

1 **A forensic population database in El Salvador: 58 STRs and 94 SNPs**

2 Ferran Casals^{1,2}, Raquel Rasal¹, Roger Anglada¹, Marc Tormo^{1,3}, Núria Bonet¹, Nury

3 Rivas⁴, Patricia Vázquez^{5*}, Francesc Calafell^{6*}

4 ¹Genomics Core Facility, Departament de Ciències Experimentals i de la Salut, Universitat
5 Pompeu Fabra, Parc de Recerca Biomèdica de Barcelona, 08003 Barcelona, Catalonia, Spain

6 ² Departament de Genètica, Microbiologia i Estadística, Universitat de Barcelona, Barcelona,
7 Spain

8 ³ Scientific IT Core Facility, Departament de Ciències Experimentals i de la Salut, Universitat
9 Pompeu Fabra, Parc de Recerca Biomèdica de Barcelona, 08003 Barcelona, Catalonia, Spain

10 ⁴ Instituto de Medicina Legal Dr. Roberto Masferrer, San Salvador, El Salvador.

11 ⁵ Asociación Pro-Búsqueda de Niñas y Niños Desaparecidos de El Salvador, 27 calle Pnte.
12 No.1329 Colonia Layco, San Salvador, El Salvador

13 ⁶Institute of Evolutionary Biology (UPF-CSIC), Department of Experimental and Health Sciences,
14 Universitat Pompeu Fabra, Barcelona, Spain

15 * Corresponding authors

16

17 **Abstract**

18 We have genotyped the 58 STRs (27 autosomal, 24 Y-STRs and 7 X-STRs) and 94 autosomal
19 SNPs in Illumina ForenSeq™ Primer Mix A in a sample of 248 men and 143 women from El
20 Salvador, Central America. Regional division (Centro, Oriente, Occidente) showed in almost all
21 cases F_{ST} values not significantly different from 0, and further analyses were applied only to the
22 undivided, country-wide population. The overall random match probability (RMP) decreased
23 from 6.79×10^{-31} in length-based genotypes in the 27 autosomal STRs to 1.47×10^{-34} in repeat-
24 sequence based genotypes. Combining the autosomal loci in this set, RMP reaches 2.97×10^{-70} .
25 In a population genetic analysis, El Salvador showed the lowest F_{ST} values with US Hispanics
26 both for autosomal and X-STRs; however, it was much closer to Native Americans for the latter
27 than for the former, in accordance with the well-known gender-biased admixture that created
28 most Latin American populations.

29

30 **Keywords:** Massive Parallel Sequencing; Repeat Sequence-Based Alleles; Missing persons

31

32 Introduction

33 El Salvador is the smallest nation (~21,000 Km²) in Central America and the only country in the
34 Central American isthmus that has no coast in the Caribbean. Its population has been
35 estimated for 2019 at 6,704,864 people, 61.7% of whom live in urban areas [1].

36 Administratively, the country is divided in three regions (Centro, Oriente and Occidente,
37 meaning respectively Center, East, and West), further subdivided into 14 departments. About
38 76% of the population is concentrated in 5 departments, namely, San Salvador, La Libertad,
39 Santa Ana, Sonsonate, and San Miguel). Three main indigenous groups have inhabited the
40 Salvadoran territory: the Nahua-Pipiles, the Lencas and the Kakawira (Cacaopera) [2].

41 According to the last census, which was conducted in 2007, 83% of the Salvadoran population
42 self-identify as Mestizo (which is a common term throughout Latin America to describe the
43 product of historic admixture mostly between Native Americans and Europeans), followed by
44 European (15%), Afro-descendant (0.13%) and only 0.23% as Native American, although the
45 latter figure could be an underestimate [3].

