

A new method for the evaluation of matches in non-recombining genomes: application to Y-chromosomal short tandem repeat (STR) haplotypes in European males[☆]

L. Roewer^{a,*}, M. Kayser^b, P. de Knijff^c, K. Anslinger^d, A. Betz^e, A. Caglià^f,
D. Corach^g, S. Füredi^h, L. Henkeⁱ, M. Hidding^j, H.J. Kärger^k, R. Lessig^l,
M. Nagy^a, V.L. Pascali^f, W. Parson^m, B. Rolf^d, C. Schmitt^j, R. Sziborⁿ,
J. Teifel-Greding^o, M. Krawczak^p

^a*Institut für Rechtsmedizin, Humboldt-Universität Berlin, Hannoversche Straße 6, D-10115 Berlin, Germany*

^b*Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany*

^c*Forensic Laboratory for DNA Research, Leiden University, Leiden, The Netherlands*

^d*Institute of Legal Medicine, University of Munich, Munich, Germany*

^e*Landeskriminalamt Baden-Württemberg, Stuttgart, Germany*

^f*Institute of Legal Medicine, Catholic University, Rome, Italy*

^g*Servicio Huellas Digitales Genetics, Buenos Aires, Argentina*

^h*Institute for Forensic Sciences, Budapest, Hungary*

ⁱ*Institut für Blutgruppenforschung, Düsseldorf, Germany*

^j*Institute of Legal Medicine, University of Cologne, Cologne, Germany*

^k*Landeskriminalamt Sachsen-Anhalt, Magdeburg, Germany*

^l*Institute of Legal Medicine, University of Leipzig, Leipzig, Germany*

^m*Institute of Legal Medicine, University of Innsbruck, Innsbruck, Austria*

ⁿ*Institute of Legal Medicine, University of Magdeburg, Magdeburg, Germany*

^o*Bayerisches Landeskriminalamt, Munich, Germany*

^p*Institute of Medical Genetics, University of Wales College of Medicine, Cardiff, UK*

Received 11 February 2000; received in revised form 27 April 2000; accepted 28 April 2000

Abstract

A 9-locus microsatellite framework (minimal haplotype), previously developed for forensic purposes so as to facilitate stain analysis, personal identification and kinship testing, has been adopted for the establishment of a large reference database of male European Y-chromosomal haplotypes. The extent of population stratification pertaining to this database, an issue crucial for its practical forensic application, was assessed through analysis of molecular variance (AMOVA) of the 20 regional samples included. Despite the notion of some significant haplotype frequency differences, which were found to correlate with known demographic and historic features of Europeans, AMOVA generally revealed a high level of genetic homogeneity among the populations analyzed. Owing to their high diversity, however, accurate frequency estimation is difficult for Y-STR haplotypes when realistic (i.e. moderately sized) datasets are being used. As expected, strong pair-wise and higher order allelic associations were found to exist between all markers studied, implying that haplotype frequencies cannot be estimated as

[☆] Y-STR database address: <http://ystr.charite.de>

* Corresponding author. Tel.: +49-30-2093-7531; fax: +49-30-2093-7240.

E-mail address: lutz.roewer@charite.de (L. Roewer)

products of allele frequencies. A new extrapolation method was therefore developed which treats haplotype frequencies as random variables and generates estimates of the underlying distribution functions on the basis of closely related haplotypes. This approach, termed frequency ‘surveying’, is based upon standard population genetics theory and can in principle be applied to any combination of markers located on the Y-chromosome or in the mitochondrial genome. Application of the method to the quality assured reference Y-STR haplotype database described herein will prove very useful for the evaluation of positive trace-donor matches in forensic casework. © 2000 Elsevier Science Ireland Ltd. All rights reserved.

Keywords: Y-STR haplotype; Haplotype diversity; Analysis of molecular variance; Haplotype frequency estimation; Reference database

1. Introduction

The potential to identify male-specific DNA has rendered Y-chromosomal short tandem repeat (STR) systems invaluable for use in forensic genetics. Particularly in cases of sexual assault and in kinship testing, Y-STR haplotyping can help to close crucial informational gaps. An exceptionally informative core set of nine Y-linked STRs, analyzable in two to three multiplex reactions, has recently been recommended for court use in Europe [1,2]. To facilitate its practical application, several forensic laboratories have joined a collaborative effort to collect haplotype data from different populations and to create an adequate reference database. Here, we report on the current status of the Y-STR haplotype database for males from Europe or of European descent. The database, which is available online via the Internet (<http://ystr.charite.de>) is maintained at the Institute of Legal Medicine, Humboldt University, Berlin, Germany, where it is constantly updated by the participating laboratories. An ever increasing number of participants and the large number of external requests addressed to the database serve to emphasize the important role of Y-STR typing and the implementation of a central data repository in modern forensic genetics.

Since the Y-STR haplotype database project was initiated in 1998, four major objectives have been pursued: first, to establish a standardized, highly informative and stain-sensitive haplotyping method for mapped, sequenced and ‘multiplexable’ Y-chromosomal STRs; second, to introduce a means of minimum quality control for forensic laboratories using Y-STRs; third, to assess the extent of population stratification among males of European extraction; and finally, to obtain reliable estimates of Y-STR haplotype frequencies for use in forensic practice.

In order to quantify the positive evidence provided by a trace-donor match, or in order to facilitate like-

lihood calculations in the case of kinship testing, haplotype frequencies inevitably must be known. However, since most Y-chromosomal loci do not recombine, the so-called ‘product rule’ is not applicable and haplotype frequencies cannot be estimated simply as products of allele frequencies. Estimates involving the inverse of the database size, on the other hand, are overly conservative. We have, therefore, developed a numerical method to obtain Y-chromosomal haplotype frequency estimates via extrapolation, based upon the similarity and frequency relationships observed among haplotypes already included in the database. This approach ensures that rare haplotypes retain their high evidential power even when the database used for estimation is of realistic, and therefore moderate, size.

