

An INDEL polymorphism at the X-STR GATA172D05 flanking region

Elzemar Martins Ribeiro Rodrigues ·
Ney Pereira Carneiro dos Santos ·
Ândrea Kely Campos Ribeiro dos Santos ·
Anderson Nonato Marinho · Marco Antonio Zago ·
Iva Gomes · António Amorim · Leonor Gusmão ·
Sidney Emanuel Batista dos Santos

Received: 14 May 2008 / Accepted: 10 November 2008 / Published online: 2 December 2008
© Springer-Verlag 2008

Abstract A new polymorphic INDEL was detected at the X-STR GATA172D05 flanking region, which corresponds to an 18-bp deletion, 141 bp upstream the TAGA repeat motif. This INDEL was found to be polymorphic in different population samples from Native Americans, Africans, and Europeans as well as in an admixed population from the Amazonia (Belém). Gene diversities varied between 37.5% in Native Americans and 49.9% in Africans. Comparison between human and chimpanzee sequences showed that the ancestral state corresponds to the presence of two copies of 18 bp, detected in both species; and the mutated allele has lost one of these two copies. The simultaneous analysis of the short tandem repeat (STR) and INDEL variation showed an association between the INDEL ancestral allele with the shorter STR

alleles. High diversities were found in all population groups when combining the information provided by the INDEL and STR variation. Gene diversities varied between 76.7% in Native Americans and 80.6% in both Portugal and Belém.

Keywords X chromosome · STRs · Indel polymorphism · Amazonia · GATA172D05

Introduction

Interest in new X chromosome genetic markers by different research groups, mainly short tandem repeats (STRs), has grown recently fostering the application in anthropological

E. M. R. Rodrigues · Â. K. C. R. dos Santos · A. N. Marinho ·
S. E. B. dos Santos
Laboratório de Genética Humana e Médica,
Departamento de Patologia, Universidade Federal do Pará,
Belém, PA, Brazil

N. P. C. dos Santos
Faculdade de Ciências Biológicas,
Campus Universitário de Altamira, Universidade Federal do Pará,
Altamira, PA, Brazil

M. A. Zago
Departamento de Clínica Médica, Departamento de Genética,
Faculdade de Medicina, Universidade de São Paulo,
Ribeirão Preto, SP, Brazil

I. Gomes · A. Amorim · L. Gusmão
IPATIMUP, Institute of Molecular Pathology and Immunology of
the University of Porto,
4200-465 Porto, Portugal

I. Gomes
Institute of Legal Medicine,
University of Santiago de Compostela,
15782 Santiago de Compostela, Spain

A. Amorim
Faculty of Sciences, University of Porto,
4050 Porto, Portugal

S. E. B. dos Santos (✉)
Laboratório de Genética Humana e Médica,
Instituto de Ciências Biológicas, Universidade Federal do Pará,
Cidade Universitária Prof. José da Silva Netto. Av. Augusto
Corrêa, 01,
Belém, PA, Brazil
e-mail: sidneysantos@ufpa.br

and forensic genetic studies [1–6]. Prior to the use of new markers in forensic or kinship testing, it is important to collect population data and to construct reference databases to document the genetic variation of these specific markers among worldwide populations. It is also important to know more about sequence variation within or around DNA sequences of repeat structures such as single nucleotide or insertion–deletion polymorphisms (SNPs/INDELs). Polymorphic sites in these regions may induce undetectable technical problems that alter genotyping, causing for example allele dropout or null alleles [7–13]. These events occur when a base pair change takes place in the primer binding region of the DNA target segment, causing failure of primer annealing and consequently failure of amplification and allele detection [10–13]. The expected amplified segment size can also change significantly in cases of insertion or deletion in short DNA fragments located in the amplified region. It is important to study the flanking regions of DNA repetitive sequences, using different primer sets, in order to prevent misinterpretation of results. INDEL polymorphisms have received less attention compared to SNPs in forensic and in population variation studies, despite accounting for approximately 16–25% of all sequence polymorphisms in the genome [14]. However, like SNPs, INDELs can increase the power of discrimination of STRs in particular cases of identity testing, e.g., in cases of degraded DNA and when only small amplicons are amplified.

