



## 2 A new self-learning computational method for footprints of early 3 human migration processes

4 Z. Juhász<sup>1</sup> · E. Dudás<sup>2</sup> · Horolma Pamjav<sup>2</sup>

5 Received: 13 February 2018 / Accepted: 28 June 2018  
6 © Springer-Verlag GmbH Germany, part of Springer Nature 2018

### 7 Abstract

8 We present a new self-learning computational method searching for footprints of early migration processes determining  
9 the genetic compositions of recent human populations. The data being analysed are 26- and 18-dimensional mitochondrial  
10 and Y-chromosomal haplogroup distributions representing 50 recent and 34 ancient populations in Eurasia and America.  
11 The algorithms search for associations of haplogroups jointly propagating in a significant subset of these populations. Joint  
12 propagations of Hgs are detected directly by similar ranking lists of populations derived from Hg frequencies of the 50 Hg  
13 distributions. The method provides us the most characteristic associations of mitochondrial and Y-chromosomal haplogroups,  
14 and the set of populations where these associations propagate jointly. In addition, the typical ranking lists characterizing these  
15 Hg associations show the geographical distribution, the probable place of origin and the paths of their protection. Compari-  
16 son to ancient data verifies that these recent geographical distributions refer to the most important prehistoric migrations  
17 supported by archaeological evidences.

18 **Keywords** Y-chromosomal and mtDNA haplogroups · Archaeogenetics · Artificial intelligence · Self-learning algorithm ·  
19 Clustering · Rank correlation

### 20 Introduction

21 Starting from the beginning of the millennium, paternal  
22 and maternal lineages based on Y-chromosomal MSY and  
23 mtDNA have been studied for population migration history  
24 in a chain reaction (Jobling and Tyler-Smith 2003; Underhill  
25 and Kivisild 2007; Karafet et al. 2008; Cinnioglu et al. 2004;  
26 Tambets et al. 2004; Semino et al. 2000; Yao et al. 2004;  
27 Pakendorf et al. 2007; Bermisheva et al. 2004; Simoni et al.  
28 2000; Quintana-Murci et al. 2004).

The mtDNA and Y-chromosomal lineages seem to sup-  
port the hypothesis that reconstructing the demographic  
history of human migration histories, which highlights a  
recent increase in effective population size, is compatible  
with admixture of both lineages between continents and  
geographic regions.

The last 10 years have witnessed a revolution in ancient  
DNA (aDNA) research. Genetic studies of ancient and  
modern populations significantly contributed to the picture  
drawn previously by archaeologists about prehistoric pro-  
cesses resulting in the contacts between different ancient  
cultures and populations by the next generation sequencing  
(NGS) methods (Skoglund et al. 2012; Lazaridis et al. 2014,  
2016; Fu et al. 2016; Haak et al. 2015; Allentoft et al. 2015).

The sequencing focus was previously limited to hyper-  
variable regions of mitochondrial DNA. Nowadays whole  
genome sequences are connected to the massive sequence  
throughput of next generation sequencing platforms able to  
target short and degraded DNA. Many ancient specimens  
being previously unsuitable for DNA analyses because of  
degradation can now successfully be used as templates for  
sequencing. At present, not only mitochondrial but also  
nuclear whole genomes have been sequenced from archaic

A1 Communicated by S. Hohmann.

A2 **Electronic supplementary material** The online version of this  
A3 article (<https://doi.org/10.1007/s00438-018-1469-7>) contains  
A4 supplementary material, which is available to authorized users.

A5 ✉ Horolma Pamjav  
A6 phorolma@hotmail.com

A7 <sup>1</sup> Centre for Energy Research, Institute of Technical Physics  
A8 and Materials Science, PO Box. 216, Budapest 1536,  
A9 Hungary

A10 <sup>2</sup> National Centre of Experts and Research, Institute  
A11 of Forensic Genetics, Budapest, Hungary

52 hominins, ancient anatomically modern humans, and pre-  
53 sent-day populations (Lazaridis et al. 2016; Fu et al. 2016;  
54 Haak et al. 2015; Allentoft et al. 2015; Batini et al. 2017;  
55 Lopopolo et al. 2016; Ilyas et al. 2015; Ermini et al. 2015;  
56 Der Sarkissian et al. 2015). Ancient DNA analysis of auto-  
57 some can provide detailed scenario of admixture. How-  
58 ever, populations in different geographic locations tend  
59 to have their own special sub-lineages of Y-chromosome  
60 and mtDNA. Therefore, the studies of Y-chromosome and  
61 mtDNA have potential to yield better resolution than that of  
62 autosome when studying the origin and migration of human  
63 population.

64 A powerful method based on PCA of the Fst distance  
65 matrix of 101 ancient individuals arising from the period of  
66 3400–200 BC indicated genetic transitions well correspond-  
67 ing to archaeological findings in Eurasia. The comparison  
68 to recent Fst data showed the connection between contem-  
69 porary and Bronze-Age populations (Allentoft et al. 2015).

70 A significantly different approach is based on cluster-  
71 ing of haplogroup (Hg) frequency distributions of recent  
72 and ancient populations, instead of analysis of pairwise Fst  
73 distances of individuals (Juhász et al. 2015, 2016). In this  
74 case, ancient and recent populations belonging to a common  
75 cluster directly point to the genetic connection of complete  
76 populations.

77 Representing populations by their Hg distributions pro-  
78 poses the following assumption: recent populations are prod-  
79 ucts of prehistoric and historic interactions, disjunctions and  
80 junctions of certain ancient source populations [for instance,  
81 the admixture of indigenous European hunter-gatherers with  
82 Neolithic farmers arising from the Middle East resulted in  
83 a new population with an Hg distribution containing both  
84 European and Middle Eastern components (Skoglund et al.  
85 2012)]. It has been verified by stepping stone simulation  
86 that such admixture processes starting from a few source  
87 areas result in strong correlations between Hgs arising from  
88 a common starting population, because they propagate nec-  
89 essarily jointly for a long time. Consequently, a search for  
90 strongly correlated Hgs in recent populations can reveal the  
91 Hg content of these ancient source populations (Juhász et al.  
92 2016; Neparáczki et al. 2017). However, iterative rank cor-  
93 relation search applied in these studies finds pairwise cor-  
94 relations, so larger subsets of the correlating Hgs are hardly  
95 detectable. In addition, that method did not utilize an impor-  
96 tant advantage of rank correlation technique, namely that the  
97 ranks attributed to the populations may refer to the emitting  
98 source populations of jointly propagating “Hg associations”.

99 In this paper, we present a new computational method  
100 aiming to reveal such groups of mitochondrial and Y-chro-  
101 mosomal haplogroups jointly propagating in a significant  
102 subset of contemporary Eurasian and American indigenous  
103 populations. In order to verify our starting assumption,  
104 namely that strong correlation of Hgs in recent populations

105 may refer to ancient source populations, we compare the  
106 results to ancient mitochondrial Hg distributions.

## 107 Materials and methods

### 108 Materials

109 We analysed 50 populations for mtDNA, as well as 50 pop-  
110 ulations for Y-chromosomal haplogroups. The frequencies  
111 for mtDNA haplogroups and Y-chromosomal haplogroups  
112 together with publication sources and a three-letter code  
113 was used to label each population as presented in Online  
114 Resource 1 (ESM\_1) and Online Resource 2 (ESM\_2). Fur-  
115 thermore, 34 ancient mtDNA haplogroup distributions were  
116 used for this study that is also included in Online Resource  
117 1 (ESM\_1). The total sums of individuals represented by  
118 mitochondrial, Y-chromosomal and ancient mitochondrial  
119 Hg distributions are 13,631; 6746 and 1266, respectively.

120 Mitochondrial and Y-chromosomal data do not per-  
121 fectly coincide in three cases, when Tuscany–Sicilian, Ser-  
122 bian–Croatian, as well as Karachay–Balkar data are com-  
123 bined, marked by the abbreviations TUS, SRB and KRC.  
124 Based on the close geographical, linguistic and historical  
125 contacts, we suppose that these couplings do not interfere  
126 the analysis in a significant manner.

127 The populations and the corresponding abbreviations of  
128 the modern data are shown in Table 1.

129 The ancient population mtDNA data, sample sizes, abbrevi-  
130 ations, places and times of origin are included in Table 2.

131 The aim of this study was to test the new method for  
132 the genetic results accepted by scientific community, so we  
133 did not focus on the resolution of haplogroups, therefore  
134 we used mainly the distribution of the basic haplogroups to  
135 compare as many populations as possible.

### 136 Methods

137 Here, we present a new computational method aiming to  
138 reveal all groups of Hgs jointly propagating in a significant  
139 subset of 50 contemporary Eurasian and American indige-  
140 nous populations. In rank correlation calculation of two Hgs,  
141 ranks are attributed to the 50 populations studied, according  
142 to the frequencies of the given Hgs in their Hg distributions.  
143 After that, the rank correlation coefficient is defined as the  
144 well-known linear correlation coefficient of the resulting two  
145 rank lists. Obviously, these 50-element rank lists of strongly  
146 correlating Hgs are necessarily similar therefore the whole  
147 set of the corresponding 50-dimensional vectors constitutes  
148 a clustered point system in its vector space. Thus, we reduce  
149 the search for groups of strongly correlating Hgs to clus-  
150 tering of their rank lists (vectors), instead of analysing the  
151 totally puzzling network of pairwise correlations.

**Table 1** The populations and the corresponding abbreviations of the modern data

Population name	Abbreviation	Population name	Abbreviation
Han Chinese	CHN	Mongolian	MNG
Kyrgyz	KYG	Chuvash	CHU
Tuscany	TUS	Bulgarian	BLG
Azeri	AZR	Turkish	TUR
Karachay	KRC	Hungarian	HUN
Slovak	SLK	Czech	CZH
Romanian	ROM	Kashubian Poles	PLK
Finnish	FIN	Norwegian	NOR
North German	GEN	South German	GES
French	FRA	Netherlands	DUT
Scottish	CO	Galiccia	GAL
Northwest Amerindian	NAW	Komi Zyryan	KOZ
Khanty	KHA	Serbian	SRB
Kurdish	KUR	Russian	RUS
Central Amerindian	NAC	Warmian	War
Poles	POL	Southern Amerindian	NAS
Greek	GRE	Estonian	EST
Saami	SAA	Karelian Finn	KRL
Ukrainian	UKR	Uyghur	UYG
Kazakh	KAZ	Mari	MRI
Tatar	TAT	Udmurt	UDM
Japanese	JPN	Székely	SEK
Altai Kazakh	ALK	Hui Chinese	HUI
Macedonian	MAC	Lithuanian	LIT
Tuvan	TUV	–	–

152 We describe the determination of jointly propagating  
153 groups of Hgs using self organizing cloud (SOC) clustering  
154 of the inverse rank vectors as follows.

155 The basic assumption of this work is that there may exist  
156 certain “Hg-associations” with characteristic compounds of  
157 mitochondrial and Y-chromosomal Hgs. It also seems a real-  
158 istic assumption that the members of these associations of  
159 Hgs were jointly emitted from certain “source populations”  
160 for a long period, therefore their correlation subsets and rank  
161 sequences became similar as a result of the migrations and  
162 admixtures in historic and prehistoric times. If this is true,  
163 the rank lists of the correlation subsets belonging to a given  
164 Hg-association may form different separable clusters.

165 In principle, the problem of finding of characteristic  
166 Hg-associations could be reduced to a clustering of the  
167  $44 \times 44$ -dimensional symmetric matrix containing the rank  
168 correlation values of the  $26 + 18 = 44$  mitochondrial and  
169 Y-chromosomal Hgs. Rank correlation is itself a similarity  
170 measure, therefore distance-based clustering algorithms like  
171 *k-medoids*, nearest neighbour, *k* nearest neighbours, maxi-  
172 mal relation probability could be applied for this purpose.  
173 However, the first experiments have shown that rank cor-  
174 relations show a rather fuzzy structure with hardly identifi-  
175 able clusters. Therefore, we developed another clustering

176 method based on training the so-called SOC algorithm by  
177 the “inverse rank vectors” derived from the iterative rank  
178 correlation algorithm. This method allowed us to simultane-  
179 ously identify both the Hg-associations propagating regu-  
180 larly together as genetic components of certain propagating  
181 source-populations, and the groups of these propagating  
182 populations themselves.