2. Materials and methods

2.1. The Y-STR haplotype reference database

The core set of loci referred to in the Y-STR haplotype database includes the following systems [1], all of which fulfill the ISFH guidelines for STR polymorphisms used in forensic practice [3]: DYS19, DYS389I, DYS389II, DYS390, DYS391, DYS392, DYS393, and DYS385I/II (also see <http://ruly70.med-fac.leidenuniv.nl/~fldo>). The above loci define a ‘minimal’ haplotype the analysis of which is obligatory for any laboratory participating in the database project. All STRs were typed as previously described [1,4–6]. Allele sizing was performed using sequenced allelic ladders and/or panels of typed control DNA provided by either of the following institutions: (i) Institute of Legal Medicine, Humboldt University, Berlin, Germany, (ii) Forensic Laboratory for DNA Research, Leiden University, The Netherlands, (iii) Institute of Legal Medicine, University of Mainz,

Germany. The system DYS385I/II detects variation at two loci simultaneously and unambiguous allele assignment by PCR is hampered by long stretches of invariant sequence flanking the two variable repeats [6]. In addition to the above STRs, roughly 50% of haplotypes logged in the database include the highly informative 2-locus system YCAII (extended haplotypes) for which the resolution of allelic origin is also not feasible [7,8].

Three prerequisites have had to be fulfilled by all laboratories participating in the Y-STR database project

1. successful passing of a quality control test that involves blind haplotyping of five samples for the nine core set STRs (plus YCAII for laboratories

submitting extended haplotypes); haplotype data from the literature are, therefore, not directly logged in the database,

2. continuous submission of 9-locus minimal or 11-locus extended haplotypes, and
3. provision of information about the ethnic and/or geographic origin of the probands typed.

By the time of the present study (September 1999), 14 laboratories have joined the project and 20 regional European population samples comprising a total of 2439 complete 9-locus minimal haplotypes have been logged in the database (Fig. 1, Table 1). Complete 11-locus extended haplotypes have been stored for 1197 males. For the purpose of comparison to more distantly related population, 78 haplotypes from Roma of

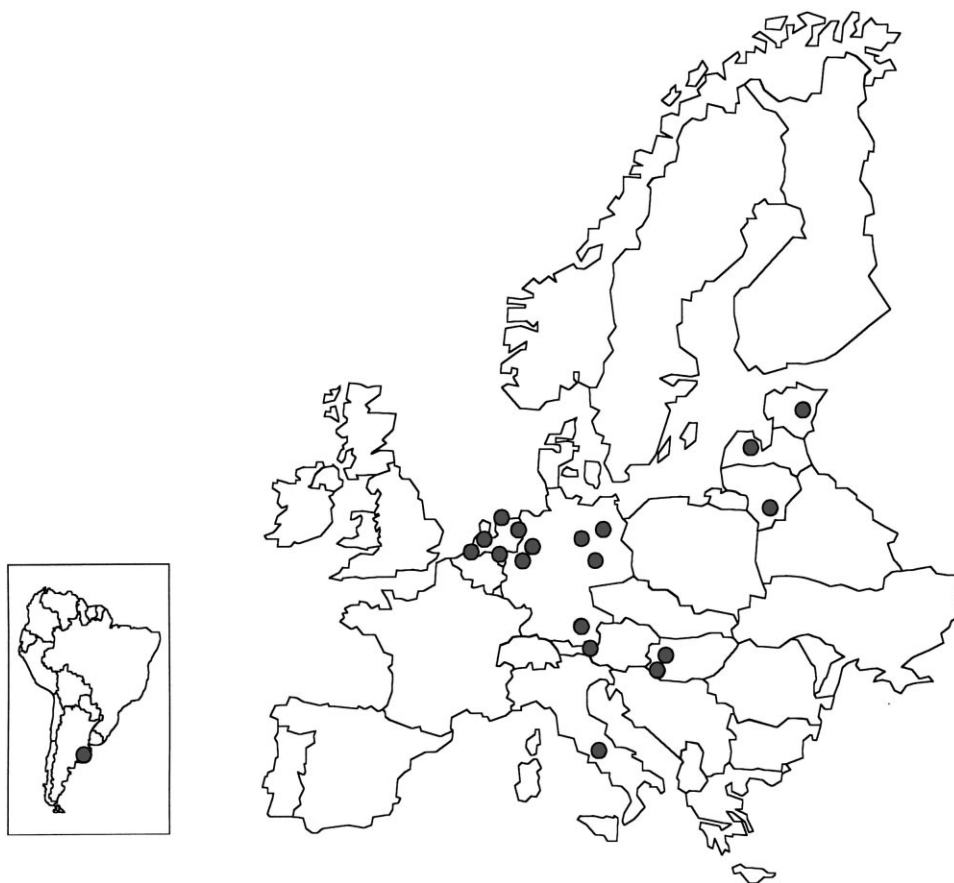


Fig. 1. Geographical origin of population samples used for Y-STR haplotype analysis (state from September 1999, for actual sampling coverage refer to <http://ystr.charite.de>).

Table 1
The European Y-STR haplotype reference database (September 1999)^a

Origin	Country	<i>N</i>	<i>M</i> ₇	<i>h</i>	S.E. (<i>h</i>)	<i>M</i> ₉
Berlin	Germany	239*	154	0.9908	0.00128	191
Düsseldorf	Germany	150	98	0.9837	0.00271	121
Cologne	Germany	135	98	0.9893	0.00235	117
Leipzig	Germany	359*	211	0.9946	0.00045	275
Magdeburg (LKA)	Germany	74	60	0.9919	0.00229	68
Magdeburg (ILM)	Germany	103*	76	0.9909	0.00200	91
Munich (LKA)	Germany	100	68	0.9836	0.00320	87
Munich (ILM)	Germany	151	107	0.9913	0.00165	128
Innsbruck	Austria	135*	94	0.9901	0.00188	109
Groningen	Netherlands	48	30	0.9628	0.00920	37
Friesland	Netherlands	44	34	0.9852	0.00463	42
Limburg	Netherlands	50	33	0.9722	0.00649	44
Zeeland	Netherlands	46	28	0.9449	0.01604	35
Unspecific	Netherlands	87	64	0.9826	0.00413	74
Rome	Italy	125	103	0.9962	0.00087	121
Budapest	Hungary	117	94	0.9958	0.00084	107
Riga	Latvia	145*	100	0.9907	0.00172	120
Vilnius	Lithuania	151*	100	0.9882	0.00205	122
Tartu	Estonia	80*	58	0.9867	0.00294	66
Buenos Aires	Argentina	100*	76	0.9877	0.00333	91
Total		2439	848	0.9930	0.00032	1419
Baranya-Roma	Hungary	78	26	0.9028	0.01388	32