As a result of a study on X-STRs variability in an Amazonian population from Brazil, we identified an 18-bp INDEL polymorphism in the upstream flanking region of GATA172D05 locus which is widely used in forensic and population genetic research [15–18].

This work focused on: (1) the molecular characterization and exact location of the INDEL in relation to the repeat structure of the GATA172D05 marker and (2) the haplotype frequencies and gene diversities observed for both polymorphisms (STR and INDEL) in different human groups, namely Africans, Europeans, and Amerindians.

Materials and methods

Samples and DNA extraction

The study was carried out in accordance with the Declaration of Helsinki (2000) of the World Medical Association. Blood samples were collected from healthy, unrelated individuals from different populations, under informed consent. Samples studied consisted of 177 individuals (83 males, 94 females) from northern Portugal [19], 224 individuals (109 males, 115 females) from Brazilian Amazonian tribes, 136 individuals (89 males, 47

females) of African origin [20], and 188 samples (106 male, 82 female) from Belém (Brazil), characterized by an inter-ethnic mixture of Europeans, Africans, and indigenous people [21]. DNA was extracted using the phenol–chloroform protocol [22].

Primers, STR, and INDEL amplification

Primers were designed using PRIMER3 (<http://www.genome.wi.mit.edu/cgi-bin/primer/primer3>) software (Fig. 1) and screened for hairpins and primer dimers with AUTO DIMER CHECK (<http://www.cstl.nist.gov/biotech/strbase/AutoDimerHomepage>). Amplification of both polymorphic regions was carried out using three primers, as shown in Fig. 1: a common reverse primer, marked with fluorochrome 6-FAM, and two forward primers, F1 for DNA segments between 294 and 340 bp and F2 for amplicons ranging between 109 and 137 bp.

PCR was performed in a final volume of 12.5 µl reaction mix containing 10–20 ng of genomic DNA, 1× PCR buffer with 0.75 mM MgCl₂, 125 µM of each dNTP, 1 U AmpliTaq Platinum DNA Polymerase (Invitrogen Life Technologies, Carlsbad, CA, USA), and 0.2 µM of each primer. PCR thermocycling conditions were: 1 min at 94°C, 1 min at 60°C, and 2 min at 70°C for ten cycles, followed by 1 min at 90°C, 1 min at 60°C, and 2 min at 70°C for 17 cycles, and a final extension of 60 min at 60°C.

Detection, typing, and analysis of PCR products

Aliquots of 1 µl of amplified PCR product was prepared for capillary electrophoresis by adding 8.5 µl of HiDi formamide (Applied Biosystems, Foster City, CA, USA) and 0.5 µl of internal size standard 500 ROX (Applied Biosystems). Samples were run in an ABI PRISM 3130 Genetic Analyzer (Applied Biosystems) and results analyzed with GeneMapper software v.3.2.2.

```

F1
ACACACTGTAGCCTGACCAATCAGATAGTAGTCAGATAGTGGTCCGACTT
      INDEL
(TATTCTTTTGTATGTACC)1-2GTAGAGACAGAAAGTGGATTAGTGGTTA
CCAGGGACTGGAGGAAGGAGGGAATGGGGCTTACTGTTTAAACCAGTACA
      F2
AAGTTTATTTTGGGAACATAAAAGTTGTGAAGATGGATAGTGGTGATGG
TTGCACAGATATA (TAGA)nGCTATATCAATACCTATATCTATAGATATA
GATCTTTTTGAATCCGGGCTTCAATTATTT
      R

```

Fig. 1 GATA172D05 sequence structure with the INDEL polymorphism represented in *bold*. F1 initial forward primer, F2 new designed forward primer, R reverse primer. In *brackets*, the TAGA repeat motif of the STR

Allele sequencing

Allele sequencing of purified STR fragments of DNA male subjects was performed using the BigDye Terminator v3.1 Cycle Sequencing Kit (Applied Biosystems). Sequencing products were visualized in an ABI PRISM 3130 Genetic Analyzer (Applied Biosystems). DNA samples from at least two carriers of each allele size were sequenced. Alleles confirmed by direct sequencing were used to construct an allelic ladder, subsequently employed in genotyping.