183 The “inverse rank vector” (IRV) of a Hg is defined as  
184 follows.

185 Firstly, we execute the iterative rank correlation search  
186 for each pair of our  $26 + 18 = 44$  Hgs. Due to the iterations  
187 eliminating the populations causing the largest decrease of  
188 the correlation, the algorithm correlating the *i*th Hg (Hgi)  
189 to 44 other Hgs (including Hgi itself) results in 44 different  
190 rank lists for Hgi. We select the rank lists belonging to a cor-  
191 relation value higher than a critical value (0.7) from this set  
192 of rank lists. After finishing the whole process, we obtain a  
193 set of rank lists, each of them belonging to a strong correla-  
194 tion, while all other couplings of Hgs having no detectable  
195 correlation are eliminated.

196 Let  $r_k(i)$  denote the rank of the *i*th population in the rank  
197 list of the *k*th member of the above rank list set. The cor-  
198 responding “inverse rank” value is defined as

$$199 \bar{r}_k(i) = 1 - r_k(i) / \max(r_k(i)) \quad (1)$$

**Table 2** The ancient populations, mtDNA sample sizes, abbreviations, places and times of origin are in the table

Middle East Neolithic-BrA	28	MEN	Middle East	11,840–1402
Iberian Neolithic	45	IBN	Iberia	10,310–3160
Near Eastern Neolithic	67	NEO	TR, IRN, SYR, JOR	8300–4000
Cisbaikalian Neolithic (Serovo)	15	SER	East Siberia	8000–4000
Early-Middle Neolithic	53	EMN	Europe	6000–3000
Starcevo	44	STR	Balkans	5700–5500
Dniepr–Donets Neolithic	17	DDO	Eastern Europe	5300–4700
Neolithic Hungary	85	NHU	Hungary	5200–4800
Yamnaya, Afanasievo	49	YAM	Russia, Ukraine	5000–2700
Kurgans (Eneolithic/Catacomb)	35	KGC	UKR, MLD, BLG	4700–2000
Baraba (UT-ODI-EK)	33	BB1	West Siberia	4000–1800
Late Neolithic-EBA Europe	56	LNB	Europe	3000–1600
Altai Bronze Age	12	ABA	South Siberia	2700–900
Tarim Basin Xiaohe	73	XIA	China	2515–1829
Sintashta-Andronovo	41	SIA	Russia, Siberia	2300–1400
Baraba (LK-FYOD-LBB)	45	BB2	West Siberia	1800–1000
Srubnaya	14	SRU	Russia	1800–1200
Bronze Age Kurgans	13	KBK	Kazakhstan	1400–1000
Baraba (Iron transition)	14	BB3	West Siberia	1000–800 BC
Iron Age Kurgans	13	KIK	Kazakhstan	800–600 BC
Tagar–Tachtyk	15	TAG	Russia	800 BC–400 AD
Scythian Iron age	14	SCI	Russia	600–200 BC
Pazyryk Scytho-Siberian	25	PAZ	Mongolia, Russia	400–200 BC
Qin China aDNA	19	QIN	East Asia	221 BC–210 AD
Egyin Gol Xiongnu	46	XIO	Inner Asia	200 BC–200 AD
Lombard early medieval	40	LOM	Hungary, Italy	500–800 AD
Vikings	65	VIK	Norway	780–790 AD
Karos	90	KAR	Hungary	850–900 AD
Hungarians 900 AD	27	AH2	Central Europe	900–1000 AD
Ancient Hungarian (tenth century)	67	AH1	Central Europe	900–1000 AD
Pre-conquest Hungary	49	HPC	Hungary	500–900 AD
Medieval Slavic	19	SLV	Slovakia	900–1200 AD
Italian medieval	27	ITM	Italy	900–1500
Cumanian	11	CUM	Hungary	1200–1300 AD
Total sample size	1266			

200 for the populations contained by the  $k$ th correlation subset,  
 201 while the populations missing from the correlation subset  
 202 have the inverse rank value  $\bar{r}_k(i) = 0$ .

203 Using this definition, the population having the highest  
 204 frequency—and consequently the lowest rank in the  
 205 ranking list of  $Hgi$ —becomes the highest “inverse rank”  
 206 value approaching 1, and the ranks of the other popula-  
 207 tions decrease with decreasing frequency of  $Hgi$ . Thus,  
 208 our “inverse rank” is really the inverse of the common rank  
 209 definition, which increases with decreasing frequency in  
 210  $Hgi$ . The inverse ranks  $\bar{r}_k(i)$  are stored in the  $N$ -dimensional  
 211 “inverse rank vectors (IRV)”  $\underline{R}(k)$ , where the  $i$ th vector ele-  
 212 ment represents the inverse rank of the  $i$ th population in the  
 213  $k$ th correlation subset. The serial numbers of the pair of  $Hgs$   
 214 belonging to the  $k$ th correlation subset are also stored by the

215 algorithm. Three examples of similar inverse rank vectors  
 216 are illustrated in Fig. S2a (ESM\_3). The horizontal axis con-  
 217 tains the serial numbers of our 50 populations in an ad hoc  
 218 order, while the inverse rank values are represented by the  
 219 corresponding column heights. (The order of the populations  
 220 has no significance in the calculations.) Populations elimi-  
 221 nated by the iterative process have zero inverse rank values.

222 Due to the operation of the iterative rank correlation  
 223 algorithm, typically 15–20 vector elements dominate in the  
 224 50-dimensional IRV  $\underline{R}(k)$ , whereas the remaining elements  
 225 are negligible or zero. The differences of these small com-  
 226 ponents are also small, reducing the Euclidean distances of  
 227 the vectors, so they damp the essential differences in the  
 228 calculation. This problem was solved by a more advanced  
 229 version of the SOC with a weighted Euclidean distance

measure highlighting the important components of each inverse rank vector (Juhász et al. 2016). These 50-dimensional weight vectors were also learned automatically during the training process. The mathematical description of the algorithm is given in ESM\_3.

The complete analysis based on our self-learning computer programs accomplishing iterative rank correlation search for strongly correlating pairs of Hgs and SOC-clustering of the resulting IRVs can be summarized in two steps:

1. Collecting all Hg pairs having strong rank correlation values for a significant set of populations, using the iterative rank correlation method. The algorithm results in two IRVs for all pairs of these correlating Hg-pairs after finishing the iteration. For instance, the leftmost black columns in Fig. S2a (ESM\_3) show the inverse rank values of mitochondrial Hg M strongly correlating with Y-chromosomal Hg O. The strong correlation between these Hgs is verified by the visible similarity of the corresponding IRV of Y-chromosomal Hg O, represented by the neighbouring grey columns.
2. The resulting vector set is used as the training set of the SOC algorithm determining the condensation centres of the corresponding 50-dimensional point system. The resulting “inverse rank vector type” (IRVT) vectors are used for clustering the whole IRV set, and the mitochondrial and Y-chromosomal Hgs belonging to IRVs assigned to a common cluster are collected into a set called the “Hg-association” of the cluster.

A more detailed mathematical description of the iterative rank correlation and self organizing cloud (SOC) algorithms is found in ESM\_3 and (Juhász et al. 2016).

The goodness of the result was characterized by a calculation based on the correlation of the distance and inference matrices: after the clustering process, the inference matrix (containing values 1 when a pair of training vectors belongs to identical cluster and 0 if not) is determined. The goodness of the clustering is characterized by the correlation coefficient of the lower triangles (without the diagonal elements) of the symmetric inference matrix and the distance matrix of the training vectors.

It proved to be very favourable to put the results into historical context by extending the analysis by ancient data. To do this, we completed our 50 recent mitochondrial data with 34 ancient mtDNA distributions, and accomplished the whole analysis with the resulting 84-dimensional inverse rank vectors. To compare the recent part of the resulting 84-dimensional IRVTs to the original 50-dimensional recent ones, we had to eliminate ancient components and recalculate the rank values of the remaining recent populations within the resulting 50-dimensional modified IRVTs.

## Results

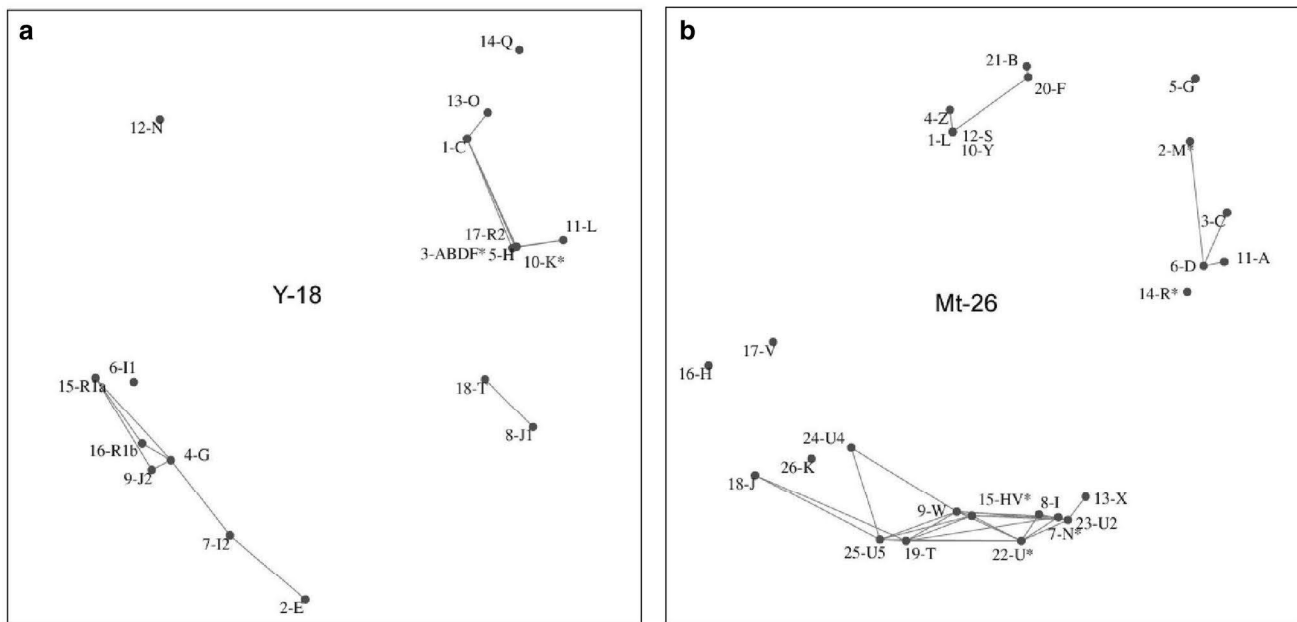
In the first step of the study, we accomplished the iterative rank correlation for all pairs of the  $26+18=44$  mtDNA and Y Hgs (including self-correlations). Subtracting the correlation coefficients from unity we obtain distance-like values approaching 0 and 2 in cases of strong positive as well as negative correlation, while this “distance” approaches 1 for uncorrelated Hg-pairs. Thus, the relationships of Hgs, determined by their correlations can be visualized by MDS maps, as it is shown in Fig. 1a, b for Y-chromosomal and mitochondrial Hgs.

High correlations refer to systematic joint propagation of pairs of Hgs within a significant subset of populations. Therefore, we selected all pairs of IRVs of the pairs of strongly correlating Hgs into an IRV set, with the constraint that the iteration resulted in a correlation exceeding 0.8 for a subset of populations exceeding the size of 15. We illustrate the results in Fig. S2 (ESM\_3).

Finally, we obtained a set of selected IRVs containing 393 elements, and we trained the self-learning cloud (SOC) algorithm to determine all the characteristic local condensation centres within the corresponding 50-dimensional point system simultaneously. SOC learning resulted in an IRVT set of 10 elements, and the t-probe showed that the distances of the closest neighbouring IRVTs is significant with probability at least 95%. The correlation of the inference and distance matrices was  $-0.52$ .

Ordering the 393 IRVs to the most similar IRVT, we obtained ten clusters. The subset of Hgs whose IRVs belong to a common cluster build the “Hg association” propagating within the populations having nonzero inverse rank values in the IRVs of the cluster.