^a *N*: number of core set (minimal) haplotypes in the database; *M*_{*k*}: number of different *k*-locus haplotypes observed (*k*=7, 9, 11), *h*: 7-locus haplotype diversity; S.E. (*h*): standard error of *h*; samples marked by an asterisk (Total=1197) have been typed for the extended haplotype (minimal+YCAII). By May 2000, the database has expanded to comprise 3859 haplotypes from 29 populations of Central, Eastern, Southern and Northern European origin.

the Baranya district in South Western Hungary have also been included in the study, but not in the database (Table 1). Haplotype data are collated in an SQL database at the Institute of Legal Medicine, Humboldt University, Berlin, Germany, and are available online via the Internet. Each observed haplotype is logged as a single entry. In addition to the respective allele designations, each record comprises a unique proband identifier (not available via the Internet) plus the proband's population affiliation. Frequency estimation currently constitutes the major practical use of the continuously growing database.

2.2. Analysis of molecular variance

Since Y-chromosomal haplotypes are confined to patriline, they are more prone to genetic drift than autosomal loci [9]. Thus, a comparatively high degree of differentiation between isolated populations can be expected to pertain to Y-chromosomal genetic data. To

assess the actual inter-population variability of Y-STR haplotypes among Europeans, analysis of molecular variance (AMOVA) [10] was performed using an in-house computer program (available from the authors upon request). AMOVA yields pair-wise Φ_{st} values which, similar to Wright's F_{st} , reflect the proportion of total molecular variance among two populations attributable to population differences. This approach had been employed before in the comparison of a Dutch and a German male population, using four Y-STRs (DYS19, DYS389I, DYS389II, DYS390) [11]. Since duplicated loci such as DYS385 and YCAII cannot be dealt with in AMOVA, the present study had to be confined to the 'classical' set of Y-STRs, comprising the seven unambiguous systems DYS19, DYS389I, DYS389II, DYS390, DYS391, DYS392, and DYS393. Distances between populations, as measured by Φ_{st} , were used to construct a similarity tree via the unweighted pair group method with arithmetic means (UPGMA) [12].

The AMOVA approach is based upon the molecular distance (d) between two haplotypes which equals the minimum number of mutational events separating them. For the present study, d values were calculated adopting a single step mutation (SSM) model. Inter-population variability as measured by Φ_{st} was tested for statistical significance via simulation. To this end, haplotypes observed in a given pair of populations were repeatedly redistributed at random between the two populations, with 1000 replications performed per population pair. The relative number of times a simulated Φ_{st} value exceeded the actual Φ_{st} represents a p value for testing $\{\Phi_{st}=0\}$ against $\{\Phi_{st}>0\}$.

2.3. Haplotype diversity

Haplotype diversity h was calculated as previously suggested by Melton et al. for genotypes from the mitochondrial control region [13]. Since h is supposed to reflect the likelihood of physical dissimilarity between two randomly drawn haplotypes, we chose to exclude the ambiguous STR systems DYS385 and YCAII from the analysis.

2.4. Allelic association

Pair-wise association between Y-STR systems was assessed by means of the standardized linkage disequilibrium coefficient, Δ' , calculated for the most frequent allele of each marker. If $+/-$ denotes the presence/absence of this allele from a given haplotype, and if f denotes the haplotype frequency, then

$$\Delta' = \frac{f_{+,+} \cdot f_{-,-} - f_{+,-} \cdot f_{-,+}}{f_{+,+} \cdot f_{-,-} + f_{+,-} \cdot f_{-,+}}.$$

The parameter Δ' takes values between -1 and $+1$, where $|\Delta'|=0$ indicates the absence of any association, and $|\Delta'|=1$ corresponds to the maximum association possible under the given allele frequencies. Pair-wise Δ' values were tested for significance (i.e. $|\Delta'|>0$) applying a χ^2 test with 1 degree of freedom to the 2×2 table defined by the respective f values.

2.5. Haplotype frequency estimation ('surveying')

Let N denote the total number of haplotypes in the Y-STR database, let M denote the number of different haplotypes logged, and let N_i denote the absolute

frequency of the i th haplotype. Since the first observation of the i th haplotype merely serves to indicate its existence, the basis for estimating its population frequency, f_i , is the observation of N_i-1 copies among $N-M$ haplotypes sampled. This implies that the maximum likelihood estimate of f_i , i.e. $(N_i-1)/(N-M)$, would be equal to zero for most haplotypes. On the other hand, using the upper 95% confidence limit of f_i for forensic purposes is overly conservative and would dramatically reduce the power of Y-STR haplotyping to evaluate trace-donor matches. At present, the upper limit of f_i would be 1:528 for every 7-locus haplotype with $N_i=1$, and 1:334 for every 9-locus haplotype with $N_i=1$.

According to classical population genetics theory [14], a Y-STR haplotype frequency f , when itself considered as a random variable, can be expected a priori to follow a β -distribution with density function

$$\varphi(f) = \frac{\Gamma(u+v)}{\Gamma(u) \cdot \Gamma(v)} f^{u-1} (1-f)^{v-1}.$$

Parameters u and v are functions of the effective population size and of the forward and backward mutation rate of the haplotype in question. Since these figures are as yet unknown for Y-STRs, we modeled u and v as functions of a related quantity, the weighted inverse molecular distance W of a haplotype from all other haplotypes in the database. This modeling approach, which invokes an SSM model via the use of W , can be interpreted as extrapolating prior distribution $\varphi(f)$ from a 'survey' of the surrounding haplotype frequencies. Let d_{ij} denote the molecular distance between the i th and j th haplotype. Then

$$W_i = \frac{1}{N} \sum_{j \neq i} \frac{N_j}{d_{ij}}.$$

Parameters u and v are related to the mean μ and the standard deviation σ of f via

$$u = \frac{\mu^2(1-\mu)}{\sigma^2} - \mu \quad \text{and} \quad v = u \left(\frac{1-\mu}{\mu} \right), \quad (1)$$

so that, instead of u and v , μ and σ could be modeled as functions of W . To this end, all European 9-locus haplotypes were divided into 15 equally sized groups, according to W value, and group-wise averages of W were determined. Intra-group μ and σ values were then calculated for the maximum likelihood estimates of f

within each group, followed by exponential regression of the two parameters on average W value.