46

47 El Salvador, like other countries in the region, has suffered repression and human rights
48 violations since colonial times. Social injustice persisted and was a major trigger of the 1980-
49 1992 civil war. In the aftermath of war, poverty and gang violence increased, which were a
50 major factor in inducing emigration, mostly to the USA, but also to other countries. These
51 historical events, as well as the current situation, have generated a need for genetic
52 identification: the Civil War produced large numbers of unidentified casualties and missing
53 persons, both adults and forcibly adopted children. After the peace agreements, and in the
54 absence of a response from the Government, in 1994 the relatives of these missing children
55 formed the Pro-Búsqueda Association (<http://www.probusqueda.org.sv/>) with the support of
56 the Jesuit priest Jon Cortina; Pro-Búsqueda manages a database of genetic profiles of
57 relatives and young people already found; in addition, in recent years it has begun including
58 young people looking for their families. On the other hand, the number of missing migrants at
59 the U.S.-Mexico border (and which fraction of those are Salvadorans) is unknown. In 2010,
60 government entities, relatives of deceased and missing migrants, and the Argentine forensic
61 anthropology team (EAAF) promoted a forensic data bank of unidentified migrants from El
62 Salvador, which is part of the Proyecto Frontera, a regional mechanism for the forensic
63 exchange between unidentified remains and missing persons along the Central America-
64 Mexico-United States of America migration corridor
65 (<https://bancoforenseelsalvador.org/quienes-somos/>). Obviously, identification of missing
66 people requires the availability of the allelic frequencies of the genetic markers used in
67 casework. Previous reports of allele or haplotype frequencies in loci of forensic interest
68 include autosomal STRs [4–9], X-STRs [10], Y-STRs [11–14], and mtDNA sequences [15].

69 Recently, several platforms have become available to apply massive parallel sequencing (MPS)
70 to the genotyping of STRs and SNPs in a forensic context. MPS offers the possibility of
71 multiplexing a much larger number of markers than capillary electrophoresis, and of
72 combining SNPs and STRs. Large numbers of markers assure that, even if results cannot be
73 obtained for part of the markers due to DNA degradation in the sample, still the remaining loci

74 can provide a likelihood ratio that can be clearly interpreted as indicating or rejecting a match.
75 And a large number of markers may be needed to resolve cases involving distant relatives as
76 references [16], as it is often the case in the identification of missing persons decades after
77 their deaths.

78 Sequencing rather than sizing STRs allows extracting more information from most of the STRs.
79 In particular, the Verogen ForenSeq™ Primer Mix A (Verogen, San Diego, CA) ,which contains
80 58 STRs and 94 SNPs, together with the with the Universal Analysis Software (UAS) provided by
81 the manufacturer yields for each successful genotype two different types of information: a
82 length-based (LB) genotype, in accordance with the numeric genotype that sizing by capillary
83 electrophoresis would have yielded, and a repeat-sequence based genotype (RSB), that is, the
84 sequence haplotype of the repeat region of each STR (the flanking region sequence is not
85 available through the UAS). A number of studies [17–20], among others, have shown that,
86 while the overall informativeness (as measured by a priori statistics) increases moderately
87 from LB to RSB genotypes, the number of different alleles and the number of rare alleles
88 register more substantial increases.

89 We hereby report the allele frequencies in 391 samples of El Salvador of the Verogen
90 ForenSeq™ Primer Mix A loci, as a resource in the quest for the identification of the missing
91 persons in the country.

92 **Methods**

93 DNA was obtained from either buccal cells or saliva for 402 samples from the general
94 population of El Salvador, after appropriate written informed consent. One sample was
95 excluded because it showed a number of STRs with three and four alleles, probably indicating
96 contamination. Ten additional samples were also removed because they were detected as
97 first- or second-degree relatives of other samples in the database. Thus, the final sample size
98 was 391 individuals (248 men and 143 women), subdivided into 196 samples from the Centro
99 region, 90 from Occidente and 105 from Oriente.