The geographical distribution of the IRVT discussed here is shown in the map of Eurasia in Fig. 2. (As the SOC algorithm ordered the serial number of 10 to this IRVT, we sign it as IRVT-10.) The heights of the columns show the inverse ranks of the corresponding populations, so the map shows a propagation from Eastern- and Inner Asia (CHN, HUI, JPN, MNG, KYG, KAZ) to the native Americans (NAW, NAC, NAS) and to the Volga region in Eastern Europe (TAT, CHU, UDM). The Y-chromosomal and mitochondrial members of this Hg association are shown in maps mirroring the correlation conditions of the Hgs in the right upper part of the figure. The columns ordered to the Hgs are proportional to the number of other Hgs strongly correlating with it. For example, Y-chromosomal O and mitochondrial A haplogroups have the most correlating partners within the Hg association (Y: O, C, Q, N; Mt: Z, B, F, G, M, C, D, A, N\*). As this Hg association was derived from the whole correlation subset belonging to IRVT-10, certain Hgs may be absent from different



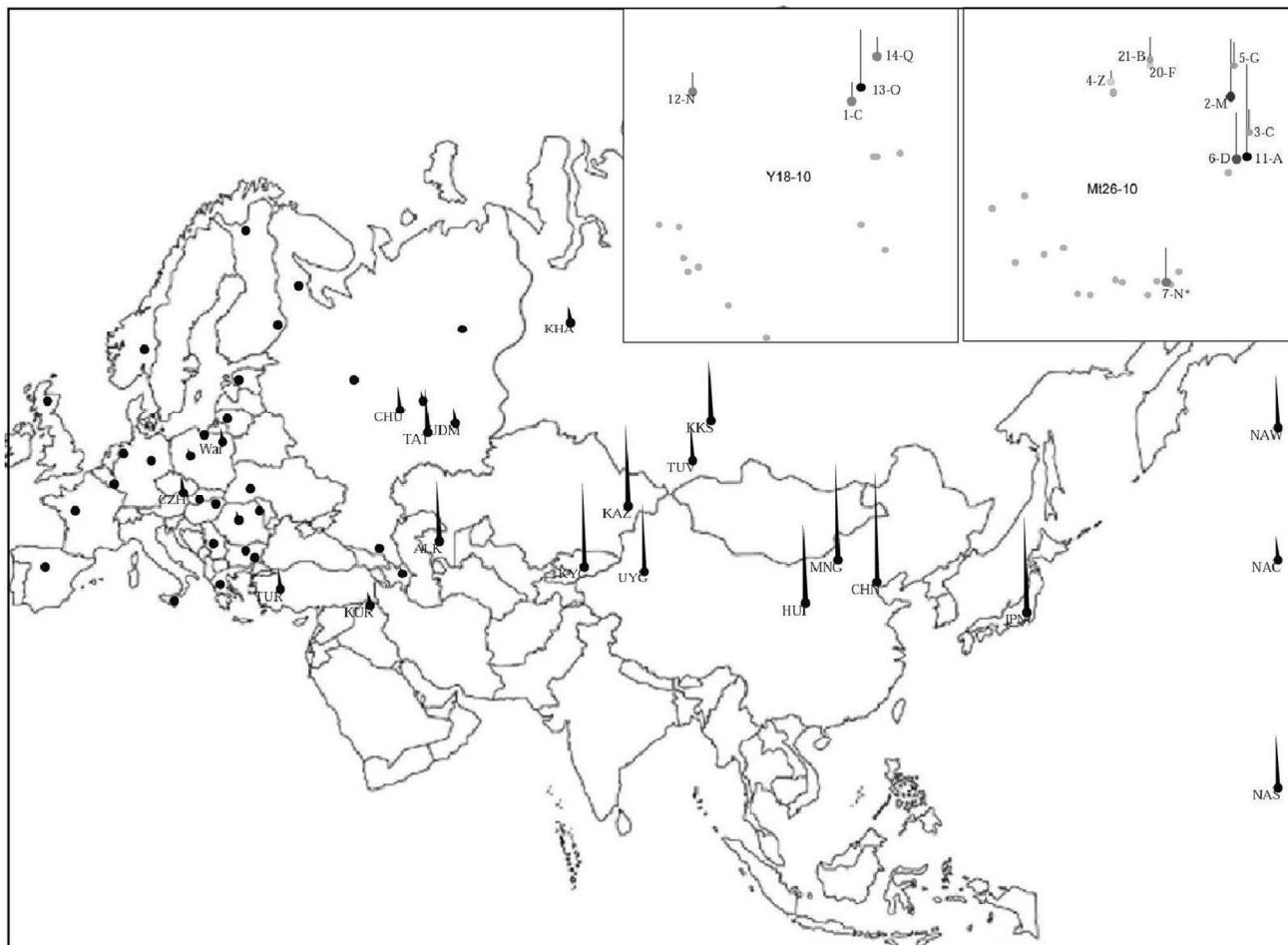
**Fig. 1** MDS map of Y-chromosomal (a) and mitochondrial Hgs (b), determined from rank correlation data

332 individual populations. For example, Central Amerindian  
 333 sample NAC contains only 3 Hgs of the Hg association (A,  
 334 B and a low rate of D). This is the reason why NAC has a  
 335 lower IR value than the neighbouring populations NAW  
 336 and NAS in the geographical map of Fig. 2.

337 To understand the historical background of the propaga-  
 338 tion of the above Hg association, we accomplished the whole  
 339 analysis on our mitochondrial database completed with 34  
 340 ancient mtDNA distributions. As this extended “historical”  
 341 database contained 84 populations in sum, the IRVs in this  
 342 analysis had 84 dimensions. We found that an appropriate  
 343 significance of the clustering was reached for 44 IRVTs.  
 344 This high number may be caused partly by the increased  
 345 dimension, partly by the small sample sizes of the ancient  
 346 mtDNA distributions resulting in a high noise level. To com-  
 347 pare the recent parts of the resulting 84-dimensional and the  
 348 originally recent 50-dimensional IRVTs, we eliminated the  
 349 ancient components of the 84-dimensional IRVTs and re-  
 350 calculated the modified inverse rank values of the remaining  
 351 50 recent components. After this, we selected the modified  
 352 IRVTs having the less Euclidean distances from the origi-  
 353 nal 50-dimensional IRVT-10. Finally, we turned back to the  
 354 84-dimensional original versions of the selected modified  
 355 IRVTs and found that the highest inverse ranks of the com-  
 356 plete versions are systematically attributed to the ancient  
 357 samples arising from South Siberia, Inner Asia and China  
 358 (ABA, Altai Bronze Age 2700–900 BC), SER (Serovo, East  
 359 Siberia 8000–4000 BC), PAZ (Pazyryk, Scytho-Siberian,  
 360 Mongolia, Russia 400–200 BC), XIA (Tarim Basin Xiaohe,  
 361 China, 2515–1829 BC), XIO (Egyin Gol Xiongnu, Inner

362 Asia, 200 BC–200 AD) and QIN (Qin China East Asia, 362  
 221 BC–210 AD). The historical migration transferring this 363  
 Hg association to Eastern Europe is also verified by the not 364  
 negligible inverse ranks in Iron-age Scythian (SCI) and ninth 365  
 century Hungarian samples (KAR, AHI). In addition, the 366  
 mitochondrial Hg content of the historical Hg associations 367  
 (Mt: B, F, G, M, C, D, A, N\*) is in a very good accordance 368  
 with those of IRVT-10 (the recent IRVTs with their most 369  
 similar historical pairs are shown in the maps of ESM\_4 370  
 also showing the mitochondrial Hg content of the corre- 371  
 sponding Hg associations. See Figs S1 and S2 in ESM\_4). 372  
 The accumulated rate of the Hg association in ancient East- 373  
 ern samples (XIO, XIA, QIN) is in the range of 75–95%. 374  
 The most Western appearance of this Hg association was 375  
 detected in early Hungarian samples arising from 800 to 376  
 1000 AD, where it takes 15–27% of the whole distributions. 377  
 These results verify that recent distribution of IRVT-10 and 378  
 the corresponding Hg association is a consequence of the 379  
 migrations of Scythians, Huns, Avars, Hungarian conquer- 380  
 ors, Cumanians and other nomadic people on the Steppe. 381

382 A totally different geographical distribution characterizes 382  
 the Hg association belonging to IRVT-1. This distribution 383  
 represented in Fig. 3 shows a propagation from the Middle 384  
 East and Asia Minor (KUR, AZE, TUR) to the Balkans and 385  
 Central Europe. The Hg association also appears in Eastern 386  
 Europe (TAT, CHU, RUS) and Inner Asia (UYG, KYG). The 387  
 Y-chromosomal and mitochondrial correlation-based maps 388  
 of the corresponding Hg association show the Hg-content 389  
 (Y: E, G, J1, J2, L, R1b, T, Mt: N\*, I, X, HV\*, U\*, U2, K, 390  
 T, J). Comparing these maps to those shown in Fig. 2 shows 391



**Fig. 2** Geographical distribution of the inverse ranks of IRVT-10. The Hg association propagating from East- and Inner Asia to Eastern Europe (CHU, TAT, UDM) and to the Native Americans (NAS, NAC, NAW) is shown in the right upper part. *Y18* Y-chromosomal, *Mt26* mitochondrial. The symbols referring to the Hgs contain a serial number and the name of the Hg (e.g. 12-N in the map “Y18-10” means Y-chromosomal Hg N having the serial number 12. 4-Z

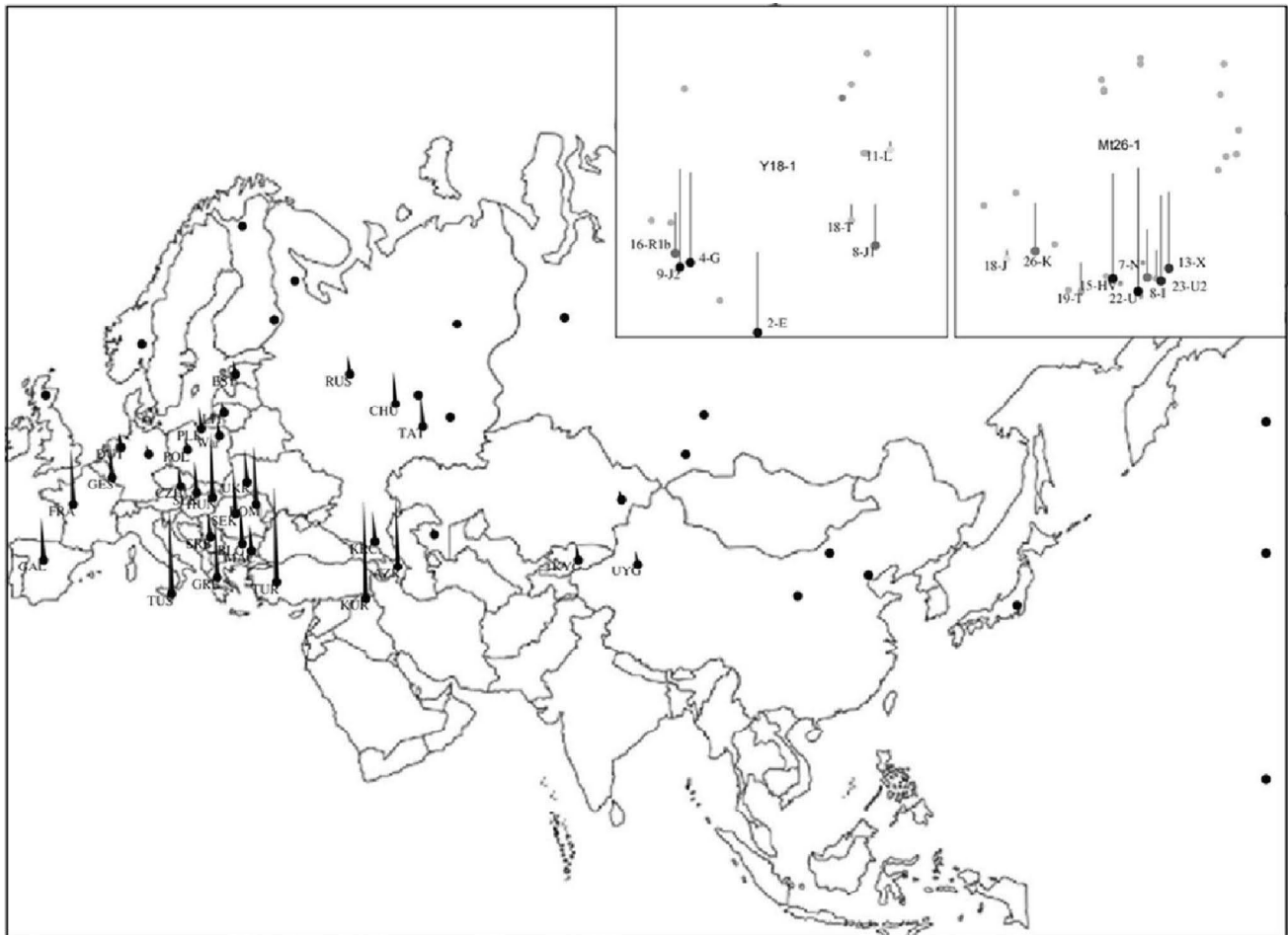
in the map “Mt26-10” means mitochondrial Hg Z having the serial number 4. The serial numbers were defined arbitrarily for the computation.) The heights of the grey columns in the Hg maps are proportional to the number of Hgs having strong correlation with the Hg belonging to the column. The black peaks in the geographical map are proportional to the IR value of the corresponding population in the given IRVT