Finally, in order to take the actually observed frequency of each haplotype into account, prior distributions φ have to be transformed into posterior distributions $\hat{\varphi}$ using a Bayesian approach. For the i th haplotype, we thus define

$$\hat{\varphi}(f_i) = \frac{P(N_i - 1|f_i)\varphi(f_i)}{\int_0^1 P(N_i - 1|g)\varphi(g) dg},$$

where $P(N_i - 1|x)$ is from a Bernoulli distribution with parameters x and $N - 1$. Note that $\hat{\varphi}$ is itself a β -distribution with parameters $u + N_i - 1$ and $v + N - N_i$.

In order to test whether the Y-STR haplotype frequencies as observed in the database fit the model developed here, we compared the observed numbers K_j of haplotypes with $N_i - 1 = j$ for $j = 0, 1, 2$ etc. with their expectations under the respective prior β -distributions. For the i th haplotype, u_i and v_i were thus calculated from W_i , using formula (1) and the regression equations for $\mu(W_i)$ and $\sigma(W_i)$. This allowed us to define densities $\varphi(f_i)$ and to calculate

$$E(K_j) = \sum_{i=1}^M \int_0^1 P(N_i - 1 = j|f_i)\varphi(f_i) df.$$

3. Results

3.1. Diversity of Y-STR haplotypes

For the ‘classical’ set of unambiguous Y-STRs (i.e. the core set except DYS385 I/II), 848 different haplotypes could be identified among the 2439 European males logged by September 1999, providing an $M:N$ ratio of 0.348 (Table 1). Upon extension to the full 9-locus core set, this number increased to 1419 ($M:N=0.582$). Further inclusion of YCAII yielded 895 different 11-locus haplotypes among 1197 fully typed males ($M:N=0.748$). The vast majority of theoretically possible haplotypes have not yet been observed in the database, and even the most frequent 9-locus haplotype is only shared by 55 males in the database. For most European populations analyzed, haplotype diversity h exceeded 0.98 even when based on the 7-locus alone. The smaller estimates of h obtained in three of the five Dutch samples may be

explicable in terms of small sample size, as is suggested by the relatively large standard errors observed. Bootstrapping analysis with 50,000 replications per population pair served to demonstrate that h was significantly smaller in Baranya-Roma than in all other European male populations (with the respective German, Dutch and non-Finno-Ugric speaking Baltic samples combined; $p < 0.0001$ in all pair-wise comparisons).

3.2. Population stratification

Whilst the largest pair-wise Φ_{st} value (0.21) was observed for Lithuania and Zeeland, all central European comparisons yielded $\Phi_{st} \leq 0.1$ (a complete list of Φ_{st} values is available for inspection via the Internet at <http://ystr.charite.de>). Indeed, most Φ_{st} values were extremely small and non-significant reflecting the genetic homogeneity of the central European male populations and subpopulations analyzed here. On the other hand, relationships between Y-STR haplotypes that were actually unraveled by AMOVA were clearly correlated with the known genetic and demographic history of the populations involved (Fig. 2). Thus, the nine Austro-German populations did not yield any significant pair-wise Φ_{st} values, but at the same time differed significantly from all other samples. The same holds true for the Dutch populations (despite a questionably high Φ_{st} for Friesland and Limburg that was possibly due to small sample size) and for males from the two non-Finno-Ugric speaking Baltic provinces of Lithuania and Latvia. Interestingly, even geographic subregions such as the Northern (Friesland, Groningen) and Southern (Zeeland, Limburg) Dutch provinces appeared to cluster. The consequent realization that 7-locus Y-STR haplotypes are extremely population-sensitive is most strikingly illustrated by the large distance observed between Baranya-Roma and other Europeans: the average Φ_{st} of 0.18 is similar to results from a global Y-STR survey (Kayser et al., in preparation) in which pair-wise Φ_{st} values between Europeans and non-Europeans were found to range between 0.2 and 0.5.

3.3. Haplotype frequencies

Most pair-wise combinations of Y-STR markers were characterized by a strong association between

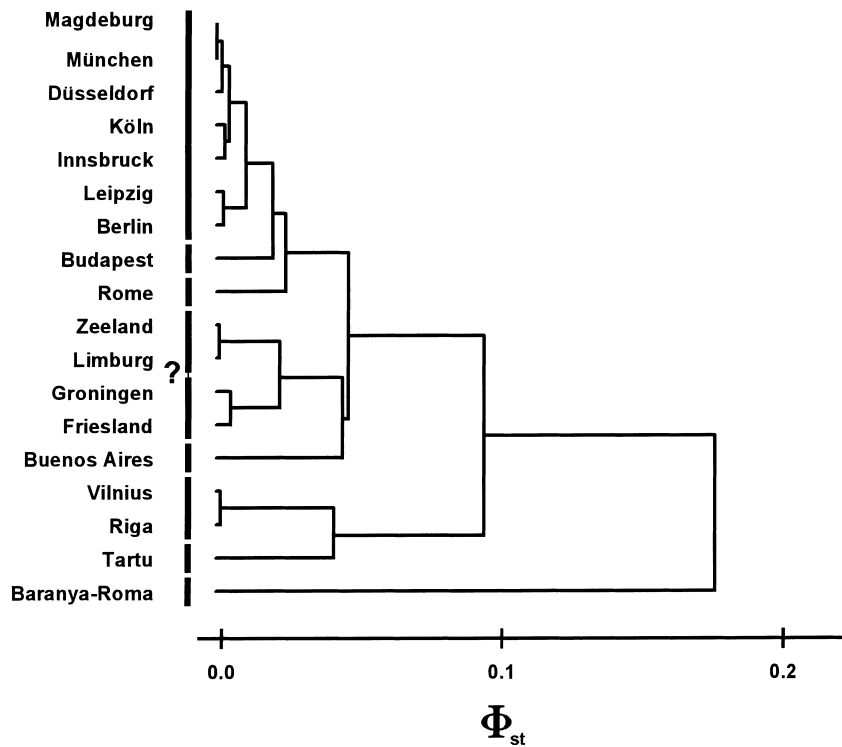


Fig. 2. Similarity tree (UPGMA) based upon pair-wise Φ_{st} values for 7-locus Y-STR haplotypes. Different population samples originating from the same geographical area have been pooled. ?: the split between the two Dutch subclusters was of border-line significance and resulted from a questionably high Φ_{st} for Friesland and Limburg.

