



## Short-tandem repeat analysis in seven Chinese regional populations

Xing-bo Song<sup>1</sup>, Yi Zhou<sup>1</sup>, Bin-wu Ying<sup>1</sup>, Lan-lan Wang<sup>1</sup>, Yi-song Li<sup>1</sup>, Jian-feng Liu<sup>2</sup>, Xiao-gang Bai<sup>3</sup>, Lei Zhang<sup>1</sup>, Xiao-jun Lu<sup>1</sup>, Jun Wang<sup>1</sup> and Yuan-xin Ye<sup>1</sup>

<sup>1</sup>Department of Laboratory Medicine, West China Hospital, Sichuan University, Chengdu, Sichuan, P.R. China.

<sup>2</sup>The Police Station of Wenzhou, Wenzhou, Zhejiang, P.R. China.

<sup>3</sup>The Police Station of Chengdu, Chengdu, Sichuan, P.R. China.

### Abstract

In the present study, we investigated the application of 13 short tandem repeat (STR) loci (D13S317, D7S820, TH01, D16S539, CSFIPO, VWA, D8S1179, TPOX, FGA, D3S1358, D21S11, D18S51 and D5S818) routinely used in forensic analysis, for delineating population relationships among seven human populations representing the two major geographic groups, namely the southern and northern Chinese. The resulting single topology revealed pronounced geographic and population partitioning, consistent with the differences in geographic location, languages and eating habits. These findings suggest that forensic STR loci might be particularly powerful tools in providing the necessary fine resolution for reconstructing recent human evolutionary history.

*Key words:* forensic medicine, population genetics, short-tandem repeat, human evolutionary history, genetic distance.

Received: October 19, 2009; Accepted: June 11, 2010.

### Introduction

The present Chinese population of around 1.4 billion is primarily divided by the Yellow River into two large groups, the southern and the northern, with diverse languages and eating habits. There is thus an immense scope to study the processes of anthropological subdivisions and microevolutionary effects in different populations groups of China. However, the traditional structure of Chinese populations is facing the imminent threat of disintegration through urbanization and increasing communication, with the consequential gene flow between subcastes through marriages. Therefore, there is a need for understanding local traditional population structure and its role in shaping human genome diversity.

A large-scale survey of autosomal variation in an ample geographic sample of human Asian populations has shown that, apart from geography, genetic ancestry is strongly correlated with linguistic affiliations (The HUGO Pan-Asian SNP Consortium 2009). A distinction between northern and southern Chinese populations (Han and minority alike) has been observed on analyzing genetic markers (Zhao and Lee, 1989; Chu *et al.*, 1998). Short tandem repeat (STR) loci are highly polymorphic loci in the human genome, are relatively small in size, and can be analyzed in a multiplex PCR fashion. Many population genetic studies

have investigated the polymorphism profile of the STR system in Chinese Han populations, this including the loci D13S317, D7S820, TH01, D16S539, CSFIPO, VWA, D8S1179, TPOX, FGA, D3S1358, D21S11, D18S51 and D5S818 (Cai *et al.*, 2005; Deng *et al.*, 2007). In the present study, these 13 STR loci in seven Chinese regional populations, comprising 3 northern (Henan, Beijing and Tianjin) and 4 southern (Sichuan, Fujian Guangdong, and Zhejiang), were analyzed by way of capillary electrophoresis on 3100 genetic analyzers.

Based on the population data of these STR polymorphisms, the forensic parameters of the respective loci were calculated in order to estimate their value in genetic identity testing. Furthermore, genomic affinities among the diverse regional population groups were evaluated. The current study contributed to supplementing the ever-increasing population-information database worldwide.

### Materials and Methods

#### Sample preparation

Whole blood was obtained by venipuncture in EDTA-coated vacutainers from unrelated, consenting donors. Community history and family disease backgrounds were recorded on blood donor cards.