392 that the Hg association of IRVT-10 and IRVT-1 have practi- 408  
393 cally no overlap. 409

394 The geographical distribution of IRVT-1 makes one think 410  
395 that this recent distribution may originate in the ancient 411  
396 migration of Neolithic farmers from the Crescent Fertile 412  
397 to Europe. To verify this suspicion, we accomplished the 413  
398 whole analysis on our extended mitochondrial database com- 414  
399 pleted with 34 ancient distributions. We found that the high- 415  
400 est inverse ranks of the resulting historical IRVT versions 416  
401 are attributed exactly to Neolithic populations in the Fer- 417  
402 tile Crescent and Central Europe NEO (Near Eastern Neo- 418  
403 lithic, TR, IRN, SYR, JOR, 8300–4000 BC), MEN (Middle 419  
404 East Neolithic-BrA, Middle East, 11,840–1402 BC), STR 420  
405 (Starcevo, Balkans, 5700–5500), NHU (Neolithic Hungary, 421  
406 5200–4800 BC), EMN (Early-Middle Neolithic, Europe, 422  
407 6000–3000), while recent populations KUR, and TUR have 423

also high inverse ranks. In addition, the mitochondrial Hg 408  
content of the historical Hg associations (Mt: N\*, X, HV\*, 409  
T, U\*, U2, K) is in a very good accordance with those of 410  
IRVT-1 (see Figs S3 and S4 in ESM\_4). The accumulated 411  
rate of the recent Hg association in ancient Neolithic sam- 412  
ples MEN, NEO, NHU, STR is in the range of 68%–80%. 413  
These results verify that recent distribution of IRVT-1 and 414  
the corresponding Hg association is a consequence of the 415  
migration of Neolithic farmers containing essentially the 416  
same mitochondrial Hgs in the past as the recent popula- 417  
tions living in the areas of the ancient migration. 418

According to Fig. 4, the geographical distribution of our 419  
next IRVT shows a Western European origin propagating 420  
to Central and Eastern Europe (IRVT-3). The correspond- 421  
ing Hg association (Y: I1, R1b, G, I2; Mt: H, V, J, K, U5, 422  
T) has some common elements with IRVT-1 (Y: R1b, G; 423

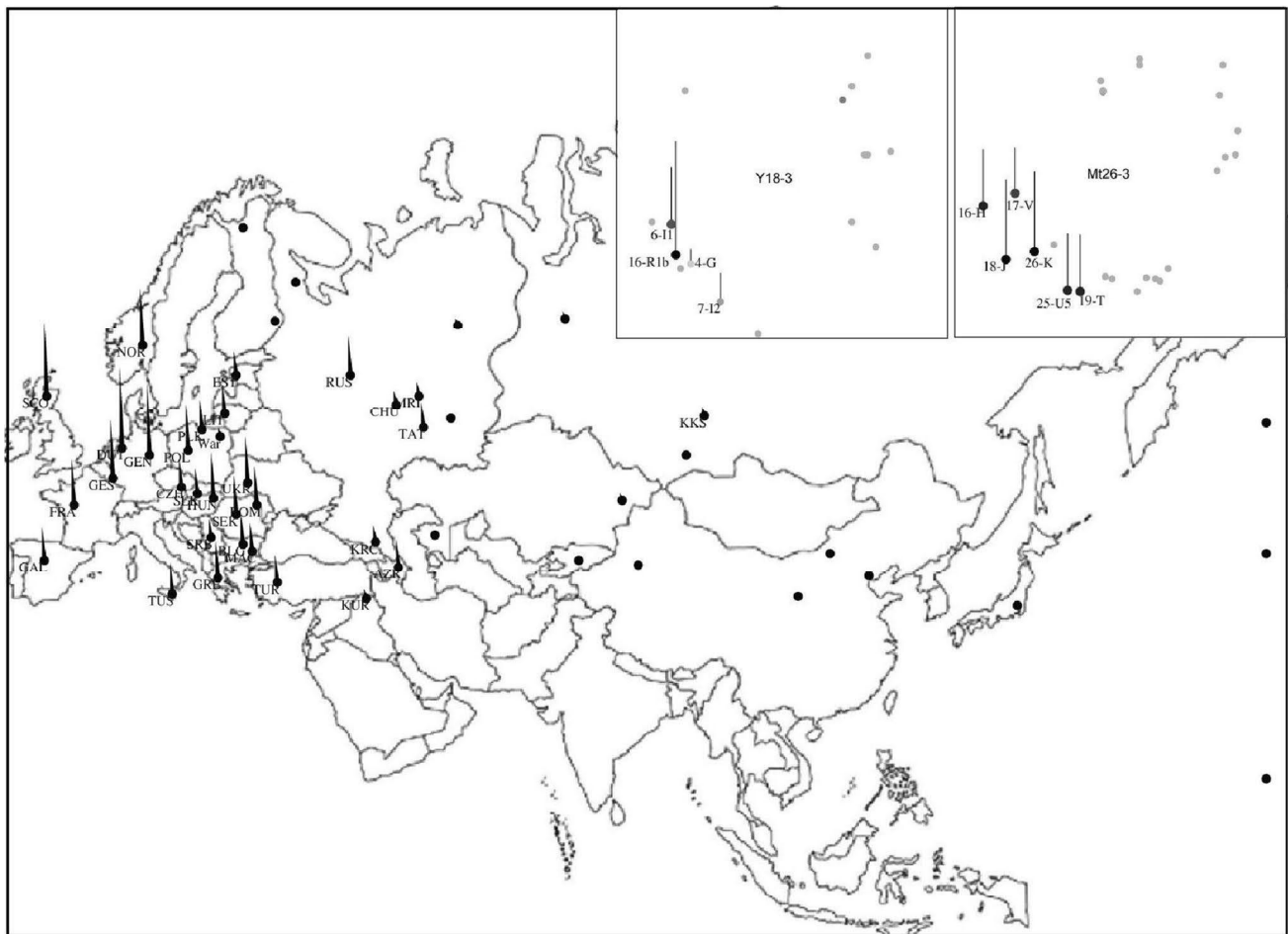


**Fig. 3** Geographical distribution of the inverse ranks of IRVT-1. The Hg association propagating from the Middle East to the Balkans and Central Europe is shown in the right upper part. For description of the symbols, see Fig. 2

424 Mt: K), but Y: I1 and the total set of the mitochondrial  
 425 Hgs except K are not found there. On the other hand, the  
 426 set (Y: J2, J1, T, L, E; Mt: HV\*, U\*, U2, N\*, X, I) in the  
 427 Hg association of IRVT-1 is not found here. The compar-  
 428 ison of the IRVTs of the extended ancient database indi-  
 429 cated the highest inverse ranks for the Neolithic European  
 430 samples STR (Starecevo, Balkans, 5700–5500) and EMN  
 431 (Early-Middle Neolithic, Europe, 6000–3000), but NHU  
 432 (Neolithic Hungary, 5200–4800 BC) and MEN (Middle  
 433 East Neolithic, Middle East, 11,840–1402 BC) have also  
 434 significant inverse ranks (see Figs S5 and S6 in ESM\_4).  
 435 Other ancient samples have zero inverse ranks for this Hg  
 436 association. The accumulated rate of the recent Hg  
 437 association in Neolithic European samples STR and EMN is  
 438 75%. Consequently, IRVT-3 may refer to the ancient Euro-  
 439 pean population preceding the migration of the farmers  
 440 from the Fertile Crescent. This is supported by the fact that  
 441 all the mitochondrial components of this Hg association  
 442 are found in the Hg distribution of EMN, taking more than  
 443 75% in sum of the whole sample.

444 The highest inverse rank values indicate Central and East-  
 445 ern Europe as the source area of the Hg association (Y: **R1a**,  
 446 I1, **I2**, **J2**, E; Mt: H, **J**, **K**, W, U\*, N\*, I, X) belonging to  
 447 IRVT-2 in Fig. 5. The mitochondrial components of this Hg  
 448 association can be divided into 2 well-defined groups—Hgs  
 449 H and J are connected to the Western European associa-  
 450 tion of IRVT-3, whereas U\*, N\* I and X are common with  
 451 IRVT-1 originating from the Fertile Crescent. The only com-  
 452 mon Hg between the two parts is Hg K. Y-chromosomal  
 453 Hgs I1 and I2 are among the most important components of  
 454 the Hg association of IRVT-3, whereas J2 and E are of high  
 455 importance in the Hg association of IRVT-1. This suggests  
 456 that this Hg association may be traced back to an admixture  
 457 of ancient Europeans and farmers arising from the Fertile  
 458 Crescent. The highest inverse rank values in the most simi-  
 459 lar historical IRVTs are assigned to neolithic samples MEN  
 460 (Iberian Neolithic, Iberia, 10,310–3160 BC) and NHU (Neo-  
 461 lithic Hungary, 5200–4800 BC). Six members of the relating  
 462 historical (ancient) Hg distribution (Mt: **J**, **K**, U4, T, U\*, I,  
 463





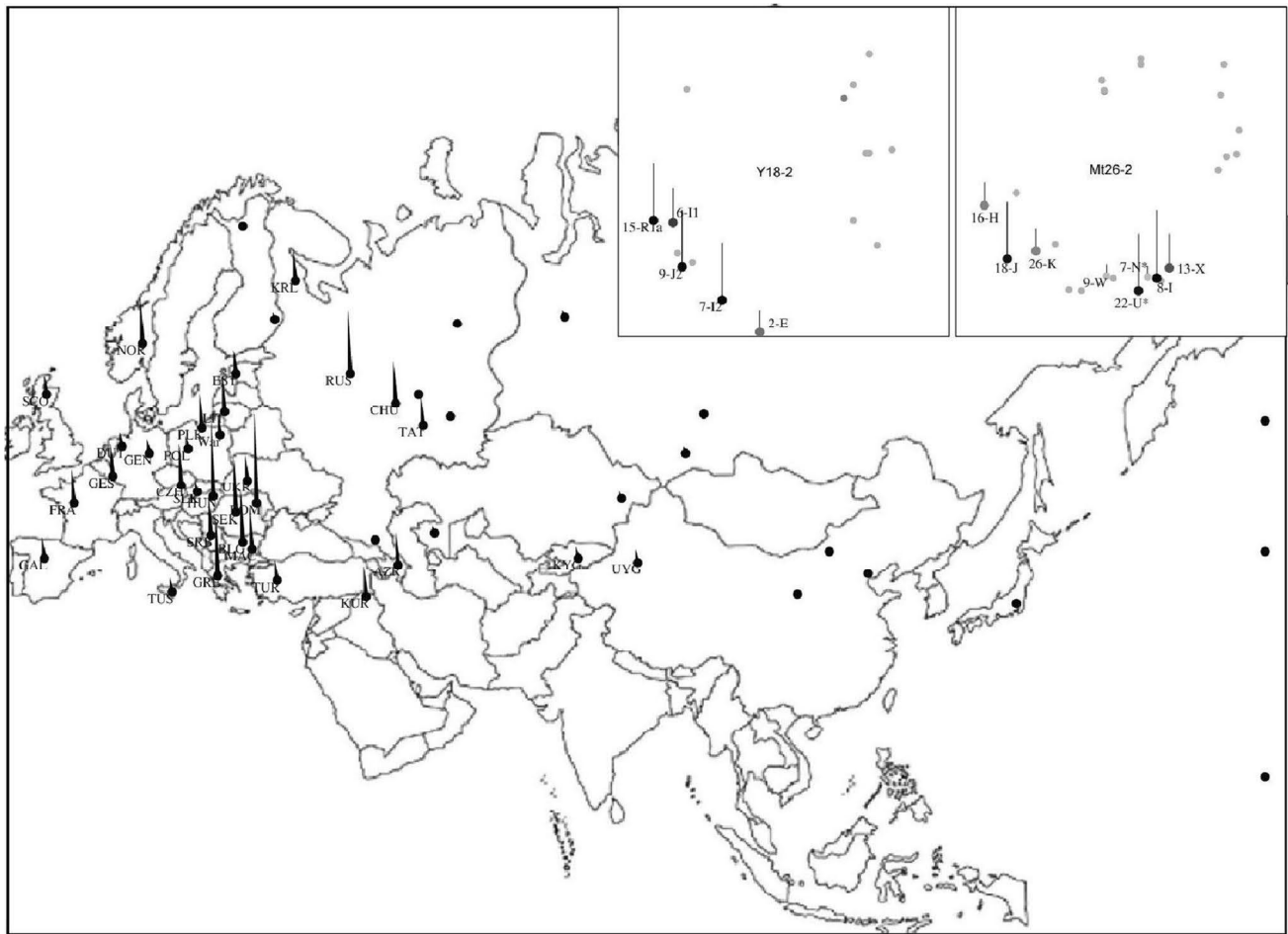
**Fig. 4** Geographical distribution of the inverse ranks of IRVT-3. The corresponding Hg association is shown in the right upper part. For description of the symbols, see Fig. 2

