic variations of insertion-deletion polymorphisms (INDELS) were discovered within the 4.7-kb region. These high-resolution maps can be used as resources of markers for association studies of complex diseases, assessment of individuals' predisposition to diseases, and tailoring of therapies, as well as research markers for population genetics and evolution.

**Keywords** Single nucleotide polymorphisms (SNPs) · Insertion-deletion polymorphisms (INDELS) · High-density SNP map · HLA-II genes · Regulatory region · Southern Chinese populations

### Introduction

The completion of a high-quality sequence of the human genome is a landmark event in this century, symbolizing the beginning of the postgenomic era. In this era, much interest has turned to genome variation (Collins et al. 2003), that is, to an understanding of how genomes change and take on new functional roles. Comparison of genome sequences from evolutionarily diverse species provides insight into the evolution of genes (Fu and Li 1999; Verrelli et al. 2002; Wooding et al. 2002) and a more comprehensive understanding of the function of important genomic elements. The study of sequence variation within species will also be important in defining the relationships between genotype and biological function, such as individual differences at health, susceptibility to diseases, drug response, and so on.

Besides the protein-coding sequences, a large amount of the noncoding portion of the human genome is also under active selection, suggesting that it is functionally important. It probably contains the bulk of the regulatory information controlling the expression of protein-coding genes as well as nonprotein-coding genes (Bamshad et al. 2002). It may contain sequence determinants of chromosome dynamics such as methylation and chromatin remodeling (Collins et al. 2003). Therefore, the noncoding portion of the human genome also becomes a focal point in the study of genetic variations.

Major histocompatibility complex (MHC) class-II antigens of human (HLA) are cell-surface molecules regulating a specific immune response to a pathogen by presenting antigens to T-cell receptors so as to mediate the activation of T lymphocytes. There are three isotypes of class-II molecules—DR, DQ, and DP—each consisting of two subunits, one  $\alpha$  and one  $\beta$  chain encoded by separate genes *DR* (*DQ*, *DP*) *A* and *B*, respectively. The abnormal expression of HLA-II genes causes certain diseases. For example, the expression of class-II molecules on inappropriate cells may change the ability

X. Liu · Z. Liu · B. Lin · Y. Liu · Z. Chen · W. He  
D. Zhong · A. Xu (✉)  
The Key Laboratory of Genetic Engineering of MOE,  
Department of Biochemistry,  
College of Life Sciences,  
Sun Yat-Sen (Zhongshan) University,  
510275 Guangzhou,  
P.R. China  
E-mail: ls36@zsu.edu.cn  
Tel.: +86-20-84113655  
Fax: +86-20-84038377

of antigen-presenting cells to present antigen (both foreign and self) to T lymphocytes, triggering an autoimmune disease (Laurie et al. 1992). A deficient expression of the MHC-II gene results in a hereditary immunodeficiency disease called bare lymphocyte syndrome (BLS) (Mach et al. 1996). Thus, the expressional regulation of HLA-II genes is crucial in the control of the immune response.

Four sequence motifs within promoter proximal regions of all class-II genes have been identified as *cis*-acting regulatory elements, termed W, X<sub>1</sub>, X<sub>2</sub>, and Y boxes, respectively (van den Elsen et al. 1998). These four boxes are highly conserved with respect to their sequences, relative positions, orientation, and spacing. Variation within these boxes could affect the gene expression level and the nuclear protein-binding affinities, which have been confirmed on the *DRB1*, *DRB3*, *DQA1*, *DQB1*, and *DPB1* genes (Emery et al. 1993; Morzycka-Wroblewska et al. 1997; Andersen et al. 1991; Varney et al. 1999). Therefore, it is very important to study the polymorphism of the regulatory regions of HLA-II genes, which is helpful to better understand the expressional regulation in association with the immune response in humans. Our study is the first attempt to study the variation of the regulatory region in the MHC-II genes of Chinese populations. The SNPs in the promoter region of HLA-II genes found in this study can be used as resources of markers for association studies of complex diseases, assessment of individuals' predisposition to diseases, and therapy tailoring, as well as markers for population genetics and evolution research.

*HLA-DPA1* and *HLA-DPB1* genes are classical HLA-II genes, and they are organized in head-to-head fashion with their 5' ends pointing toward each other resulting in the sequence between them functioning as promoters of both genes. Therefore, we selected an approximately 5-kb-region at these two loci containing the promoter region, the exon1, and partial intron1 of both *DPA1* and *DPB1* genes. To cover as much polymorphism as possible, sequence data were obtained from seven different ethnic populations, including both ethnic populations of southern origin and those of northern origin, in China (Yao et al. 2002). They are Jing, Lahu, Yao, Pumi, Naxi, Li, and Guangdong Han, mainly from southwest China.

## Materials and methods

Fourteen healthy and unrelated peripheral blood samples with different *HLA-DPB1* alleles (including 02012, 0202, 03011, 0401, 0402, 0501, 1401, 2201, 6201, 2801, 6301, 5101, 5601, 8001) based on our works before were collected from southwest China populations for studying the 5-kb region. Genomic DNA was extracted from whole blood containing ACD anticoagulant by the modified salting-out method, as indicated by the International Histocompatibility Work Group (IHWG) (<http://www.ihwg.org/protocols>).