Seven geographically targeted populations, encompassing the major biogeographical zones and representing the two main Han populations (southern and northern),

were selected. These included 4 southern, the Sichuan ( $n = 260$ , Ying *et al.*, 2005), Fujian ( $n = 150$ ), Guangdong ( $n = 522$ ) and Zhejiang ( $n = 147$ ), and 3 northern, the Henan ( $n = 101$ ), Tianjin ( $n = 150$ ) and Beijing ( $n = 216$ ). Their respective location is shown in Figure 1

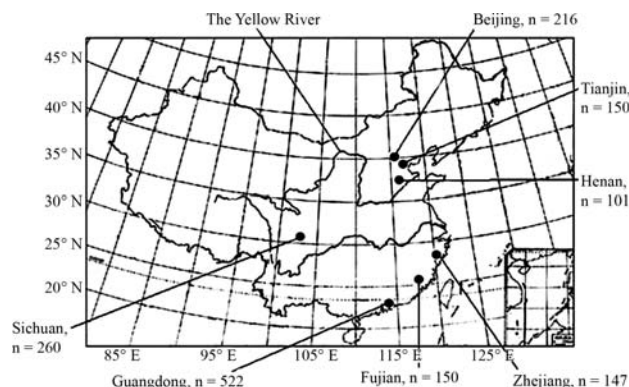
DNA was extracted using the Chelex method (Walsh *et al.*, 1991).

### PCR amplification

PCR amplification was carried out on a thermal cycler, using primers with the same sequences as those in the "PowerPlex 16 System" kit (Krenke *et al.*, 2002). Each PCR reaction was performed with 2.5  $\mu\text{L}$  of template DNA (5-250 ng), 0.5  $\mu\text{M}$  of each primer, 2.5  $\mu\text{L}$  of Taq buffer (10PCR Buffer, Applied Biosystems), 2  $\mu\text{L}$  of  $\text{MgCl}_2$  (25  $\mu\text{M}$ , Applied Biosystems), 0.5  $\mu\text{L}$  of a dNTPs mix (10  $\mu\text{M}$  PCR nucleotide Mix, Promega), and 1U Taq polymerase (DyNAzyme, DNA Polymerase, Finnzymes) in a total volume of 25  $\mu\text{L}$ . A total of 30 cycles were run, with an initial incubation (preliminary denaturation) step at 96 $^\circ\text{C}$  for 2 min, followed by 10 cycles of 94 $^\circ\text{C}$  for 1 min, 60 $^\circ\text{C}$  for 1 min and 70 $^\circ\text{C}$  for 1.5 min, followed by 20 cycles of 90 $^\circ\text{C}$  for 1 min, 60 $^\circ\text{C}$  for 1 min and 70 $^\circ\text{C}$  for 1.5 min, ending with a final extension at 60 $^\circ\text{C}$  for 30 min.

### Electrophoresis and analysis

The PCR product (1.5  $\mu\text{L}$ ), as well as GeneScan-400HD-ROX Size Standard (Applied Biosystems) (0.5  $\mu\text{L}$ ), were added to 24.5  $\mu\text{L}$  of deionized formamide, and subsequently denatured for 3 min at 95 $^\circ\text{C}$ . Alleles were then separated by capillary electrophoresis in POP-4 polymer (Applied Biosystems) with the GS STR POP4 D Module (1 mL), using an ABI PRISM 3100 Genetic Analyzer (Applied Biosystems). Samples were injected into the capillaries in batches of 16 samples, directly from the microtitre plate, for 10 s at 3 kV. Electrophoresis was performed at 15 kV and 60 $^\circ\text{C}$  for 45 min under routine running conditions. Alleles were identified by means of GeneScan Analysis 3.7 Software (Applied Biosystems), whereupon the analyzed data were automatically genotyped using



**Figure 1** - Geographical location of the seven populations in China.

Genotyper 3.6 Software (Applied Biosystems) and a template specially made for this specific multiplex system. The Peak Amplitude Threshold adopted was more than 150 RFU (relative fluorescence units).