Statistical analysis

Exact tests for Hardy–Weinberg equilibrium among females and exact tests of population differentiation between allele frequencies in males and females were performed using Arlequin software 3.1 [23]. Human and chimpanzee sequences were compared using the BLAT search genome tool (<http://genome.brc.mcw.edu>).

Relative age of mutation was estimated by using a Bayesian approach incorporated in the program Batwing [24], considering a model of exponential growth from initial constant population size, using the effective population size and the population growth rate priors specified in Weale et al. [25]. Mutation rate was set at 0.0014, as an average estimated from data at the American Association of Blood Banks 2003 Annual Report (<http://www.cstl.nist.gov/biotech/strbase/mutation.htm>). Generation time was set at 25 years.

The power of discrimination of the male sample (PD_M) and the female sample (PD_F) and mean exclusion chance for X-STR in standard trios (MEC_T) and father/daughter duos (MEC_D) lacking maternal genotype information (motherless) were calculated according to Desmarais et al. [26].

Results and discussion

INDEL identification

When studying the genetic variation of GATA172D05 in a Brazilian Amazon population, using primers F1 and R (see Fig. 1), three new large intermediate alleles not previously described in any other population were found with significantly high frequencies. These alleles were characterized by the presence of more than 13 repeats, and apparently containing two extra pairs of nucleotides (13.2–15.2). Sequencing of at least two samples of each new allele was performed and results were compared with those obtained for the three most frequent alleles (alleles 9–11). The results revealed that these new alleles do not correspond to a variation inside the repeat motif but to an

18-bp insertion/deletion polymorphism, located 141 bases upstream to the beginning of the X-STR repeat structure (the insertion sequence is shown in bold type in Fig. 1). Therefore, the presence of this duplication induced alterations to the correct genotyping of the STR allele, increasing in 4.2 units the real number of repeats present in the chromosome.

The new set of primers used in this work (with the forward primer 190 bp away from the upstream flanking region of the repeat; F1 in Fig. 1) allowed the identification of this previously undescribed variation. Sequence data were submitted to GenBank (<http://www.ncbi.nlm.nih.gov>; GenBank accession FJ386492).

Ancestral state and allele nomenclature

The 306-bp sequence of the GATA172D05 (allele 9 PCR fragment amplified by primers F1 and R; Fig. 1) was compared with the whole chimpanzee genome. The sequence with higher similarity was found at the X chromosome, with an alignment of 281 bp (95.9% identity). Two copies of the 18-bp sequence are present in the chimpanzee genome, which supports the hypothesis that this INDEL polymorphism arose with the loss of one of the copies during human evolution.

For the bi-allelic marker described here, we have adopted the nomenclature INDEL*L (long) for the ancestral duplication and INDEL*S (short) in order to designate the absence of one of the 18-bp tracts [14].

The nomenclature of the STR alleles was based on the number of repeats, as recommended by the International Society for Forensic Genetics [27] and is in accordance with Edelfmann et al. [28].

DNA of the cell line NA9947 (Promega Corporation, Madison, WI, USA) was also genotyped to be used as typing reference sample. The number of repeats is in agreement with the data previously published [29] and it does not contain the 18-bp duplication described in this work.