Based on the contig NT\_033951 (gi: 27498326) in GenBank containing the complete sequence of the HLA region, two 28-nt primers were designed to amplify the 5-kb target fragment:

5'-AGGGCTTGAGGGCTGTATTCAAGGAGAT-3' and 5'-AGCTGGGTCTGGACTTCAAACCTGGCTC-3'. PCR amplification was performed in a 20- $\mu$ l reaction volume containing 0.75 mmol/l each dNTP, 0.25  $\mu$ mol/l each primer, 1U Extaq polymerase (Takara) and 50 ng genomic DNA. A two-step PCR program of 35 cycles in total was carried out: 95°C for 3 min; 10 cycles of 94°C for 40 s, 68°C for 4 min; and 25 cycles of 94°C for 40 s, 68°C for 4 min (increasing 5 s each cycle) followed by 72°C for 10 min at the end. The products were cloned into the pGEM-T Easy Vector (Promega, USA). Six positive plasmids for each consensus sequence were sequenced from both directions on an ABI 3700 sequencer using Bigdye reagent (Applied Biosystems, USA).

All segment sequences were assembled automatically using SeqMan in DNASTAR software package and then were carefully checked manually using the same program. All sequences were aligned with the Clustalx program (Thomson et al. 1997). Singletons and doubletons were verified by reamplifying and resequencing in both directions.

## Results and discussion

From the 14 samples, 25 cloned sequences of 4751 to 4759-bp-long fragments were obtained. There were 170 polymorphisms found, eight of which were insertions and deletions (INDELs) and all of which were shorter than 12 bp, with three INDELs in the intron 1 of the *HLA-DPB1* gene and five in the region between *HLA-DPB1* and *DPA1* genes. The detailed sequence information about the polymorphisms is listed in Table 1.

After exclusion of all INDELs, 4735 bp remained and were used to position the polymorphic sites (Table 1). Within the 4735-bp region, we identified 162 SNPs: 49 in intron1 of the *HLA-DPB1* gene, three in exon1 of the *HLA-DPB1* gene, two in the 5'-untranslated region of the *HLA-DPB1* gene, five in the 5'-flanking region of the *HLA-DPB1* gene, 83 in the region between the *HLA-DPB1* and *DPA1* genes, two in the 5'-flanking region of the *HLA-DPA1* gene, two in the 5'-untranslated region of the *HLA-DPA1* gene, one in exon1 of the *HLA-DPA1* gene, and 15 in intron 1 of the *HLA-DPA1* gene. The distribution was one SNP per 29 bp on average, and much denser than the average level of the human genome (0.1~1%). Frequencies of substitutions by types were 38.9% for A/G (63), 32.7% for C/T (53), 9.3% for A/T (15), 7.4% for C/G (12), 5.6% for A/C (9), 4.9% for G/T (8), and 0.6% for both C/T/G (1) and C/T/A (1). The ratio of transition and transversion was 2.6, being close to the 2.3.

Regarding the four cSNPs (SNP in coding regions), two were synonymous and the other two were non-synonymous. One synonymous substitution in exon1 of the *HLA-DPB1* gene was C/T at +117 and encodes alanine; another synonymous substitution was A/T at +40 in exon1 of the *HLA-DPA1* gene and encodes the same proline. Both nonsynonymous substitutions were in exon 1 of the *DPB1* gene: one was a C/T transition at the +140 position, leading to Thr/Met change at 16; another was T/C transition at +152, leading to Met/Thr change at 20. None were reported in either the dbSNP database (<http://www.ncbi.nlm.nih.gov/entrez>)

**Table 1** Characterization of variations in the 4.7-kb region in southern Chinese ethnic populations. Variation is shown by capital letter. The number of nucleotides in the coding sequences, 5'-untranslated region, and 5'-flanking regions is according to the sequence information of NT\_033951 (gi: 27498326) from NCBI. *DPA1* exon 1 and *DPA1* exon 1, 100 bp each; 5'-untranslated region of *DPA1*, 59 bp; 5'-flanking region of *DPA1*, 41 bp; 5'-flanking region of *DPB1*, 34 bp. The position of variations in the region between two genes is according to the distance to the *DPA1* gene's transcription initiation code. Frequency, the number of the major single nucleotide polymorphism (SNP) allele; the number of the minor SNP allele; e.g., the variation with 1D170 has 20 Gs and 51 Ts in our samples. *INDEL*, insertion and deletion polymorphism

Table 1 (Continued)



Table 1 (Continued)

query.fcgi?db=snp) or the IGMT/HLA sequence database (<http://www.ebi.ac.uk/imgt/hla/>) until September 2003.