### Statistical analysis

Individual locus frequency was calculated from the number of each genotype in the sample set. Unbiased estimates of expected heterozygosity were computed as described by Edwards *et al.* (1992). Possible divergence from Hardy-Weinberg equilibrium (HWE) was determined by calculating an unbiased estimate of expected homozygote/heterozygote frequencies (Nei and Roychoudhury, 1974; Chakraborty *et al.*, 1988; ), through likelihood-ratio testing (Weir, 1992; Buscemi *et al.*, 1995). The Chi-square test was applied for comparing the genotype and allelic frequency of each STR locus among the studied populations. We also calculated certain parameters of genetic and forensic interest, *i.e.*, the power of discrimination (Grunbaum *et al.*, 1978), the chance of exclusion (Ohno *et al.*, 1982), polymorphism information content (PIC) (Botstein *et al.*, 1980) and heterozygosity. Distance was estimated using the Nei formula (Nei and Roychoudhury, 1972; Li and Nei, 1977), whereas phylogeny was inferred by UPGMA and Neighbor-Joining methods in Mega 2.1.

## Results

### Polymorphisms of 13 STR loci in seven Chinese Han populations

Details on polymorphism exhibited at the 13 loci with respect to the allele frequencies in the seven Chinese populations are listed in Tables S1-S13.

Despite the wide range of allelic variation in the 13 STR loci, a discernable pattern depicting mutual geographical affiliation is apparent. Generally speaking, frequency was high in only few alleles (*e.g.*, allele 9 of TH01, allele 14 of VWA, allele 14 of D16S539, allele 30 of D21S11, and allele 10 of TPOX) (Tables S1-S13). 13 STR loci among seven Chinese populations showed similar trends. Furthermore, both genotype and allele distribution were not significantly different among the seven Chinese populations ( $p > 0.05$ ). These results are thought to reflect the influence of gene flow due to geographic proximity.

### Phenotype distribution and value in forensic application

The distribution of observed allele frequencies in the 13 loci (D13S317, D7S820, TH01, D16S539, CSFIPO, VWA, D8S1179, TPOX, FGA, D3S1358, D21S11, D18S51, D5S818), as well as the results from the various analytical procedures for testing the correspondence of genotype frequencies with Hardy-Weinberg equilibrium, are shown in Tables S1-S13.

All the 13 loci complied with Hardy-Weinberg equilibrium, with no evidence of association of alleles among the 13 loci. The parameters for both forensic efficiency and genetic variability, such as MP, PD, PIC, PE and heterozygosity, were calculated and subsequently listed for each population in the supplementary tables.

### Analysis of genetic distances

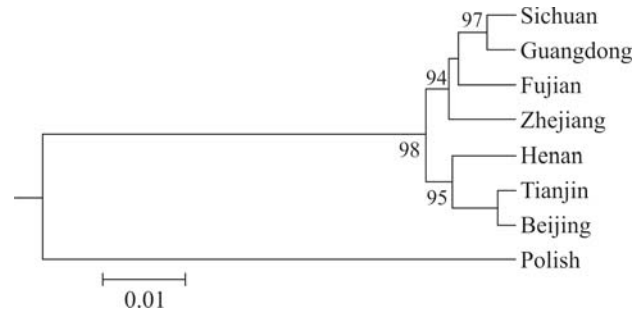
In order to ascertain relationships among the seven Chinese populations, we have calculated the Nei measure of pairwise genetic distances using allele frequency data from the 13 STR markers. Polish population data (Pepinski *et al.*, 2005) was included in the analysis as outgroup reference.

The longest distance (0.0320) was noted between the Fujian (a southern) and Henan (a northern) populations, whereas the lowest (0.0041) was observed between Beijing (a northern) and Tianjin (also a northern) populations (Table 1).