New primers design

In order to simultaneously type both STR and INDEL alleles, a new forward primer was designed (F2; see Fig. 1) which along with the same reverse primer amplifies a DNA fragment excluding the INDEL polymorphism. This new pair of primers (F2 and R) was used to type GATA172D05 STR alleles. The use of the three primers in a single PCR reaction allowed the simultaneous genotyping of the GATA172D05 STR and INDEL alleles (see Fig. 2). With this approach, it is possible to identify the genotype's gametic phase and determine haplotype frequencies, even in double heterozygote females (see for an example Fig. 2).

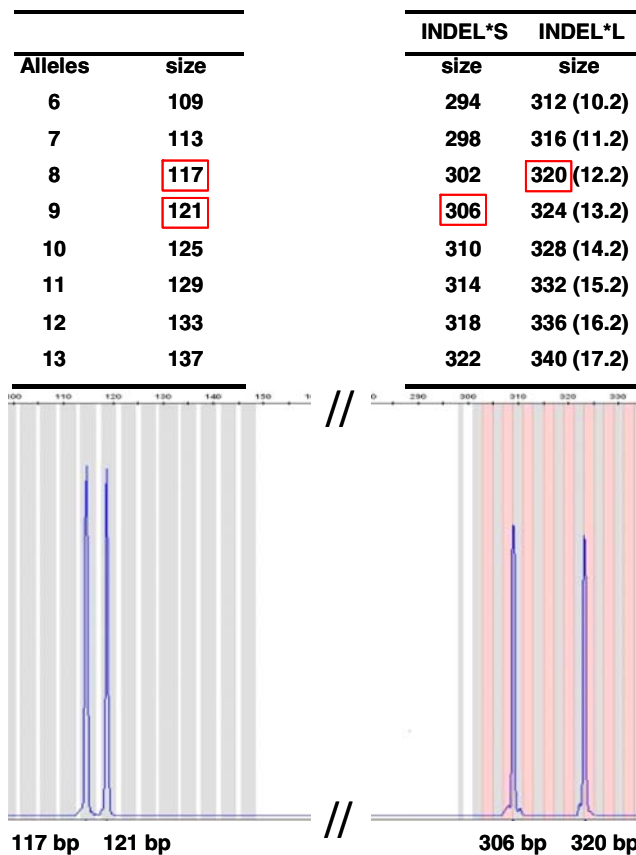


Fig. 2 Electropherogram of a double heterozygote female (STR*8–9; INDEL*L–S) illustrating how to determine the gametic phase based on the PCR product sizes, when using the three primers described in this work in a single PCR reaction. The table indicates the expected size for the different STR-INDEL alleles and haplotypes. In this case, the size of the two shorter amplicons (117 and 121 bp) allows the STR alleles typing (eight and nine repeats); and based on the size of the larger fragments (306 and 329 bp), it is possible to determine the INDEL state of each STR alleles (STR-INDEL*8-L/9-S). Note that the sizes in the table were obtained using the conditions described in the “Materials and methods” section, and different primer label, size standards, or apparatus usually alter these values

GATA172D05 INDEL variation

In a first screening of GATA172D05 INDEL alleles in a population sample from Belém, this marker was found to be polymorphic, with a high frequency of 70% for the INDEL*S-derived allele. Since this population is the result of miscegenation between Europeans, Africans, and Amerindians [21], in order to know whether this polymorphism is population specific, we expanded the investigation of genetic variation to other population groups, namely from European, African, and Amerindian ancestry. The results obtained (see Table 1) show that this polymorphism is present in all population groups studied. The derived allele is the most frequent in all populations except in Africans, but only marginally.