their most frequent alleles (Table 2). Since this finding could have been due potentially to stratification among the male European populations studied, the analysis was repeated for the Austro-German haplotypes alone. Again, most pair-wise associations were

strong and highly significant (Table 2). Strong associations were also found between most other pairs of alleles and also at higher levels than just marker pairs (data not shown). These results highlight the fact that non-recombining genetic markers, such as Y-STRs,

Table 2
Pair-wise association (Δ') between the most common Y-STR alleles^a

Marker	DYS19	DYS389I	DYS389II	DYS390	DYS391	DYS392	DYS393
DYS19	–	0.009	0.316 ^d	0.178 ^d	–0.191 ^d	–0.717 ^d	0.136 ^c
DYS389I	–0.044	–	0.643 ^d	0.364 ^d	–0.356 ^d	–0.334 ^d	0.338 ^d
DYS389II	0.332 ^d	0.618 ^d	–	0.111 ^b	–0.120 ^c	–0.531 ^d	0.048
DYS390	0.110 ^b	0.422 ^d	0.129 ^b	–	–0.303 ^d	–0.395 ^d	0.186 ^d
DYS391	–0.364 ^d	–0.361 ^d	–0.149 ^c	–0.316 ^d	–	0.562 ^d	–0.162 ^d
DYS392	–0.778 ^d	–0.436 ^d	–0.575 ^d	–0.476 ^d	0.641 ^d	–	0.105 ^b
DYS393	0.234 ^d	0.371 ^d	–0.012	0.145 ^b	–0.245 ^d	–0.031	–

^a Upper right half: all male haplotypes ($N=2439$); lower left half: Austro-German male haplotypes only ($N=1446$).

^b $p<0.05$.

^c $p<0.01$.

^d $p<0.001$.

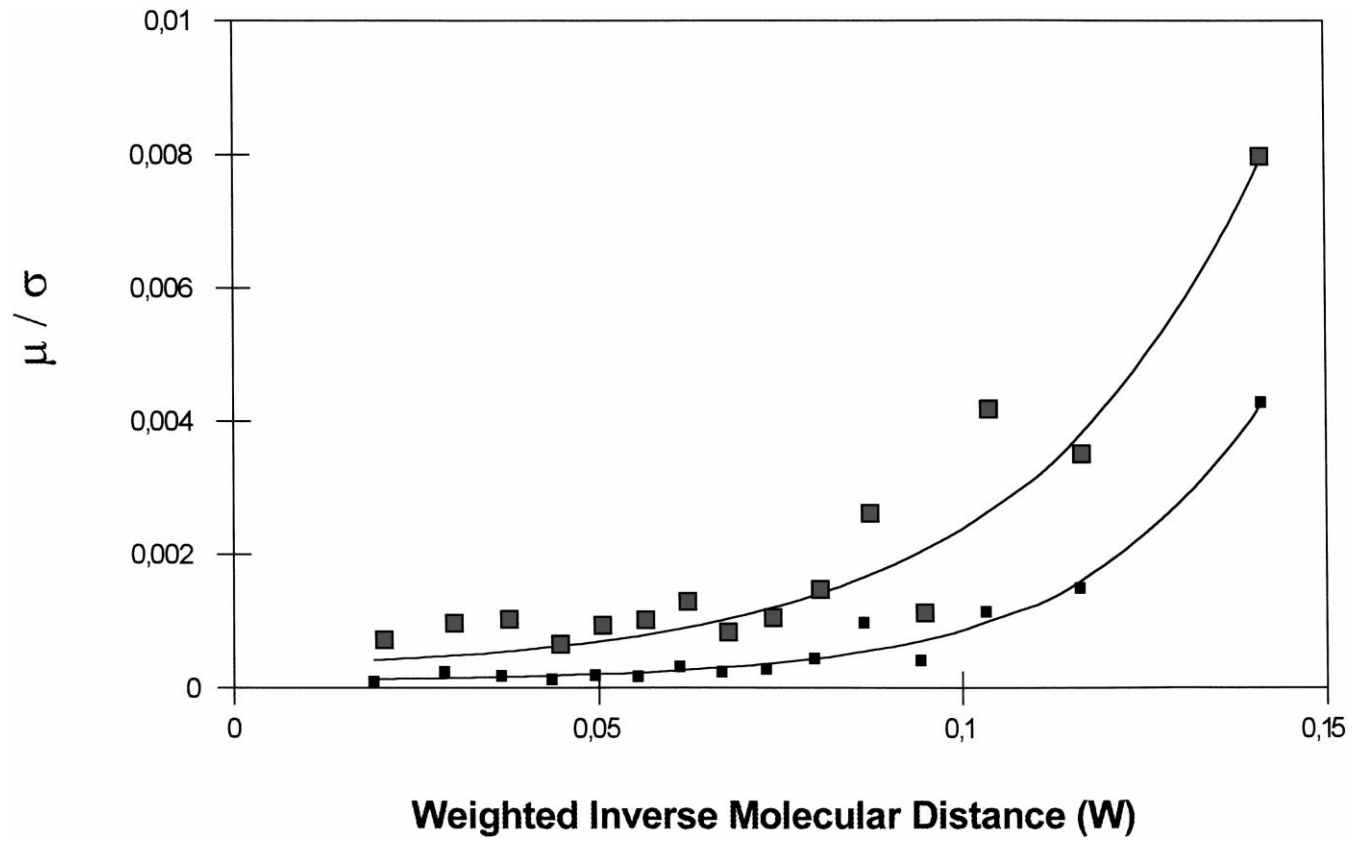


Fig. 3. Regression of 9-locus Y-STR haplotype frequency (f) on haplotype similarity to the reference database (W). Solid squares: mean (μ) of f in haplotype group; hatched squares: standard deviation (σ) of f in haplotype group.