464 **N\***, **X**) are common with those of IRVT-2 (see Figs S7 and  
 465 S8 in ESM\_4). The accumulated rates of the recent Hg asso-  
 466 ciation in the Neolithic samples are in the range of 68–82%.  
 467 These results may really be explained by an admixture of  
 468 early Europeans and farmers in the Balkans, the Carpathian  
 469 Basin and Eastern Europe. The relatively high importance  
 470 of the resulting Hg association in the Volga region (TAT,  
 471 CHU), the ninth and tenth century Hungarian samples  
 472 (AH2, Hungarians 900–1000 AD; AH1, Ancient Hungar-  
 473 ian, 900–1000 AD) as well as HPC (pre-Conquest Hungary,  
 474 500–900 AD) may need a further explanation.

475 Figure 6 shows the geographic distribution of IRVT-4  
 476 having the largest inverse rank values in North-Eastern  
 477 Europe. The corresponding Hg association is composed by  
 478 (Y: **N**, **I1**, **R1a**, **R1b**, **I2**, **E**; Mt: **H**, **V**, **J**, **U4**, **U5**, **T**, **W**, **U\***,  
 479 **I**, **N\***, **U2**, **X**). The overlaps between this Hg association  
 480 with those of IRVT-3 and IRVT-2 are (Y: **I1**, **R1b**; Mt: **H**,  
 481 **V**, **J**, **U5**, **T**) and (**I1**, **R1a**, **I2**; Mt: **H**, **J**, **U\***, **N\***, **X**). The  
 482 Hgs being found in both overlaps (Y: **I1**; Mt: **H**, **J**) are more  
 483 important components of IRVT-3 than IRVT-2 (see the

484 column heights in the Hg-maps in Figs. 4, 5). The overlap  
 485 with the Hg association of IRVT-1 (Y: **R1b**, **E**; Mt: **T**, **U\***, **I**,  
 486 **N\***, **U2**) contains no important components in the Hg asso-  
 487 ciation of IRVT-4 (see the Hg-maps in Fig. 6).

488 The most similar historical IRVTs show high inverse rank  
 489 values for Neolithic Western European sample LNB (Late  
 490 Neolithic, Europe, 3000–1600 BC), as well as Copper-age  
 491 Eastern European YAM (Yamnaya, Afanasievo, Russia,  
 492 Ukraine, 5000–2700 BC) (see Figs S9–12 in ESM\_4). Early  
 493 medieval Viking (Norway, 780–790 AD) sample has also  
 494 high inverse rank of this Hg association. The complete set of  
 495 the components of the historical IRVTs (Mt: **H**, **J**, **U4**, **U5**, **T**,  
 496 **W**, **U\***, **I**, **X**) is very similar to that of IRVT-4. The accumu-  
 497 lated rates of the Hg association in LNB, YAM and VIK are  
 498 81, 88, and 74%. These results imply an admixture of Neo-  
 499 lithic European hunters as well as a population composed by  
 500 Neolithic hunters and farmers arising from Eastern Europe.  
 501 The geographical distribution of the inverse ranks shows that  
 502 the resulting complex Hg association has the highest weight  
 503 in Eastern and Northern Europe.



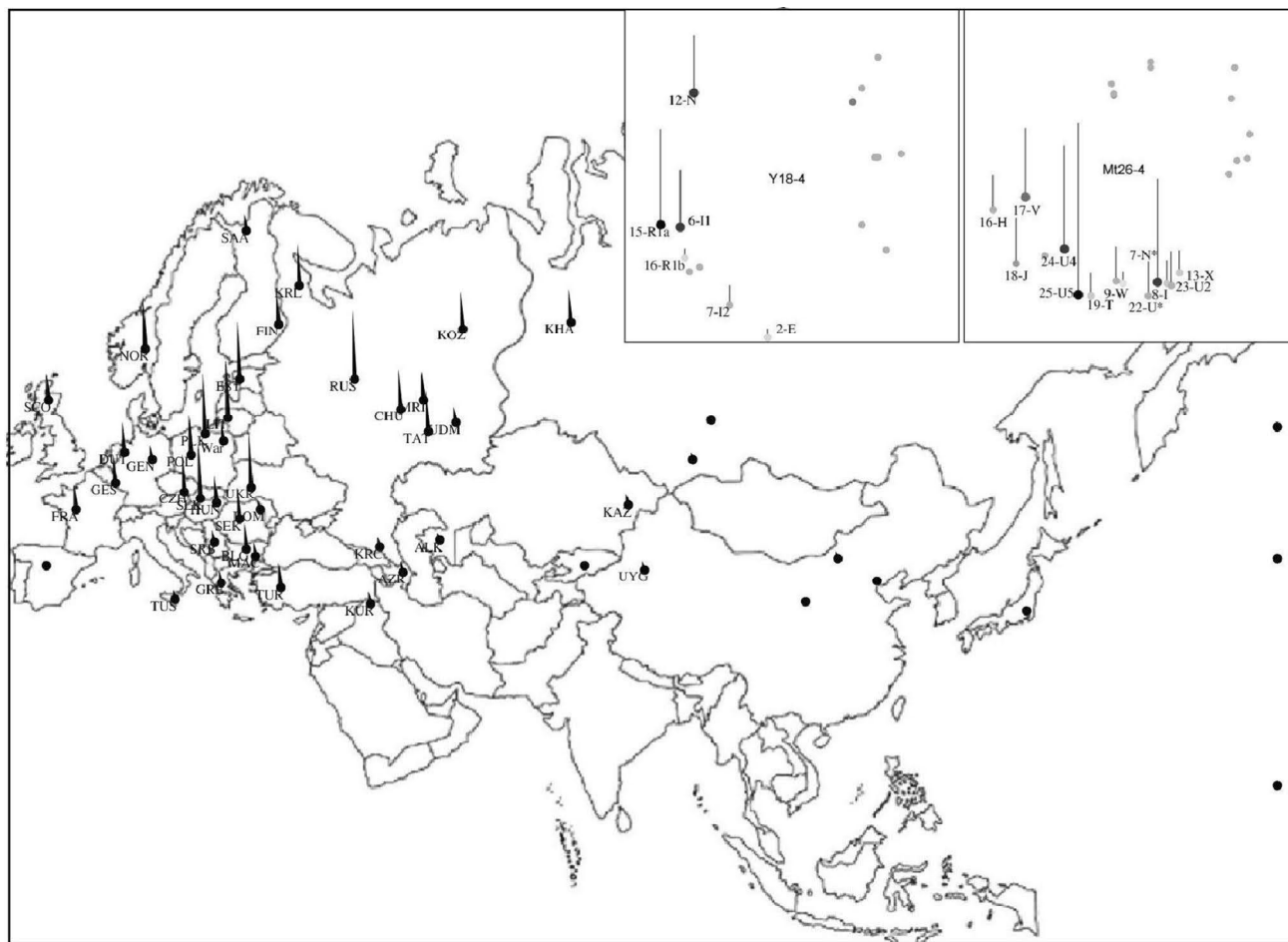
**Fig. 5** Geographical distribution of the inverse ranks of IRVT-2. The corresponding Hg association is shown in the right upper part. For description of the symbols, see Fig. 2

504 The highest inverse rank values of our last example  
 505 IRVT-5 show the Carpathian Basin and the Balkans as  
 506 source area of the Hg association (Y: R1a, R1b, J2, **G**, **I2**,  
 507 E; Mt: **H**, **U4**, J, K, T, **W**, **HV\***, U\*) (see Fig. 7). The largest  
 508 overlaps (Y: R1a, R1b, J2, **I2**, E; Mt: **H**, J, K) and (Y: R1a,  
 509 R1b, **I2**, E; Mt: **H**, **U4**, J, T, **W**) connect this Hg association  
 510 to IRVT-2 and IRVT-4. The most important Y-chromosomal  
 511 and mitochondrial Hgs G and I2 as well as H, U4 and W are  
 512 of similar importance in IRVT-2 an IRVT-4, while HV\* has  
 513 similar importance in IRVT-1.

514 For the first sight, these results imply again an admixture  
 515 of Neolithic farmers and European hunters in Central- and  
 516 Eastern Europe, like IRVT-2. However, the most similar  
 517 historical IRVTs show here the highest inverse rank values  
 518 for the West Asian BB3 and KBK, as well as early Medi-  
 519 eeval Hungarian HPC (See Figs S13-S16 in ESM\_4). The  
 520 accumulated rates of the recent Hg association are 85, 77,  
 521 and 78% in these ancient distributions, respectively. These  
 522 results imply a more complex interpretation: first, the admix-  
 523 ture of Neolithic European and Near-Eastern populations,

524 detected in IRVT-2, migrated from Eastern Europe to West-  
 525 ern Asia. (See the historical origin of the Andronovo culture.)  
 526 This may be the reason of the high cumulated rates of the  
 527 recent Hg association within Bronze-age samples BB3  
 528 (Baraba, Iron transition), West Siberia (1000–800 BC) and  
 529 KBK (Bronze Age Kurgans, Kazakhstan, 1400–1000 BC)  
 530 representing the populations of the late Andronovo culture  
 531 and Western Asia. Later, the resulting West Asian people—  
 532 Scythians, Sarmatians, Huns, Avars, Hungarians, Cumanians,  
 533 etc.—invaded Eastern Europe and the Carpathian Basin.  
 534 This is the reason, why early medieval Hungarian samples  
 535 HPC and AH2 also fit into this IRVT with significant inverse  
 536 weights.

537 We have found that the clusters belonging to the remain-  
 538 ing four IRVTs are very small and the corresponding Hg sets  
 539 contain only a few elements. This shows that these IRVTs  
 540 do not represent realistic Hg associations, but they proved to  
 541 be useful to separate outlier IRVs from the realistic clusters.  
 542 The relationships between the 10 IRVTs are shown in the  
 543 map constructed by the SOC algorithm in ESM\_3 (Fig. S3).