In the highly conserved X<sub>1</sub>, Y, and W' box within the promoter of the *DPB1* gene (van den Elsen et al. 1998), there was one substitution per box, respectively. These three substitutions were all G/A transitions and had been reported before by Varney et al. (1999) who named the allele containing these three G/A substitutions as DP-PRO4. More interestingly, Varney et al. found this allele in seven individuals with eastern Asian origin. All these data suggest that DP-PRO4 allele containing these 3 G/A substitutions in the X<sub>1</sub>, Y, and W' box may originate from China. Their competitive binding assay (Varney et al. 1999) showed that the substitutions in the W' and X<sub>1</sub> boxes had no effect on binding affinity, while a single substitution at the site immediately adjacent to the inverted CCAAT motif in the Y box reduced binding affinity. However, whether this substitution can influence the transcription of the *DPB1* gene *in vivo* should be further studied by experiments *in vivo*, since the Y box has not the same importance as the X<sub>1</sub> box in regulating gene expression.

By comparing our data with SNPs deposited in the dbSNP database in the NCBI, we found that 145 (89.5%) of 162 SNPs were novel as of August 2003. However, three SNPs found in the dbSNP database (rs2071349, rs2856830, and rs4279481) in GenBank within this region have not been found in our 25 sequences. In short, these high-resolution genome variation maps with an unusually high density of SNPs can be used as resources of markers for association studies of complex diseases, assessment of individuals' predisposition to diseases, and therapy tailoring, as well as research markers for population genetics and evolution.

**Acknowledgements** This research was supported by projects (Nos. 30178073 and 30100275) and the key project (No. 69935020) of the National Natural Science Foundation of China, the key project (No. 021691) of the Guangdong Natural Science Foundation, and the projects (No. 2001AA224021-04) of the State High-Tech Development Project of the Ministry of Science & Technology, Guangdong Provincial Department of Science and Technology, and the Bureau of Science and Technology of Guangzhou City.

## References

- Andersen LC, Beaty JS, Nettles JW, Seyfried CE, Nepom GT, Nepom BS (1991) Allelic polymorphism in transcriptional regulatory regions of *HLA-DQB1* genes. *J Exp Med* 173:181–192
- Bamshad MJ, Mummidi S, Gonzalez E, Ahuja SS, Dunn DM, Watkins WS, Wooding S, Stone AC, Jorde LB, Weiss RB, Ahuja SK (2002) A strong signature of balancing selection in the 5'cis-regulatory region of *CCR5*. *Proc Natl Acad Sci USA* 99:10539–10544
- Collins FS, Green ED, Guttmacher AE, Guyer MS (2003) A vision for the future of genomics research: A blueprint for the genomic era. *Nature* 422:835–847
- Emery P, Mach B, Reith W (1993) The different level of expression of *HLA-DRB1* and -*DRB3* genes is controlled by conserved isotypic differences in promoter sequence. *Hum Immunol* 38:137–147
- Fu YX, Li WH (1999) Coalescing into the 21st century: an overview and prospects of coalescent theory. *Theor Popul Biol* 56:1–10
- Laurie HG, Catherine JK (1992) Sequences and factors: a guide to MHC class-II transcription. *Annu Rev Immunol* 10:13–49
- Mach B, Steimle V, Martinez-Soria E, Reith W (1996) Regulation of MHC class II genes: lessons from a disease. *Annu Rev Immunol* 14:301–331
- Morzycka-Wroblewska E, Munshi A, Ostermayer M, Harwood JI, Kagnoff MF (1997) Differential expression of *HLA-DQA1* alleles associated with promoter polymorphism. *Immunogenetics* 45:163–170
- Thomson JD, Gibson TJ, Plewniak F, Jeanmougin F, Higgins DG. (1997) The CLUSTAL\_X windows interface: flexible strategies for multiple sequences alignment aided by quality analysis tools. *Nucleic Acids Res* 25:4876–4882
- Varney MD, Gavrilidis A, Tait BD (1999) Polymorphism in the regulatory regions of the *HLA-DPB1* gene. *Hum Immunol* 60:955–961
- Van den Elsen PJ, Peijnenburg A, Van Eggermond MCJA, Gobin SJP (1998) Shared regulatory elements in the promoters of MHC class I and II genes. *Immunol Today* 19:309–312
- Verrelli BC, McDonald JH, Argyropoulos G, Destro-Bisol G, Froment A, Drousiotou A, Lefranc G, Helal AN, Loiselet J, Tishkoff SA (2002) Evidence for balancing selection from nucleotide sequence analyses of human *G6PD*. *Am J Hum Genet* 71:1112–1128
- Wooding SP, Watkins WS, Bamshad MJ, Dunn DM, Weiss RB, Jorde LB (2002) DNA sequence variation in a 3.7-kb noncoding sequence 5' of the *CYP1A2* gene: implications for human population history and natural selection. *Am J Hum Genet* 71:528–542
- Yao YG, Nie L, Harpending H, Fu YX, Yuan ZG, Zhang YP (2002) Genetic relationship of Chinese ethnic populations revealed by mtDNA sequence diversity. *Am J Phys Anthropol* 118:63–76