Table 3
Exponential regression analysis of Y-STR haplotype frequencies^a

Parameter	Regression equation	Root mean square error
μ	$1.11 \times 10^{-4} + e^{41.20W - 11.30}$	1.67×10^{-4}
σ	$2.37 \times 10^{-4} + e^{30.86W - 9.22}$	6.19×10^{-4}

^a μ : mean of haplotype frequency estimates f ; σ : standard deviation of f ; W : weighted inverse molecular distance.

are not statistically independent entities and that haplotype frequencies are not normally equal to the products of allele frequencies.

Depicted in Fig. 3 is the relationship between the weighted inverse molecular distance (W) and the mean (μ) and standard deviation (σ) of the maximum likelihood haplotype frequency estimates in the 9-locus haplotype groups defined. Exponential regression gave a perfect fit for both parameters (Table 3) indicating that the population frequency of a given haplotype is positively correlated with the combined

frequency of closely related (surrounding) haplotypes in the same population. The observed numbers K_j of haplotypes with $N_i - 1 = j$ ($j=0, 1, 2$ etc.) also fit their expectations exceptionally well (Fig. 4). These findings suggest that our approach of haplotype frequency ‘surveying’, i.e. the assumption of Y-chromosomal haplotype frequencies a priori following a β -distribution combined with the exponential regression of the distribution parameters on haplotype similarities (W values), is indeed valid.

4. Discussion

Discriminative genetic analysis of the Y-chromosome has become an important tool in forensic science where it is used either as a complementary or as an alternative to other profile techniques. Regular quality control exercises such as that obligatory for laboratories participating in the present survey have helped to ensure that Y-STR systems can safely be

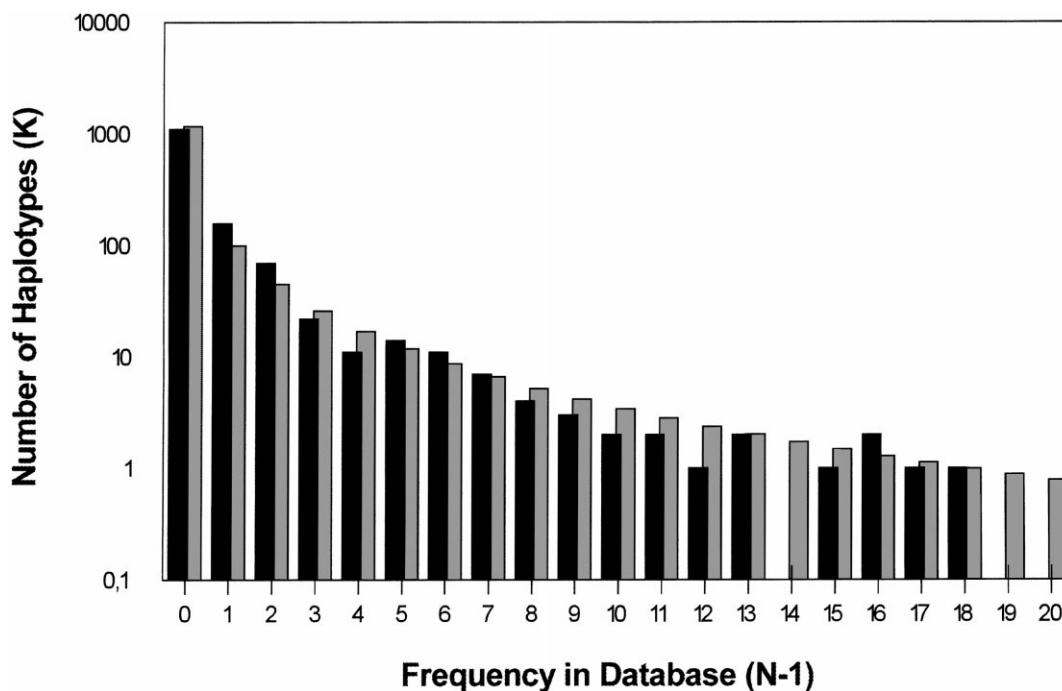


Fig. 4. Observed vs. expected frequency distribution of 9-locus Y-STR haplotypes. Given are the observed (solid bars) and expected (hatched bars) numbers of haplotypes K_j for which the observed frequency in the reference database equals $j+1$.

standardized between different forensic laboratories. The relevance of Y-STR haplotyping is most clearly evidenced by the increasing number of case studies and courtroom presentations published in the literature [15–19] that report on the practical use of Y-STRs following national guidelines (e.g. the TWGDAM program in the USA [18] or the EDNAP interlaboratory exercise in Europe (www.STADNAP.uni-mainz.de, [20])), and by the numerous regional population studies aimed at gathering standardized Y-STR haplotype data [19,21–26].

Since previous projects were generally based upon samples of as few as 100–200 haplotypes, their practical value in terms of reliable frequency estimation was limited. At best, only the most frequent haplotypes were appropriately represented. The present paper shows that, in contrast to local initiatives, international collaboration appears to represent a much more efficient way of obtaining sufficiently large samples for a major population such as the Europeans. Our database includes a representative number of regional subsets of European extraction, revealing however only a comparatively small extent of population stratification, and the large number of haplotypes logged allows frequency estimation to be carried out with high accuracy. In addition, the database is flexible and, if necessary, can easily be expanded beyond the nine STR systems currently used.

The project described herein serves to highlight the enormous potential of Y-STR haplotyping for forensic casework. It has been demonstrated that Y-STRs are powerful enough to allow unambiguous inter-individual (i.e. inter-lineage) discrimination with >99% probability in virtually all European populations studied. Non-identity of male trace donors is therefore an issue that could in principle be solved by Y-chromosomal genetic analysis alone. Even samples showing full identity for their 11-locus haplotypes could potentially be discriminated between using additional Y-STRs, such as those recently described by White et al. [27], or by SNP haplogrouping exploiting population-sensitive point mutations (Jobling et al., in preparation). In order to allow Y-STR or Y-SNP haplotyping to be performed at a later stage, aliquots of DNA samples have been stored by all groups participating in the current study.