**Fig. 6** Geographical distribution of the inverse ranks of IRVT-4. The corresponding Hg association is shown in the right upper part. For description of the symbols, see Fig. 2

## 544 Validation

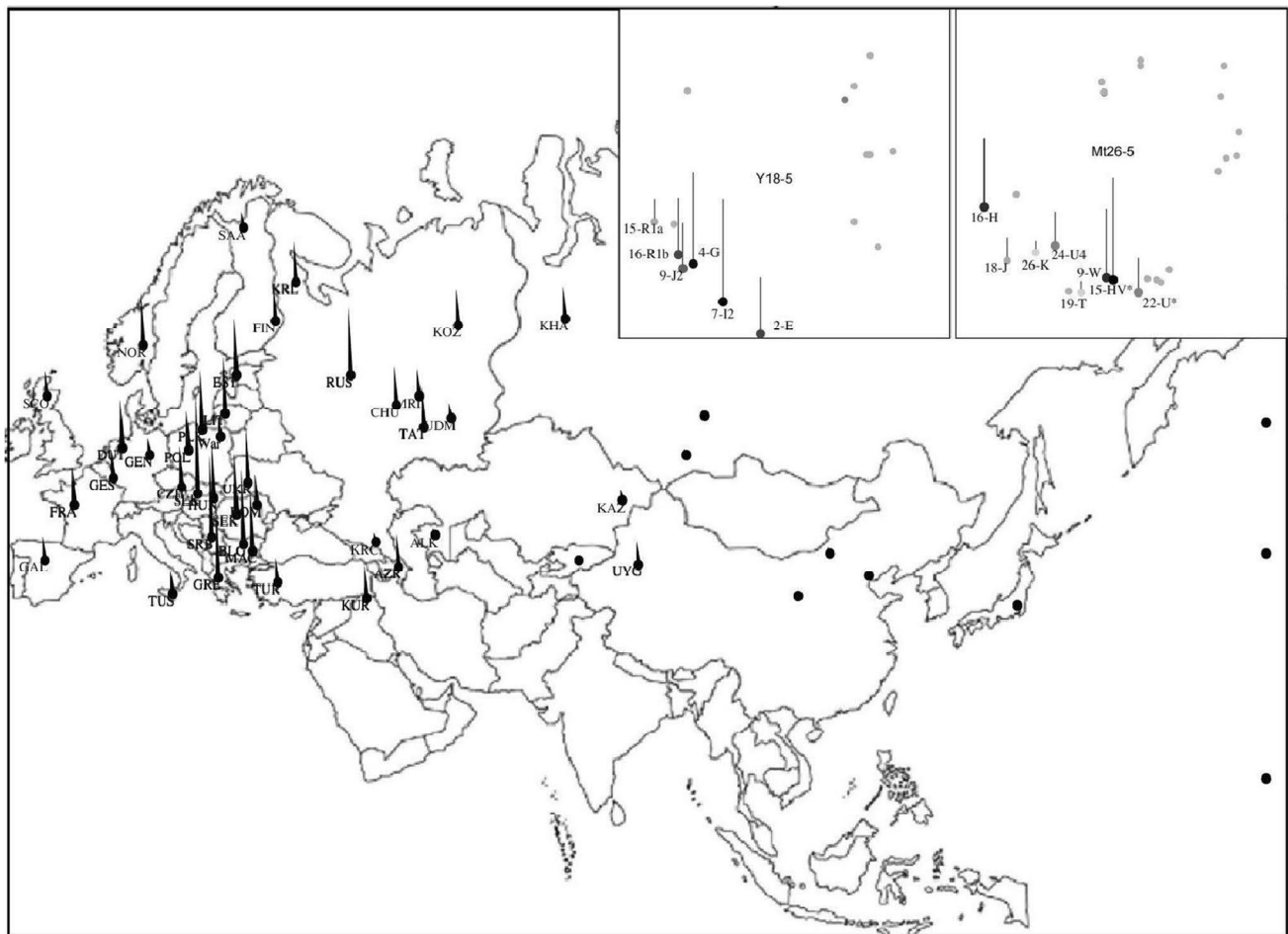
545 The standard of the goodness of the inverse rank correlation  
 546 method is the resulting correlation value itself obtained after  
 547 finishing the iteration. The highest and lowest values 1 and  
 548  $-1$  indicate totally correlating as well as totally anti-correlating  
 549 pairs of Hgs, while total un-correlation is indicated by  
 550 a correlation value of 0. Therefore, our constraints selecting  
 551 pairs of Hgs showing a correlation higher than 0.8 for a cor-  
 552 relation subset counting at least 15 populations clearly define  
 553 the requirements regarding the goodness of the iterating rank  
 554 correlation algorithm at the same time.

555 The goodness of the SOC-clustering of the IRV set fulfill-  
 556 ing the above-mentioned constraints was measured by the  
 557 correlation coefficients of the distance- and inference mat-  
 558 rices of the IRV sets.

559 To validate the method with independent data, we accom-  
 560 plished the whole process with separated mitochondrial and  
 561 Y-chromosomal Hg distributions. The results are found in  
 562 ESM\_3.

## 563 Discussion

564 Our first main contribution was to describe a new method  
 565 based on a self-learning algorithm searching for systemat-  
 566 ically jointly propagating sets of Hgs in a significant  
 567 subset of populations. The basic idea of the method is  
 568 that the inverse rank vectors of jointly propagating Hgs  
 569 are necessarily similar, so the complete set of the IRVs  
 570 belonging to Hgs having at least one strongly correlated  
 571 pair construct a clustered point system in their 50-dimen-  
 572 sional vector space. The local condensation centres of  
 573 these local condensations (IRVTs) were determined as  
 574 the learning vectors of the self-learning SOC algorithm  
 575 trained by the complete set of IRVs belonging to strongly  
 576 correlating Hgs. In addition, clustering the training IRV  
 577 set using these IRVT vectors determines the “Hg asso-  
 578 ciations” as the corresponding subsets of Hgs. Thus, this  
 579 method provides us the associations of jointly propagating  
 580 Hgs and the paths of their propagations simultaneously  
 581 and immediately.



**Fig. 7** Geographical distribution of the inverse ranks of IRVT-5. The corresponding Hg association is shown in the right upper part. For description of the symbols, see Fig. 2

582 It has been shown in previous works that a high frequency  
 583 of a Hg does not necessarily indicate its source population,  
 584 because bottleneck and founder effects may cause drastic  
 585 changes in Hg frequencies (Cinnioğlu et al. 2004; Bíró et al.  
 586 2009). However, these cases result in the loss of the correlation  
 587 with other Hgs, so our iterative rank correlation method  
 588 automatically eliminates them from the correlation subset of  
 589 populations, while the remaining subset still indicates the  
 590 real correlation.

591 Till now, Y-chromosomal and mitochondrial Hgs were  
 592 studied usually separately in human population genetics.  
 593 The novelty of our method lies in the possibility of studying  
 594 jointly propagating associations of mitochondrial and  
 595 Y-chromosomal Hgs. Moreover, joint propagation of genetic  
 596 and/or cultural (e.g. linguistic or musical) characteristics  
 597 could also be studied using IRV clustering.

598 The consideration that migrating human populations nec-  
 599 cessarily contain male and female components suggests the  
 600 idea of studying all correlations including both mitochon-  
 601 drial and Y-chromosomal Hgs. Our first example (IRVT-10,

Fig. 2) shows the geographical distribution of an IRVT hav-  
 ing the highest values in Eastern and Inner Asia as well as  
 American indigenous people, and shows a gradual reduction  
 in Western direction. This clear correlation of IRVT-10 with  
 the geographical conditions is itself an independent evidence  
 of the goodness of our method, since the geographical condi-  
 tions are totally ignored in the analysis. The Hg association  
 derived from IRVT-10 clearly contains the set of male and  
 female Hgs of well-known Eastern Asian origin (Yao et al.  
 2004; Derenko et al. 2007a, b, c; Forster et al. 1996; Kim  
 et al. 2011; Zegura et al. 2004). As the possible areas of  
 origin of the Hgs are also totally ignored from the analysis,  
 this result is a further independent evidence supporting our  
 method. In a good accordance with these results obtained  
 from recent data, we also found high inverse ranks for the  
 same mitochondrial Hg association in ancient Inner Asian  
 populations.

Similarly, good accordance between geographical distri-  
 butions of inverse ranks and places of origin of the corre-  
 sponding Hg association was experienced for IRVT1 (Fig. 3)

622 showing the Fertile Crescent as source area of the corre-  
 623 sponding Hg association propagating to Europe through  
 624 Asia Minor, the Balkans and the Carpathian Basin. As the  
 625 mitochondrial part of this Hg association was also detected  
 626 with high inverse ranks in ancient populations in the Fer-  
 627 tile Crescent, IRVT1 can be attributed to the well-known  
 628 migration of Neolithic farmers starting from the Middle East  
 629 (Bramanti et al. 2009; Malmström et al. 2009; Skoglund  
 630 et al. 2012).

631 Also, very clear correlation between the geographical  
 632 distribution and the place of origin of the corresponding  
 633 Hg association was found for IRVT-3 (Fig. 4) playing the  
 634 most important role in recent Western people. As the highest  
 635 inverse ranks of the same mitochondrial Hg association are  
 636 found in ancient European populations, IRVT-3 can be con-  
 637 sidered as the Hg association of Neolithic indigenous Euro-  
 638 peans (Gamba et al. 2014; Szécsényi-Nagy 2015; Kivisild  
 639 2017; Wong et al. 2017).

640 The three Hg associations discussed above can be consid-  
 641 ered as “pure” descendants of early populations preceding  
 642 the later admixture processes generated by the migrations  
 643 in the Neolithic period, the Bronze- and Iron ages, the late  
 644 Antiquity as well as early Middle Age.

645 The Hg associations derived from the remaining 3  
 646 IRVTs clearly mirror the admixture of the Hg associations  
 647 of IRVT-3 and IRVT-1 representing indigenous Europeans  
 648 and Neolithic Farmers. The geographical distribution of  
 649 IRVT-2 (Fig. 5) shows that the most probable stages of this  
 650 admixture were the Balkans and the Carpathian Basin. The  
 651 geographic distribution also implies the propagation of the  
 652 resulting population to Eastern Europe, in good accordance  
 653 with earlier studies of the Yamnaya culture (Anthony 2007;  
 654 Kristiansen and Larsson 2005; Kristiansen 2007; Wong et al.  
 655 2017). The high inverse ranks of the mitochondrial part of  
 656 this Hg association in early Medieval Hungarian samples  
 657 also support the Eastern European origin of this population.

658 The similarity of the geographical distribution and Hg  
 659 association of IRVT-5 to IRVT-2 intimates that the popu-  
 660 lation attributed above to IRVT-2 may play an important  
 661 role in IRVT-5, too. However, the most similar IRVTs of  
 662 ancient mitochondrial Hgs indicate the presence of the Hg  
 663 association derived from IRVT-5 (Fig. 7) in Bronze-Age  
 664 Western Asian samples, too. The explanation of this may  
 665 be the known migration of Bronze-Age Eastern Europeans  
 666 to Western Asia, and a further admixture with populations  
 667 arising from the Middle East (Antony 2007; Hanks et al.  
 668 2007) (note that the Hg associations of both IRVT-2 and  
 669 IRVT-5 contain numerous Hgs arising from the Middle East,  
 670 but these sets are not identical). The resulting population  
 671 is attributed to the Andronovo culture (Keyser et al. 2009;  
 672 Allentoft et al. 2015). It is also supported by archaeological  
 673 results that the descendants of the Andronovo culture were  
 674 found in the Eurasian Steppe after 1700–1500 BC, so the

675 high rate of the Hg distribution of IRVT-5 in early medieval  
 676 Hungarians can be traced back to the expansions of Scyth-  
 677 ian, Sarmatian, Hun, Avar, etc., empires all of them reaching  
 678 the Carpathian Basin, as well as the Hungarian conquest  
 679 (Neparáczki et al. 2017; Cynarski and Maciejewska 2016;  
 680 Szécsényi-Nagy 2015; Gamba et al. 2014; Korjakova and  
 681 Epimakhov 2014).

682 The origin of the Hg association derived from IRVT-4 can  
 683 also be traced back to the population attributed to IRVT-2  
 684 on the one hand, and early indigenous Europeans attributed  
 685 to IRVT-3 on the other hand. The clear geographical dis-  
 686 tribution of IRVT-4 (Fig. 6) shows the propagation of the  
 687 resulting complex Hg distribution of IRVT-4 from North-  
 688 Eastern Europe into Southern and Western directions. The  
 689 role of the Eastern European Yamnaya culture in the evolu-  
 690 tion of the Corded Ware culture in Northern Europe has also  
 691 been shown previously (Harrison and Heyd 2007; Vandkilde  
 692 2007; Wong et al. 2017; Allentoft et al. 2015).