Based on genetic distance data, population trees were constructed using the UPGMA and Neighbor-Joining methods. As both methods revealed the same pattern, UPGMA results were preferred for display. Bootstrap values for the trees were high (Figure 2). The Sichuan (southern) and the Guangdong (also southern) populations first clustered together with a high bootstrap value (97%), to then cluster with the other two southern populations, the Zhejiang and Fujian, with bootstrap values of 94%. The three northern populations (Beijing, Tianjin and Henan) formed a single cluster with bootstrap values of 95%. The two major populations (the northern and southern) clustered together with bootstrap values of 98%. As expected, on comparing the Polish population, as outgroup control, with any pair of the Chinese populations, the distance was greater.

### Discussion

Owing to the several advantages, such as high polymorphism, ease and low-cost, STR markers have been widely used for fine-scale genetic mapping (Edwards *et al.*, 1991, 1992; Hearne *et al.*, 1992), intra-species phylogen-



**Figure 2** - Genetic affinities between seven Chinese populations based on 13 STR loci by DA distance and UPGMA clustering methods.

etic reconstruction (Bowcock *et al.*, 1994; Jorde *et al.*, 1998), maternity/paternity determination (Hammond *et al.*, 1994), and forensic analysis (Edwards *et al.*, 1991; Hearne *et al.*, 1992). Consistent with previous studies (Cai *et al.*, 2005; Deng *et al.*, 2007; Ying *et al.*, 2005, 2006), all the 13 STR loci were highly polymorphic in the seven population samples and exhibited desirable values in the forensic analysis and genetic analysis.

Over the past decades, and based on STR polymorphisms, important information has contributed to elucidating the history of human populations (Jorde *et al.*, 1997; Shriver *et al.*, 1997), as well as genetic microdifferentiation among local subdivided populations (Reddy *et al.*, 2001). In the current study, seven Chinese Han populations, with three representative groups from the northern portion and four from the southern, were investigated, by comparing the allele frequency of 13 STR loci, whereby the following consequential information was obtained. First, the 13 loci exhibited high polymorphism in all the seven populations, but with no significant difference in allele distribution in any. It was inferred that both geographical and ethnic affiliations in Chinese Han populations are close. A single STR-based comparison of the population was insufficient to detect the delicate mutual difference among these populations. A method integrating polymorphic information on all the 13 STR loci of each population is essential for determining respective genetic distances. In addition, the specific parameters revealed the high forensic efficiency of the 13 STR loci. Heterozygosity among these ranged from

**Table 1** - Genetic distances of 8 populations using UPGMA software.

Population	Sichuan	Fujian	Guangdong	Tianjin	Zhejiang	Beijing	Henan	Polish
Sichuan								
Fujian	0.0132							
Guangdong	0.0071	0.0146						
Tianjin	0.0188	0.0213	0.0119					
Zhejiang	0.0142	0.0195	0.0149	0.0230				
Beijing	0.0166	0.0193	0.0121	0.0041	0.0210			
Henan	0.0285	0.0320	0.0241	0.0153	0.0318	0.0155		
Polish	0.1202	0.1255	0.1141	0.1099	0.1254	0.1066	0.0980	

0.5248 (TPOX in the Henan population) to 0.8989 (D8S1179 in the Zhejiang), whereas the number of alleles observed ranged from 8 (TPOX) to 20 (D18S51). The data presented herein will facilitate calculating matching probabilities in forensic casework, in the event of Chinese individuals being considered as the source of DNA evidence. Furthermore, by using the UPGMA and Neighbor-Joining methods, it was possible to calculate genetic distances on the basis of data from all the 13 STR locus polymorphisms in each population, whereby a population tree was constructed to reflect mutual evolutionary relationships. The results indicated that genetic distances among these populations correspond to their geographic location, Whereas three northern populations formed one cluster, the four southern ones formed another cluster, as confirmed through UPGMA and Neighbor-Joining methodology. Although the distances among the studied populations were only short, clustering remained distinct in certain groups, this being consistent with their ethnohistory and geographic location. Compared to the outgroup control (Polish population), Chinese southern and northern populations clustered together. While clustering tended to occur between two populations with smallest geographic distance, it was notable that the Guangdong population first clustered with that of Sichuan, instead of doing so with the two geographically nearer populations of Fujian and Zhejiang, thereby providing evidence for historical records that the earliest Sichuan population most likely emigrated from Guangdong.