GATA172D05 STR variation

The STR allele frequencies observed in the four populations analyzed are depicted in Table 1. All alleles have already been described in other populations [15–19, 30, 31] and showed a variation between 6 and 13 TAGA repeat units. For all studied populations, genotype frequency distributions in females showed no significant deviation from Hardy–Weinberg equilibrium ($p > 0.05$) and exact tests of population differentiation demonstrated that the allele

Table 1 Allele and haplotype frequencies for two X chromosome loci in Belém (total number of chromosomes, $N=270$) and in Portuguese ($N=271$), Native American ($N=339$), and African ($N=183$) samples

GATA172D05	Belém	Native Americans	Portugal	Africans
INDEL				
*S	0.700	0.751	0.660	0.465
*L	0.300	0.249	0.340	0.535
STR				
6	0.092	0.091	0.162	0.022
7	0.022	0.021	0.004	0.016
8	0.122	0.141	0.167	0.081
9	0.122	0.003	0.070	0.325
10	0.341	0.336	0.317	0.349
11	0.219	0.303	0.174	0.131
12	0.078	0.096	0.107	0.065
13	0.004	0.009	0.000	0.011
INDEL-STR				
S-6	0.018	0.003	0	0
S-7	0	0	0.004	0
S-8	0.007	0.012	0.011	0.005
S-9	0.041	0	0.055	0.016
S-10	0.337	0.330	0.313	0.300
S-11	0.215	0.300	0.170	0.077
S-12	0.078	0.096	0.107	0.054
S-13	0.004	0.009	0	0.011
L-6	0.074	0.088	0.162	0.022
L-7	0.022	0.021	0	0.016
L-8	0.115	0.129	0.156	0.076
L-9	0.081	0.003	0.015	0.309
L-10	0.004	0.006	0.004	0.049
L-11	0.004	0.003	0.004	0.054
L-12	0	0	0	0.011
HET	0.805	0.765	0.812	0.791
PD_M	0.806	0.767	0.808	0.793
PD_F	0.940	0.912	0.938	0.932
MEC_T	0.784	0.733	0.820	0.768
MEC_D	0.664	0.601	0.660	0.645

INDEL alleles nomenclature is according to Mills et al. [14] where “L” stands for the long ancestral allele and “S” for the short allele (carrying the 18-bp deletion)

HET expected heterozygosity, MEC_T mean exclusion chance in trios involving daughters, MEC_D mean exclusion chance in mother/son duos, PD_M power of discrimination in males, PD_F power of discrimination in females

frequencies observed among males are not significantly different from those in females ($p > 0.05$).

GATA172D05 INDEL vs. STR variation

The frequencies of INDEL-STR haplotypes were determined in all population samples (Table 1).

Although almost all STR alleles could be detected in both INDEL backgrounds, a significant association was observed between alleles in these two linked markers (linkage disequilibrium p values < 0.001) in all populations. Chromosomes carrying the ancestral INDEL state present a high frequency of short STR alleles (from six to nine repeats) and large alleles, with more than ten repeats, are highly represented in the group of chromosomes with the derived allele.

The high INDEL diversity in all three European, Native American, and African groups is compatible with an old origin for this mutation in Africa, which is also supported by the high STR diversity inside each INDEL alleles, namely at the derived allele. The gene diversities inside both sets of INDEL*L and INDEL*S chromosomes in Africa are 63% and 53%, respectively. These values not only support the ancient origin of this polymorphism but also confirm the ancestral state of the allele carrying the 18-bp duplication. Indeed, the age estimate obtained from the STR variation associated with GATA172D05 INDEL*S chromosomes in Africa (see “Materials and methods” section for details) points to approximately 80,000 years for the time since the most recent common ancestor (mean = 79,830; posterior 95% interval = 19,634–250,117).

Haplotype diversities and forensic relevant parameters

When using the typing strategy described in this work, it could be possible to identify a high number of different haplotypes in all populations, varying between 11 and 13, in European and African samples, respectively. The relevant forensic parameters estimated based on the haplotype frequencies (Table 1) demonstrate its high diversity and usefulness to identification and kinship analysis in a very wide range of population backgrounds.