In non-exclusion constellations, haplotype frequency estimates derived from a Y-STR reference database could in principle lend strong evidential power to the observation that an alleged father shares his Y-chromosomal haplotype with a putative son, or that a suspect's haplotype matches that of a stain. However, hitherto employed methods of haplotype frequency estimation based upon mere haplotype counting seemed inadequate since they either lead to meaningless results (e.g. $f_i=0$) or cause an unacceptably strong reduction in evidential power (e.g. by consideration of upper confidence limits). On the other hand, the strong allelic association between all Y-STRs implies that haplotype frequencies cannot be estimated simply as products of allele frequencies. The novel method of frequency 'surveying' described herein generates estimates of the prior and posterior frequency distributions of haplotypes instead of single quantities. Since the haplotypes in question do not enter into the estimation process of the prior distributions, they do not necessarily have to be present in the reference database themselves. In such cases, the observed frequency N_i would have to be set to 1 for the generation of the respective posterior distributions. Knowledge of the posterior distribution functions in turn allows various statistical measures of the haplotype frequency in question to be provided (e.g. mean, variance, percentiles), and should therefore be useful in assisting the decision making process ensuing from the observation of positive matches. The fact that the prior distributions result from extrapolation does not invalidate the approach. On the contrary, since it relies upon well-established population genetics principles, the method warrants as much credit as, for example, the estimation of composite autosomal genotype frequencies from the multiplication of genotype frequencies over loci, the latter also representing an extrapolation-type approach. Although frequency 'surveying' has been developed here for Y-chromosomal data, it can in principle also be applied to mitochondrial haplotype data. The methodology has been tested and validated on the 2439 haplotypes from 20 European populations logged by September 1999, meanwhile it also proven to yield consistent results on a largely extended version of the database.

In order to demonstrate the practical utility of the Y-STR haplotype reference database and the frequency surveying method, we have applied the combination

Table 4
Three forensic casework examples, investigated by Y-STR haplotyping

Case	Sample	DYS19	DYS389I	DYS389II	DYS390	DYS391	DYS392	DYS393	DYS385	f_{obs}	Mean f_{ext} (80% confidence interval)
1	Stain	16	13	29	–	10	11	13	–	56/2439	2.7×10^{-2} (2.3×10^{-2} – 3.2×10^{-2})
	Suspect	16	13	29	25	10	11	13	–	–	–
2	Stain	15	13	31	23	10	11	12	–	0/2439	8.0×10^{-5} (6.4×10^{-11} – 2.4×10^{-4})
	Suspect	15	13	31	23	10	11	12	–	–	–
3	Child	14	13	28	23	12	13	13	11, 14	0/2439	9.9×10^{-4} (1.3×10^{-4} – 2.2×10^{-3})
	Putative father	14	13	28	23	12	13	13	11, 14	–	–

of both to three casework examples: (1) In a rape case, a stain from the underwear of the victim was analyzed. No sperm cells could be identified, and autosomal STR analysis generated a female but failed to identify a male profile. With the Y-STR systems, an incomplete male Y-STR haplotype was typed in the stain and compared to a suspects DNA. The rapist confessed to the crime (Table 4, case 1). (2) In another rape case, a mixed stain containing epithelial and sperm cells as well as leukocytes was analyzed. Differential lysis has been applied but autosomal STR analysis failed to identify a male profile. Using the Y-STR systems, a 7-locus male haplotype was generated from the stain matching a suspects Y-STR profile (Table 4, case 2, described in [28]). (3) Paternity had to be investigated for a male child whose mother or other relatives were not available for testing. Identical minimal Y-STR haplotypes were observed in the child and alleged father (Table 4, case 3). Observed and estimated frequency values were determined in all three cases using an in-house computer program (available from the authors upon request). The more frequent a haplotype was, the closer was the mean frequency to the observed one, and the more narrow was the confidence interval for the frequency. In the case of unobserved (i.e. rare) haplotypes, confidence intervals were comparatively broad but nevertheless still conclusive.

Population analysis using AMOVA showed that, despite their high diversity, Y-STR haplotypes were found to yield small pair-wise Φ_{st} values, suggesting that haplotype distributions are comparatively homogeneous within European subpopulations. Although statistically significant differences could be discerned and related to historic/demographic relationships between subgroups, the actual level of stratification was found to be so low that, at least in Central Europe, the problem of choosing the correct population as a basis for evaluating a trace-donor match should be irrelevant to forensic practice. This robustness of Y-STRs against population differences parallels that previously demonstrated for autosomal microsatellites used for identification purposes [29,30]. In order to facilitate further investigation of this matter, however, the integration of existing and planned Y-chromosomal genetic databases and the establishment of generally agreed quality criteria are urgently required, and should represent a major concern of the forensic scientific community.

Acknowledgements

The authors wish to thank D.N. Cooper (Cardiff) for helpful discussion, and two anonymous reviewers for constructive comments. C. Krüger, M. Knopf and M. Meyer are acknowledged for technical assistance. M. Krawczak is supported by the Deutsche Forschungsgemeinschaft through a Heisenberg grant (Kr 1093/5-2).

References

- [1] M. Kayser, A. Cagliá, D. Corach, N. Fretwell, C. Gehrig, G. Graziosi, F. Heidorn, S. Herrmann, B. Herzog, M. Hidding, K. Honda, M.A. Jobling, M. Krawczak, K. Leim, S. Meuser, E. Meyer, W. Oesterreich, A. Pandya, W. Parson, A. Piccinini, A. Perez-Lezaun, M. Prinz, C. Schmitt, P.M. Schneider, R. Szibor, J. Teifel-Greding, G. Weichhold, P. de Knijff, L. Roewer, Evaluation of Y-chromosomal STRs: a multicenter study, *Int. J. Legal Med.* 110 (1997) 125–133.
- [2] V.L. Pascali, M. Dobosz, B. Brinkmann, Coordinating Y-chromosomal STR research for the Courts, *Int. J. Legal Med.* 112 (1998) 1.
- [3] Editorial, DNA recommendations-further report of the DNA Commission of the ISFH regarding the use of short tandem repeat systems, *Forensic Sci. Int.* 87 (1997) 179–184.
- [4] A.J. Redd, S.L. Clifford, M. Stoneking, Multiplex DNA typing of short-tandem-repeat loci on the Y chromosome, *Biol. Chem.* 378 (1997) 923–927.
- [5] M. Prinz, K. Boll, H. Baum, B. Shaler, Multiplexing of Y chromosome specific STRs and performance for mixed samples, *Forensic Sci. Int.* 85 (1997) 209–218.
- [6] P.M. Schneider, S. Meuser, W. Waiyawuth, C. Rittner, Tandem repeat structure of the duplicated Y-chromosomal STR locus DYS385 and frequency studies in the German and three Asian populations, *Forensic Sci. Int.* 97 (1998) 61–70.
- [7] N. Mathias, M. Bayes, C. Tyler-Smith, Highly informative compound haplotypes for the human Y chromosome, *Hum. Mol. Genet.* 3 (1994) 115–123.
- [8] H. Matsumoto, S.-I. Tsuruya, R. Tsuda, Y. Orihara, S.-I. Kubo, Japanese population study of a Y-linked dinucleotide repeat DNA polymorphism, *J. Forensic Sci.* 44 (1999) 588–591.
- [9] M.A. Jobling, A. Pandya, C. Tyler-Smith, The Y chromosome in forensic analysis and paternity testing, *Int. J. Legal Med.* 110 (1997) 118–124.
- [10] L. Excoffier, P.E. Smouse, J.M. Quattro, Analysis of molecular variance inferred from metric distances among DNA haplotypes: application to human mitochondrial DNA restriction data, *Genetics* 131 (1992) 479–491.
- [11] L. Roewer, M. Kayser, P. Dieltjes, M. Nagy, E. Bakker, M. Krawczak, P. de Knijff, Analysis of molecular variance (AMOVA) of Y-chromosome-specific microsatellites in two closely related human populations, *Hum. Mol. Genet.* 5 (1996) 1029–1033.