693 We have shown certain ancient populations where the  
 694 cumulated rate of the mitochondrial Hg associations derived  
 695 from recent data is extremely high. In principle, the fre-  
 696 quency of an Hg in a population can be summed up by more  
 697 Hg associations constructing the given population, because  
 698 of the overlaps of their Hg contents (Zerjal et al. 2002;  
 699 Sharma et al. 2009). Therefore, these cumulated rates refer  
 700 merely to a possible maximal rate of the Hg associations,  
 701 and the actual rates may be lower. However, an extremely  
 702 high cumulated rate of a Hg association in an ancient popu-  
 703 lation may refer to a situation antedating later admixtures.

704 Unfortunately, we could not collect all ancient Y-chro-  
 705 mosomal data exactly corresponding to our ancient mito-  
 706 chondrial distributions. However, joint propagation of con-  
 707 temporary male and female haplogroups is itself a strong  
 708 validation of past human population migrations. In addition,  
 709 the validation of our method is also based on simultaneous  
 710 search for both mitochondrial and Y-chromosomal IRVTs.

711 These considerations clearly show the importance of  
 712 ancient Hg distributions in credible interpretation of the  
 713 results. The mathematically correct estimation of the rates  
 714 of the Hg associations in an actual Hg distribution and  
 715 the completion of our ancient mitochondrial data by their  
 716 Y-chromosomal counterparts would result in a much clearer  
 717 insight into the early migration processes.

718 First and last, the above discussions illustrate that our  
 719 method based on the clustering of the inverse rank vectors of  
 720 Hgs provides a good insight into the most effective migration  
 721 processes and the prehistory of the mankind. The accord-  
 722 ance with prior knowledge regarding genetic and archaeo-  
 723 logical footprints of Neolithic and Bronze-age migrations  
 724 validates our method in itself (Allentoft et al. 2015). The  
 725 results also improve that correlations of jointly propagat-  
 726 ing Hgs in contemporary populations can be traced back  
 727 to prehistoric migration processes. In addition, the method

728 could be extended to study the correlations of cultural and  
729 genetic characteristics, to validate linguistic, archaeologi-  
730 cal, ethnomusicological, etc., theories and hypotheses by  
731 genetic evidences. Other biological applications, i.e. corre-  
732 lation analysis of frequencies of different species also could  
733 reveal joint propagations of different associations of plants  
734 and/or animals.

735 **Acknowledgements** The authors are grateful to Tibor Fehér for the  
736 collection of mitochondrial and Y-chromosomal Hg distributions from  
737 prior published results. No funding was received.

### 738 Compliance with ethical standards

739 **Conflict of interest** The authors declare no conflict of interest.

740 **Research involving human participants and/or animals** All procedures  
741 performed in studies involving human participants were in accordance  
742 with the ethical standards of the institutional and/or national research  
743 committee and with the 1964 Helsinki declaration and its later amend-  
744 ments or comparable ethical standards.

### 745 References

- 746 Allentoft ME, Sikora M, Sjögren K, Rasmussen S, Rasmussen M,  
747 Stenderup J, Damgaard PB, Schroeder H, Ahlström T, Vinner L,  
748 Malaspinas A, Margaryan A, Higham T, Chivall D, Lynnerup N,  
749 Harvig L, Baron J, Della Casa P, Dabrowski P, Duffy PR, Ebel  
750 AV, Epimakhov A, Frei K, Furmanek M, Gralak T, Gromov A,  
751 Gronkiewicz S, Grupe G, Hajdu T, Jarysz R, Khartanovich V,  
752 Khokhlov A, Kiss V, Kolár J, Kriiska A, Lasak I, Longhi C, Mc-  
753 Lynn G, Merkevicius A, Merkyte I, Metspalu M, Mkrtychyan R,  
754 Moiseyev V, Paja L, Pálfi GY, Pokutta D, Pospieszny L, Price D,  
755 Saag L, Sablin M, Shishlina N, Smrčka V, Soenov VI, Szeverényi  
756 V, Tóth G, Trifanova SV, Varul L, Vicze M, Yepiskoposyan L,  
757 Zhitenev V, Orlando L, Sicheritz-Ponté T, Brunak S, Nielsen  
758 R, Kristiansen K, Eske Willerslev E (2015) Population genomics  
759 of bronze age Eurasia. *Nature* 522(7555):167–172
- 760 Anthony D (2007) The horse, the wheel and language. How bronze-  
761 age riders from the Eurasian Steppes Shaped the Modern World.  
762 Princeton Univ. Press, Princeton
- 763 Batini C, Hallast P, Vágné AJ, Zadik D, Eriksen HA, Pamjav H, Sajantila  
764 A, Wetton JH, Jobling MA (2017) Population resequencing of  
765 European mitochondrial genomes highlights sex-bias in Bronze  
766 Age demographic expansions. *Sci Rep* 7(1):12086
- 767 Bermisheva MA, Kutuev IA, Korshunova TY, Dubova NA, Villems R,  
768 Khusnutdinova EK (2004) Phylogeographic analysis of mitochon-  
769 drial dna in the Nogays: a strong mixture of maternal lineages from  
770 Eastern and Western Eurasia. *Mol Biol* 38:516–523
- 771 Bíró AZ, Zalán A, Völgyi A, Pamjav H (2009) A Y-chromosomal com-  
772 parison of the Madjars (Kazakhstan) and the Magyars (Hungary).  
773 *Am J Phys Anthropol* 139(3):305–310
- 774 Bramanti B, Thomas MG, Haak W, Unterlaender M, Jores P, Tambets  
775 K, Antanaitis-Jacobs I, Haidle MN, Jankauskas R, Kind CJ, Lueth  
776 F, Terberger T, Hiller J, Matsumura S, Forster P, Burger J (2009)  
777 Genetic discontinuity between local hunter-gatherers and Central  
778 Europe's first farmers. *Science* 326:137–140
- 779 Cinnioglu C, King R, Kivisild T, Kalfoglu E, Atasoy S, Cavalleri GL,  
780 Lillie AS, Roseman CC, Lin AA, Prince K, Oefner PJ, Shen P,  
781 Semino O, Cavalli-Sforza LL, Underhill PA (2004) Excavating  
Y-chromosome Haplotype Strata in Anatolia. *Hum Genet* 114:127–148
- Cynarski WJ, Maciejewska A (2016) The proto-Slavic warrior in  
Europe: the Scythians, Sarmatians and Lekhs. Ido movement for  
culture. *J Martial Arts Anthropol* 16(3):1–14
- Der Sarkissian C, Allentoft ME, Ávila-Arcos MC, Barnett R, Campos  
PF, Cappellini E, Ermini L, Fernández R, da Fonseca R, Ginolhac  
A, Hansen AJ, Jónsson H, Korneliusen T, Margaryan A, Martin  
MD, Moreno-Mayar JV, Raghavan M, Rasmussen M, Velasco  
MS, Schroeder H, Schubert M, Seguin-Orlando A, Wales N,  
Gilbert MT, Willerslev E, Orlando L (2015) Ancient genomics.  
*Philos Trans R Soc Lond B Biol Sci* 370(1660):20130387
- Derenko M, Malyarchuk B, Grzybowski T, Denisova G, Dambueva  
I, Perkova M, Dorzhu C, Luzina F, Lee HK, Vaneeck T, Villems  
R, Zakharov I (2007a) Phylogeographic analysis of mitochon-  
drial DNA in northern Asian populations. *Am J Hum Genet*  
81(5):1025–1041
- Derenko M, Malyarchuk B, Denisova G, Wozniak M, Grzybowski  
T, Dambueva I, Zakharov I (2007b) Y-chromosome haplo-  
group N dispersals from south Siberia to Europe. *J Hum Genet*  
52(9):763–770
- Derenko MV, Malyarchuk BA, Wozniak M, Denisova GA, Dambueva  
IK, Dorzhu CM, Grzybowski T, Zakharov IA (2007c) Distribution  
of the male lineages of Genghis Khan's descendants in northern  
Eurasian populations. *Genetika* 43(3):422–426
- Ermini L, Der Sarkissian C, Willerslev E, Orlando L (2015) Major  
transitions in human evolution revisited: a tribute to ancient DNA.  
*J Hum Evol* 79:4–20
- Forster P, Harding R, Torroni A, Bandelt HJ (1996) Origin and evolu-  
tion of Native American mtDNA variation: a reappraisal. *Am J  
Hum Genet* 59(4):935–945
- Fu Q, Posth C, Hajdinjak M, Petr M, Mallick S, Fernandes D, Furt-  
wängler A, Haak W, Meyer M, Mittnik A, Nickel B, Peltzer A,  
Rohland N, Slon V, Talamo S, Lazaridis I, Lipson M, Mathieson  
I, Schiffels S, Skoglund P, Derevianko AP, Drozdov N, Slavinsky  
V, Tsybankov A, Cremonesi RG, Mallegni F, Gély B, Vacca E,  
Morales MR, Straus LG, Neugebauer-Maresch C, Teschler-Nicola  
M, Constantin S, Moldovan OT, Benazzi S, Peresani M, Cop-  
pola D, Lari M, Ricci S, Ronchitelli A, Valentin F, Thevenet C,  
Wehrberger K, Grigorescu D, Rougier H, Crevecoeur I, Flas D,  
Semal P, Mannino MA, Cupillard C, Bocherens H, Conard NJ,  
Harvati K, Moiseyev V, Drucker DG, Svoboda J, Richards MP,  
Caramelli D, Pinhasi R, Kelso J, Patterson N, Krause J, Pääbo S,  
Reich D (2016) The genetic history of Ice Age Europe. *Nature*  
534(7606):200–205
- Gamba C, Jones ER, Teasdale MD, McLaughlin RL, Gonzalez-Fortes  
G, Mattiangeli V, Domboróczki L, Kóvári I, Pap I, Anders A,  
Whittle A, Dani J, Raczky P, Higham TFG, Hofreiter M, Bradley  
DG, Pinhasi R (2014) Genome flux and stasis in a five millennium  
transect of European prehistory. *Nat Commun* 5:5257
- Haak W, Lazaridis I, Patterson N, Rohland N, Mallick S, Llamas  
B, Brandt G, Nordenfelt S, Harney E, Stewardson K, Fu Q,  
Mittnik A, Bánffy E, Economou C, Francken M, Friederich S,  
Pena RG, Hallgren F, Khartanovich V, Khokhlov A, Kunst M,  
Kuznetsov P, Meller H, Mochalov O, Moiseyev V, Nicklisch N,  
Pichler SL, Risch R, Rojo Guerra MA, Roth C, Szécsényi-Nagy  
A, Wahl J, Meyer M, Krause J, Brown D, Anthony D, Cooper  
A, Alt KW, Reich D (2015) Massive migration from the steppe  
was a source for Indo-European languages in Europe. *Nature*  
522(7555):207–211
- Hanks BK, Epimakhov AV, Renfrew AC (2007) Towards a refined  
chronology for the Bronze Age of the southern Urals. *Russ Antiq-  
uity* 81:353–367
- Harrison R, Heyd V (2007) The Transformation of Europe in the third  
millennium BC: the example of 'Le Petit-Chasseur IIII' (Sion,  
Valais, Switzerland). *Praehistorische Zeitschrift* 82:129–214