## Acknowledgments

We thank Dr. Junping Xin (Loyola University Medical Center) for critical review and editorial assistance during manuscript revision. This study was supported by Grants #30900658 from the National Natural Science Foundation of China.

## References

- Botstein D, White RL, Scolnick M and Davis RW (1980) Construction of a genetic linkage map in man using restriction fragment length polymorphisms. *Am J Hum Genet* 32:314-331.
- Bowcock AM, Ruiz-Linares A, Tomfohrde J, Minch E, Kidd JR and Cavalli-Sforza LL (1994) High resolution of human evolutionary trees with polymorphic microsatellites. *Nature* 368:455-457.
- Buscemi L, Cucurachi N, Mencarelli R, Tagliabracci A, Wiegand P and Ferrara SD (1995) PCR analysis of the short tandem repeat (STR) system HUMVWA31. Allele and genotype frequencies in an Italian population sample. *Int J Legal Med* 107:171-173.
- Cai GQ, Chen LX, Tong DY, Ou JH and Wu XY (2005) Mutations of 15 short tandem repeat loci in Chinese population. *Zhonghua Yi Xue Yi Chuan Xue Za Zhi* 22:507-509.
- Chakraborty R, Smouse PE and Neel JV (1988) Population amalgamation and genetic variation: Observations on artificially agglomerated tribal populations of Central and South America. *Am J Hum Genet* 43:709-725.
- Chu JY, Huang W, Kuang SQ, Wang JM, Xu JJ, Chu ZT, Yang ZQ, Lin KQ, Li P, Wu M, *et al.* (1998) Genetic relationship of populations in China. *Proc Natl Acad Sci USA* 95:11763-11768.
- Deng YJ, Yan JW, Yu XG, Li YZ, Mu HF, Huang YQ, Shi XT and Sun WM (2007) Genetic analysis of 15 STR loci in Chinese Han population from West China. *Genomics Proteomics Bioinform* 5:66-69.
- Edwards A, Civitello A, Hammond HA and Caskey CT (1991) DNA typing and genetic mapping with trimeric and tetrameric tandem repeats. *Am J Hum Genet* 49:746-756.
- Edwards A, Hammond HA, Jin L, Caskey CT and Chakraborty R (1992) Genetic variation at five trimeric and tetrameric tandem repeat loci in four human population groups. *Genomics* 12:241-253.
- Grunbaum BW, Selvin S, Pace N and Black DM (1978) Frequency distribution and discrimination probability of twelve protein genetic variants in human blood as functions of race, sex, and age. *J Forensic Sci* 23:577-587.
- Hammond HA, Jin L, Zhong Y, Caskey CT and Chakraborty R (1994) Evaluation of 13 short tandem repeat loci for use in personal identification applications. *Am J Hum Genet* 55:175-189.
- Hearne CM, Ghosh S and Todd JA (1992) Microsatellites for linkage analysis of genetic traits. *Trends Genet* 8:288-294.
- Jorde LB, Bamshad M and Rogers AR (1998) Using mitochondrial and nuclear DNA markers to reconstruct human evolution. *Bioessays* 20:126-136.
- Jorde LB, Rogers AR, Bamshad M, Scott WW, Krakowiak P, Sung S, Kere J and Harpending HC (1997) Microsatellite diversity and the demographic history of modern humans. *Proc Natl Acad Sci USA* 94:3100-3103.
- Krenke BE, Tereba A, Anderson SJ, Buel E, Culhane S, Finis CJ, Tomsey CS, Zachetti JM, Masibay A, Rabbach DR *et al.* (2002) Validation of a 16-locus fluorescent multiplex system. *J Forensic Sci* 47:773-785.
- Li WH and Nei M (1977) Persistence of common alleles in two related populations or species. *Genetics* 86:901-914.
- Nei M and Roychoudhury AK (1972) Gene differences between Caucasian, Negro, and Japanese populations. *Science* 177:434-436.
- Nei M and Roychoudhury AK (1974) Sampling variances of heterozygosity and genetic distance. *Genetics* 76:379-390.
- Ohno Y, Sebetan IM and Akaishi S (1982) A simple method for calculating the probability of excluding paternity with any number of codominant alleles. *Forensic Sci Int* 19:93-98.
- Pepinski W, Niemcunowicz-Janica A, Skawronska M, Janica J, Koc-Zorawska E, Aleksandrowicz-Bukin M and Soltyszewski I (2005) Genetic data on 15 STR loci in the ethnic group of Polish Tatars residing in the area of Podlasie (Northeastern Poland). *Forensic Sci Int* 49:263-265.
- Reddy BM, Pfeffer A, Crawford MH and Langstieh BT (2001) Population substructure and patterns of quantitative variation among the Gollas of southern Andhra Pradesh, India. *Hum Biol* 73:291-306.
- Shriver MD, Jin L, Ferrell RE and Deka R (1997) Microsatellite data support an early population expansion in Africa. *Genome Res* 7:586-591.