- [12] W.H. Li, D. Graur, *Fundamentals of Molecular Evolution*, Sinauer, Sunderland, 1991, pp. 106–108.
- [13] T. Melton, R. Peterson, A.J. Redd, N. Saha, A.S.M. Sofro, J. Martinson, M. Stoneking, Polynesian genetic affinities with Southeast Asian populations as identified by mtDNA analysis, *Am. J. Hum. Genet.* 57 (1995) 403–414.
- [14] J.S. Gale, *Theoretical Population Genetics*, Unwin Hyman, London, 1990, pp. 313–316.
- [15] D. Corach, A. Sala, G. Penacino, N. Ianucci, P. Bernardi, M. Doretto, L. Fondevbrider, A. Ginarte, A. Inchaurregui, C. Somigliana, S. Turner, E. Hagelberg, Additional approaches to DNA typing of skeletal remains: the search for ‘missing’ persons killed during the last dictatorship in Argentina, *Electrophoresis* 18 (1997) 1608–1612.
- [16] K. Honda, L. Roewer, P. de Knijff, Male DNA typing from 25-year-old vaginal swabs using Y chromosomal STR polymorphisms in retrial request case, *J. Forensic Sci.* 44 (1999) 868–872.
- [17] M. Kayser, C. Krüger, M. Nagy, G. Geserick, P. de Knijff, L. Roewer, Y-chromosomal DNA-analysis in paternity testing: experiences and recommendations, in: B. Olaisen, et al. (Eds.), *Progress in Forensic Genetics*, Vol. 7, Elsevier, Amsterdam, 1998, pp. 494–496.
- [18] M. Prinz, A. Ishii, M. Sansone, H.J. Baum, R.C. Shaler, Y chromosome specific STR testing and the US legal system, In: G. Sensabaugh, et al. (Eds.), *Progress in Forensic Genetics*, Vol. 8, Elsevier, Amsterdam, 2000, pp. 591–594.
- [19] C. Gehrig, M. Hochmeister, B. Budowle, Swiss allele frequencies and haplotypes of 7 Y-specific STRs, *J. Forensic Sci.* 45 (2000) 436–439.
- [20] P.M. Schneider, E. d’Aloja, B.M. Dupuy, B. Eriksen, A. Jangblad, A.D. Kloosterman, A. Kratzer, M.V. Lareu, H. Pfitzinger, S. Rand, R. Scheithauer, H. Schmitter, I. Skitsa, D. Syndercombe-Court, M.C. Vide, Results of a collaborative study regarding the standardization of the Y-linked STR system DYS385 by the European DNA profiling (EDNAP) group, *Forensic Sci. Int.* 102 (1999) 159–165.
- [21] A. Caglià, M. Dobosz, I. Boschi, E. d’Aloja, V.L. Pascali, Increased forensic efficiency of a STR-based Y-specific haplotype by addition of the highly polymorphic DYS385 locus, *Int. J. Legal Med.* 111 (1998) 142–146.
- [22] R. Lessig, J. Edelmann, Y chromosome polymorphisms and haplotypes in West Saxonia (Germany), *Int. J. Legal Med.* 111 (1998) 215–218.
- [23] E. Rossi, B. Rolf, M. Schürenkamp, B. Brinkmann, Y-chromosome STR haplotypes in an Italian population sample, *Int. J. Legal Med.* 112 (1998) 78–81.
- [24] C. Pestoni, M.L. Cal, M.V. Lareu, M.S. Rodriguez-Calvo, A. Carracedo, Y chromosome STR haplotypes: genetic and sequencing data of the Galician population (NW Spain), *Int. J. Legal Med.* 112 (1998) 15–21.
- [25] M. Gene, N. Borrego, A. Xifro, E. Pique, P. Moreno, E. Huguet, Haplotype frequencies of eight Y-chromosome STR loci in Barcelona (North-East Spain), *Int. J. Legal Med.* 112 (1999) 403–405.
- [26] M. Nata, B. Brinkmann, B. Rolf, Y-chromosomal STR haplotypes in a population from north west Germany, *Int. J. Legal Med.* 112 (1999) 406–408.
- [27] P.S. White, L.T. Owathan, L.L. Deaven, J.L. Longmire, New male-specific microsatellite markers from the human Y chromosome, *Genomics* 57 (1999) 433–437.
- [28] H.J. Kärger, H. Sackewitz, Beispiele der Anwendung von STR-Systemen in der Spurenkunde unter besonderer Berücksichtigung von Y-chromosomal Systemen, *Arch. Kriminol.* 204 (1999) 175–185 (abstract in English).
- [29] P. Gill, I. Evett, Population genetics of short tandem repeat (STR) loci, *Genetica* 96 (1995) 69–87.
- [30] R. Deka, M.D. Shriver, L.M. Yu, R.E. Ferrell, R. Chakraborty, Intra- and interpopulation diversity at short tandem repeat loci in diverse populations of the world, *Electrophoresis* 16 (1995) 1659–1664.