- 848 Ilyas M, Kim JS, Cooper J, Shin YA, Kim HM, Cho YS, Hwang S, Kim  
849 H, Moon J, Chung O, Jun J, Rastogi A, Song S, Ko J, Manica A,  
850 Rahman Z, Husnain T, Bhak J (2015) Whole genome sequencing  
851 of an ethnic Pathan (Pakhtun) from the north-west of Pakistan.  
852 *BMC Genom* 16:172
- 853 Jobling MA, Tyler-Smith C (2003) The human Y chromosome: an  
854 evolutionary marker comes of age. *Nat Rev Genet* 4:598–612
- 855 Juhász Z (2007) Analysis of melody roots in Hungarian folk music  
856 using self-organizing maps with adaptively weighted dynamic  
857 time warping. *Appl Artif Intell* 21(1):35–55
- 858 Juhász Z, Fehér T, Bárány G, Zalán A, Németh E, Pádár Z, Pamjav  
859 H (2015) New clustering methods for population comparison on  
860 paternal lineages. *Mol Genet Genom* 290(2):767–784
- 861 Juhász Z, Fehér T, Németh E, Pamjav H (2016) mtDNA analysis of  
862 174 Eurasian populations using a new iterative rank correlation  
863 method. *Mol Genet Genom* 291(1):493–509
- 864 Karafet TM, Mendez FL, Meilerman MB, Underhill PA, Zegura SL,  
865 Hammer MF (2008) New binary polymorphisms reshape and  
866 increase resolution of the human Y chromosomal haplogroup tree.  
867 *Genome Res* 18:830–838
- 868 Keyser C, Bouakaze C, Crubézy E, Nikolaev VG, Montagnon D, Reis  
869 T, Ludes B (2009) Ancient DNA provides new insights into the  
870 history of south Siberian Kurgan People. *Hum Genet* 126:395–410
- 871 Kim SH, Kim KC, Shin DJ, Jin HJ, Kwak KD, Han MS, Song JM, Kim  
872 W, Kim W (2011) High frequencies of Y-chromosome haplogroup  
873 O2b-SRY465 lineages in Korea: a genetic perspective on the peo-  
874 pling of Korea. *Investig Genet* 2(1):10
- 875 Kivisild T (2017) The study of human Y chromosome variation through  
876 ancient DNA. *Hum Genet* 136(5):529–546
- 877 Korjakova L, Epimakhov AV (2014) The Urals and western Siberia in  
878 the Bronze and Iron ages. Cambridge University Press, Cambridge
- 879 Kristiansen K (2007) In the world system and the earth system. In:  
880 Hornborg B, Crumley C (eds) *Global socioenvironmental change  
881 and sustainability since the neolithic*. Left Coast Press, Walnut  
882 Creek
- 883 Kristiansen K, Larsson T (2005) The rise of bronze age society. travels,  
884 transmissions and transformations. Cambridge Univ. Press,  
885 Cambridge
- 886 Lazaridis I, Patterson N, Mittnik A, Renaud G, Mallick S, Kirsanov  
887 K, Sudmant PH, Schraiber JG, Castellano S, Lipson M, Berger B,  
888 Economou C, Bollongino R, Fu Q, Bos KI, Nordenfält S, Li H, de  
889 Filippo C, Prüfer K, Sawyer S, Posth C, Haak W, Hallgren F, For-  
890 nander E, Rohland N, Delsate D, Francken M, Guinet JM, Wahl  
891 J, Ayodo G, Babiker HA, Baillif G, Balanovska E, Balanovsky  
892 O, Barrantes R, Bedoya G, Ben-Ami H, Bene J, Berrada F, Bravi  
893 CM, Brisighelli F, Busby GB, Cali F, Churnosov M, Cole DE,  
894 Corach D, Damba L, van Driem G, Dryomov S, Dugoujon JM,  
895 Fedorova SA, Gallego Romero J, Gubina M, Hammer M, Henn  
896 BM, Hervig T, Hodoglugil U, Jha AR, Karachanak-Yankova S,  
897 Khusainova R, Khusnutdinova E, Kittles R, Kivisild T, Klitz W,  
898 Kučinskás V, Kushniarevich A, Laredj L, Litvinov S, Loukidis T,  
899 Mahley RW, Melegh B, Metspalu E, Molina J, Mountain J, Näk-  
900 käläljärvi K, Nesheva D, Nyambo T, Osipova L, Parik J, Platonov  
901 F, Posukh O, Romano V, Rothhammer F, Rudan I, Ruizbakiev R,  
902 Sahakyan H, Sajantila A, Salas A, Starikovskaya EB, Tarekegn A,  
903 Toncheva D, Turdikulova S, Uktveryte I, Utevska O, Vasquez R,  
904 Villena M, Voevoda M, Winkler CA, Yepiskoposyan L, Zalloua  
905 P, Zemunik T, Cooper A, Capelli C, Thomas MG, Ruiz-Linares  
906 A, Tishkoff SA, Singh L, Thangaraj K, Vilems R, Comas D,  
907 Sukernik R, Metspalu M, Meyer M, Eichler EE, Burger J, Slatkin  
908 M, Pääbo S, Kelso J, Reich D, Krause J (2014) Ancient human  
909 genomes suggest three ancestral populations for present-day Euro-  
910 peans. *Nature* 513(7518):409–413
- 911 Lazaridis I, Nadel D, Rollefson G, Merrett DC, Rohland N, Mallick  
912 S, Fernandes D, Novak M, Gamarra B, Sirak K, Connell S,  
913 Stewardson K, Harney E, Fu Q, Gonzalez-Forbes G, Jones ER,  
914 Roodenberg SA, Lengyel G, Bocquentin F, Gasparian B, Monge  
915 JM, Gregg M, Eshed V, Mizrahi AS, Meiklejohn C, Gerritsen  
916 F, Bejenaru L, Blüher M, Campbell A, Cavalleri G, Comas D,  
917 Froguel P, Gilbert E, Kerr SM, Kovacs P, Krause J, McGettigan  
918 D, Merrigan M, Merriwether DA, O'Reilly S, Richards MB,  
919 Semino O, Shamoon-Pour M, Stefanescu G, Stumvoll M, Tönjes  
920 A, Torroni A, Wilson JF, Yengo L, Hovhannisyann NA, Patterson  
921 N, Pinhasi R, Reich D (2016) Genomic insights into the origin  
922 of farming in the ancient Near East. *Nature* 536(7617):419–424
- 923 Lopopolo M, Børsting C, Pereira V, Morling N (2016) A study of  
924 the peopling of Greenland using next generation sequencing  
925 of complete mitochondrial genomes. *Am J Phys Anthropol*  
926 161(4):698–704
- 927 Malmström H, Gilbert MT, Thomas MG, Brandström M, Storå J,  
928 Molnar P, Andersen PK, Bendixen C, Holmlund G, Götherström  
929 A, Willerslev E (2009) Ancient DNA reveals lack of continuity  
930 between Neolithic hunter-gatherers and contemporary Scandi-  
931 navians. *Curr Biol* 19:1758–1762
- 932 Neparáczi E, Juhász Z, Pamjav H, Fehér T, Csányi B, Zink A,  
933 Maixner F, Pálfi G, Molnár E, Pap I, Kustár Á, Révész L,  
934 Raskó I, Török T (2017) Genetic structure of the early Hungar-  
935 ian conquerors inferred from mtDNA haplotypes and Y-chro-  
936 mosome haplogroups in a small cemetery. *Mol Genet Genom*  
937 292(1):201–214
- 938 Pakendorf B, Novgorodov IN, Osakovskij VL, Stoneking M (2007)  
939 Mating patterns amongst Siberian reindeer herders: inferences  
940 from mtDNA and Y-chromosomal analyses. *Am J Phys Anthro-  
941 pol* 133:1013–1027
- 942 Quintana-Murci L, Chaix R, Wells RS, Behar DM, Sayar H, Scozzari  
943 R, Rengo C, Al-Zahery N, Semino O, Santachiara-Benerecetti  
944 AS, Coppa A, Ayub Q, Mohyuddin A, Tyler-Smith C, Qasim  
945 Mehdi S, Torroni A, McElreavey K (2004) Where west meets  
946 east: the complex mtDNA landscape of the southwest and Cen-  
947 tral Asian corridor. *Am J Hum Genet* 74:827–845
- 948 Semino O, Passarino G, Oefner PJ, Lin AA, Arbuzova S, Beck-  
949 man LE, De Benedictis G, Francalacci P, Kouvatsi A, Lim-  
950 borska S, Marcikiae M, Mika A, Mika B, Primorac D, San-  
951 tachiara-Benerecetti AS, Cavalli-Sforza LL, Underhill PA  
952 (2000) The genetic legacy of Paleolithic Homo sapiens  
953 in extant Europeans: a Y chromosome perspective. *Science*  
954 290:1155–1159
- 955 Sharma S, Rai E, Sharma P, Jena M, Singh S, Darvishi K, Bhat AK,  
956 Bhanwer AJ, Tiwari PK, Bamezai RN (2009) The Indian origin  
957 of paternal haplogroup R1a1\* substantiates the autochthonous  
958 origin of Brahmans and the caste system. *J Hum Genet* 54:47–55
- 959 Simoni L, Calafell F, Pettener D, Bertranpetit J, Barbujani G (2000)  
960 Geographic patterns of mtDNA diversity in Europe. *Am J Hum  
961 Genet* 66:262–278
- 962 Skoglund P, Malmström H, Raghavan M, Storå J, Hall P, Willerslev  
963 E, Gilbert MT, Götherström A, Jakobsson M (2012) Origins  
964 and genetic legacy of Neolithic farmers and hunter-gatherers in  
965 Europe. *Science* 336(6080):466–469
- 966 Szécsényi-Nagy A (2015) Molecular genetic investigation of the Neo-  
967 lithic population history in the western Carpathian Basin. PhD  
968 Dissertation
- 969 Tambets K, Rootsi S, Kivisild T, Help H, Serk P, Loogva'li EL, Tolk  
970 HV, Reidla M, Metspalu E, Pliss L, Balanovsky O, Pshenichnov  
971 A, Balanovska E, Gubina M, Zhadanov S, Osipova L, Damba L,  
972 Voevoda M, Kutuev I, Bermisheva M, Khusnutdinova E, Gusar  
973 V, Grechanina E, Parik J, Pennarun E, Richard C, Chaventre A,  
974 Moisan JP, Bara'c L, Pericic' M, Rudan P, Terzić' R, Mikerezi I,  
975 Krumina A, Baumanis V, Koziel S, Rickards O, De Stefano GF,  
976 Anagnou N, Pappa KI, Michalodimitrakis E, Fera'k V, Fu'redi S,  
977 Komel R, Beckman L, Vilems R (2004) The western and eastern  
978 roots of the Saami—the story of genetic "outliers" told by mito-  
979 chondrial DNA and Y chromosomes. *Am J Hum Genet* 74

- 980 Underhill PA, Kivisild T (2007) Use of Y chromosome and mitochon- 992  
981 drial DNA population structure in tracing human migrations. 993  
982 *Annu Rev Genet* 41:539–564 994  
983 Vandkilde H (2007) Culture and change in the Central European Pre- 995  
984 history, 6th to 1st millennium BC. Aarhus Univ. Press, Aarhus 996  
985 Wong EHM, Khrunin A, Nichols L, Pushkarev D, Khokhrin D, Ver- 997  
986 benko D, Evgrafov O, Knowles J, Novembre J, Limborska S, 998  
987 Valouev A (2017) Reconstructing genetic history of Siberian and 999  
988 Northeastern European populations. *Genome Res* 27(1):1–14 1000  
989 Yao YG, Kong QP, Wang CY, Zhu CL, Zhang YP (2004) Different 1001  
990 matrilineal contributions to genetic structure of ethnic groups in 1002  
991 the silk road region in China. *Mol Biol Evol* 21(12):2265–2280 1003
- Zegura SL, Karafet TM, Zhivotovsky LA, Hammer MF (2004) High- 992  
resolution SNPs and microsatellite haplotypes point to a single, 993  
recent entry of Native American Y chromosomes into the Ameri- 994  
cas. *Mol Biol Evol* 21(1):164–175 995  
Zerjal T, Wells RS, Yuldasheva N, Ruzibakiev R, Tyler-Smith C (2002) 996  
A genetic landscape reshaped by recent events: Y-chromosomal 997  
insights into central Asia. *Am J Hum Genet* 71(3):466–482 1000

UNCORRECTED PROOF



Journal:	<b>438</b>
Article:	<b>1469</b>

## Author Query Form

**Please ensure you fill out your response to the queries raised below and return this form along with your corrections**

Dear Author

During the process of typesetting your article, the following queries have arisen. Please check your typeset proof carefully against the queries listed below and mark the necessary changes either directly on the proof/online grid or in the 'Author's response' area provided below

Query	Details Required	Author's Response
<a href="#">AQ1</a>	Author: Since this article follows "UK English" kindly check and clarify if the word "medieval" can be changed to "mediaeval" globally.	
<a href="#">AQ2</a>	Author: Kindly check and provide page range for Tambets et al. (2004).	
<a href="#">AQ3</a>	Author: Fig1, Fig2, Fig3, Fig4, Fig5, Fig6, Fig7 - Figure is poor in quality as its labels are not readable. Please supply a new version of the said figure with legible labels preferably in .eps, .tiff or .jpeg format with 600 dpi resolution.	