The HUGO Pan-Asian SNP Consortium (2009) Mapping human genetic diversity in Asia. *Science* 326:1541-1545.

Walsh PS, Metzger DA and Higuchi R (1991) Chelex 100 as a medium for simple extraction of DNA for PCR-based typing from forensic material. *Biotechniques* 10:506-513.

Weir BS (1992) Independence of VNTR alleles defined as fixed bins. *Genetics* 130:873-887.

Ying BW, Wei YG, Sun XM, Liu TT and Hou YP (2005) STR data for the AmpFISTR profiler plus from western China. *J Forensic Sci* 50:716-717.

Ying BW, Fan H, Liu TT, Zhao ZH, Liang ZH, Feng S, Yuan WA and Yun LB (2006) Genetic variation for five short tandem repeat loci in a Central China population sample. *J Forensic Sci* 51:1201.

Zhao TM and Lee TD (1989) Gm and Km allotypes in 74 Chinese populations: A hypothesis of the origin of the Chinese nation. *Hum Genet* 83:101-110.

## Supplementary Material

The following online material is available for this article:

Table S1: Genetic polymorphism at the D3S1358 locus for the seven Chinese population groups.

Table S2: Genetic polymorphism at the D16S539 locus for the seven Chinese population groups.

Table S3: Genetic polymorphism at the TPOX locus for the seven Chinese population groups.

Table S4: Genetic polymorphism at the TH01 locus for the seven Chinese population groups.

Table S5: Genetic polymorphism at the CSF1PO locus for the seven Chinese population groups.

Table S6: Genetic polymorphism at the D7S820 locus for the seven Chinese population groups.

Table S7: Genetic polymorphism at the VWA locus for the seven Chinese population groups.

Table S8: Genetic polymorphism at the FGA locus for the seven Chinese population groups.

Table S9: Genetic polymorphism at the D8S1179 locus for the seven Chinese population groups.

Table S10: Genetic polymorphism at the D21S11 locus for the seven Chinese population groups.

Table S11: Genetic polymorphism at the D18S51 locus for the seven Chinese population groups.

Table S12: Genetic polymorphism at the D5S818 locus for the seven Chinese population groups.

Table S13: Genetic polymorphism at the D13S317 locus for the seven Chinese population groups.

This material is available as part of the online article from <http://www.scielo.br/gmb>.

*Associate Editor: Francisco Mauro Salzano*

License information: This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.