Conclusions

In this work, we have described the existence of and a convenient typing method for an 18-bp INDEL polymorphism, located upstream to the repeat structure of the GATA172D05 STR, which is widely employed both in anthropological and forensic studies. This INDEL is polymorphic in distinct ethnic groups, including Native Americans, Africans, and Europeans, which points to the

establishment of the polymorphism well before an out-of-Africa event. It has not escaped our attention that this INDEL polymorphism displays some very unusual features, namely the duplicated nature of the ancestral allele, the lack of homology of the repeat motif region between human and chimpanzee, and the correlation between the INDEL state and STR allele length, that can be interesting for evolutionary studies.

The forensic implications of our findings are: (1) the simultaneous typing of both STR and INDEL polymorphisms increases, without extra typing efforts, the informative power of GATA172D05 locus and (2) the straightforward haplotyping provided by our method can be very useful in kinship analyses, particularly in the so-called deficiency cases.

Acknowledgements This work was partially supported by Fundação para a Ciência e a Tecnologia (grant SFRH/BD/21647/2005 and POCI 2010, VI Programa-Quadro 2002–2006), Financiadora de Estudos e Projetos (FINEP), Milênio/CNPq (Conselho Nacional de Desenvolvimento Científico e Tecnológico), and UFPA.

References

1. Bedoya G, Montoya P, Garcya J et al (2006) Admixture dynamics in Hispanics: a shift in the nuclear genetic ancestry of a South American population isolate. *PNAS* 103:7234–7239
2. Laan M, Wiebe V, Khusnutdinova E, Remm M, Paabo S (2005) X-chromosome as a marker for population history: linkage disequilibrium and haplotype study in Eurasian populations. *Eur J Hum Genet* 13:452–462
3. Pereira RW, Pena SD (2006) Phylogeography of haplotypes of five microsatellites located in a low-recombination region of the X chromosome: studies worldwide and in Brazilian populations. *Genetica* 126:243–250
4. Toni C, Presciuttini S, Spinetti I, Domenici R (2003) Population data of four X-chromosome markers in Tuscany, and their use in a deficiency paternity case. *Forensic Sci Int* 137:215–216
5. Caine LM, Pontes L, Abrantes D, Lima G, Pinheiro F (2007) Genetic data of four X-chromosomal STRs in a population sample of Santa Catarina, Brazil. *J Forensic Sci* 52:502–503
6. Edelmann J, Lessig R, Klintschar M, Szibor R (2004) Advantages of X-chromosomal microsatellites in deficiency paternity testing: presentation of cases. *Int Congr Ser* 1261:257–259
7. Ballard DJ, Phillips C, Wright G, Thacker CR, Robson C, Revoir AP, Syndercombe Court D (2005) A study of mutation rates and the characterisation of intermediate, null and duplicated alleles for 13 Y chromosome STRs. *Forensic Sci Int* 155:65–70
8. Gusmão L, Amorim A, Prata MJ, Pereira L, Lareu MV, Carracedo A (1996) Failed PCR amplifications of MBP-STR alleles due to a polymorphism in the primer annealing region. *Int J Legal Med* 108:313–315
9. Gusmão L, Alves C, Costa S, Amorim A, Brion M, González-Neira A, Carracedo A (2002) Point mutations in the flanking regions of the Y-chromosome specific STRs DYS391, DYS437 and DYS438. *Int J Legal Med* 116:322–326
10. Clayton TM, Hill SM, Denton LA, Watson SK, Urquhart AJ (2004) Primer binding site mutations affecting the typing of STR loci contained within the AMPFISTR1 SGM Plus™ kit. *Forensic Sci Int* 139:255–259

11. Alves C, Gusmão L, Damasceno A, Soares B, Amorim A (2004) Contribution for an African autosomic STR database (AmpF/STR Identifiler and Powerplex 16 System) and a report on genotypic variations. *Forensic Sci Int* 139:201–205
12. Whitaker JP, Cotton EA, Gill P (2001) A comparison of the characteristics of profiles produced with the AMPFISTR[®] SGM Plus[™] multiplex system for both standard and low copy number (LCN) STR DNA analysis. *Forensic Sci Int* 123:215–223
13. Budowle B, Masibay A, Anderson SJ et al (2001) STR primer concordance study. *Forensic Sci Int* 124:47–54
14. Mills RE, Luttig CT, Larkins CE, Beauchamp A, Tsui C, Pittard WS, Devine SE (2006) An initial map of insertion and deletion (INDEL) variation in the human genome. *Genome Res* 16:1182–1190
15. Zarrabeitia MT, Alonso A, Martin J et al (2006) Study of six X-linked tetranucleotide microsatellites: population data from five Spanish regions. *Int J Legal Med* 120:147–150
16. Aler M, Sanchez-Diz P, Gomes I, Gisbert M, Carracedo A, Amorim A, Gusmão L (2007) Genetic data of 10 X-STRs in a Spanish population sample. *Forensic Sci Int* 173:193–196
17. Robino C, Giolitti A, Gino S, Torre C (2006) Development of two multiplex PCR systems for the analysis of 12 X-chromosomal STR loci in a northwestern Italian population sample. *Int J Legal Med* 120:315–318
18. Asamura H, Sakai H, Ota M, Fukushima H (2006) Japanese population data for eight X-STR loci using two new quadruplex systems. *Int J Legal Med* 120:303–309
19. Pereira R, Gomes I, Amorim A, Gusmão L (2007) Genetic diversity of 10 X chromosome STRs in northern Portugal. *Int J Legal Med* 121:192–197
20. Silva WA, Bortolini MC, Schneider MP, Marrero A, Elion J, Krishnamoorthy R, Zago MA (2006) MtDNA haplogroup analysis of black Brazilian and sub-Saharan populations: implications for the Atlantic slave trade. *Hum Biol* 78:29–41
21. Santos SEB, Guerreiro JF (1995) The indigenous contribution to the formation of the population of the Brazilian Amazon Region. *Brazil J Genet* 18:311–315
22. Sambrook J, Fritsch EF, Maniatis T (1989) Isolations of DNA from mammalian cells. In: Ford N, Nolan C, Ferguson M (eds) *Molecular cloning*. Cold Spring Harbor Laboratory Press, New York, pp 916–919
23. Excoffier L, Laval G, Schneider S (2006) Arlequin 3.1, An integrated software package for population genetics data analysis. University of Berne, Switzerland
24. Geisler WS, Diehl RL (2003) A Bayesian approach to the evolution of perceptual and cognitive systems. *Cogn Sci* 27:379–402
25. Weale ME, Yepiskoposyan L, Jager RF et al (2001) Armenian Y chromosome haplotypes reveal strong regional structure within a single ethno-national group. *Hum Genet* 109:659–674
26. Desmarais D, Zhong, Chakraborty R, Perreault C, Busque L (1998) Development of a highly polymorphic STR marker for identity testing purposes at the human androgen receptor gene (HUMARA). *J Forensic Sci* 43:1046–1049
27. Bär W, Brinkmann B, Budowle B et al (1997) DNA recommendations. Further report of the DNA Commission of the ISFH regarding the use of short tandem repeat systems. International Society for Forensic Haemogenetics. *Int J Legal Med* 110:175–176
28. Edelmann J, Deichsel D, Hering S, Plate I, Szibor R (2002) Sequence variation and allele nomenclature for the X-linked STRs DXS9895, DXS8378, DXS7132, DXS6800, DXS7133, GATA172D05, DXS7423 and DXS8377. *Forensic Sci Int* 129:99–103
29. Szibor R, Edelmann J, Hering S et al (2003) Cell line DNA typing in forensic genetics—the necessity of reliable standards. *Forensic Sci Int* 138:37–43
30. Shin SH, Yu JS, Park SW, Min GS, Chung KW (2005) Genetic analysis of 18 X-linked short tandem repeat markers in Korean population. *Forensic Sci Int* 147:35–41
31. Gomes I, Prinz M, Pereira R et al (2007) Genetic analysis of three US population groups using an X-chromosomal STR decaplex. *Int J Legal Med* 121:198–203