100 **DNA was extracted from saliva or FTA cards.** For saliva samples we used a standard organic
101 method with proteinase K digestion, followed by phenol-chloroform extraction. DNA was
102 extracted from FTA cards, using 8 punches (1.2mm each) per sample and the PrepFiler BTA kit
103 (ThermoFisher Scientific, Waltham MA, USA). Few changes were made from the original
104 protocol (https://assets.thermofisher.com/TFS-Assets/LSG/manuals/cms_099065.pdf) such as
105 the elution buffer volume which was 40ul instead of 50ul and the incubation time that was
106 increased to 20min. DNA was quantified using Qubit (ThermoFisher Scientific) according to
107 manufacturer recommendations. This project was reviewed and approved by the National
108 Committee for Health Research, El Salvador (reference CNEIS/2018/030), on July 24th, 2018.

109 Samples were sequenced for the Verogen ForenSeq™ Primer Mix A loci according to the
110 manufacturer's protocol. Sample volume for amplification and subsequent library preparation
111 was 5 ml, at a DNA concentration of 0.2 ng/ul. The pooled libraries were sequenced in a
112 351×31 cycles run with the MiSeq FGx™ instrument (Illumina, San Diego, CA, USA) following
113 the supplier's protocol. We performed six sequencing runs in a standard flow cell, with 22, 88,

114 96, 96, 94 and 96 samples, and one run in a micro flow cell with 16 samples, plus the
115 manufacturer-supplied positive and negative controls in each run.

116 STR allele sequences were retrieved from the report generated by the Forenseq UAS interface
117 and inspected by means of an in-house R script (IFator for autosomal STRs, YIFator for Y-STRs,
118 and XIFator for X-STRs, available from github <https://github.com/fcalafell/>) [17]. These scripts
119 allow uncovering much more sequence diversity than that reported by the Forenseq UAS
120 interface, which only highlights sequence variants when they are found in isometric
121 heterozygotes, that is, in individuals carrying two alleles of the same length but different
122 sequence. Note that the Forenseq UAS provides exclusively the repeat region sequence, and
123 thus all of our subsequent analyses are based on the repeat region sequence (as processed
124 with our scripts) and cannot include the flanking region. We used a shorthand notation for
125 sequence-based alleles which was based on the called number of repeats plus a lower-case
126 letter indicating an approximation to the repeat structure, consistently with [17]. For instance,
127 STR D3S1358 has the general structure TCTA [TCTG]_x [TCTA]_y; length is given by 1 + x + y,
128 which we supplement with *a* if x = 1, *b* if x = 2, *c* if x = 3 or *d* if x = 4. Thus, allele TCTA [TCTG]₁
129 [TCTA]₁₃ is denoted 15a. See [17] for further details. The full list of RSB variants and their
130 notation can be found in Supplementary Table 1. Allele length information was taken directly
131 from the Forenseq UAS. Hardy-Weinberg equilibrium (HWE) and F_{ST} , were computed with
132 Arlequin 3.5 [21]. We computed two a priori informativeness statistics: power of
133 discrimination [22] and the chance of excluding a putative father in a paternity trio if the
134 mother is known [23], by direct calculation from allele frequencies using MS Excel.
135 Information in the F_{ST} distance matrix was extracted and plotted by means of multidimensional
136 scaling (MDS), which was computed with the isoMDS function in the MASS library in R. Y
137 haplogroups (i.e., the main branches of the Y-SNP tree) were predicted from Y-STR haplotypes
138 using the nevgen (<http://nevgen.org>) Bayesian predictor, and adapting the nomenclature of
139 the resulting haplogroups to that suggested by [24] and used in <http://yhrd.org>. For SNP
140 analyses, the analytical threshold (e.g., the lower limit of detection) was set to 0 and the
141 interpretation (allele calling) threshold was set to 2.5% [17], to avoid false negative calls.

142

143

144 **Results**

145 Allele frequency and forensic informativeness data were generated for 27 autosomal STRs, 7 X-
146 STRs, 24 Y-STRs and 94 SNPs in a sample of 248 men and 143 women from the general
147 population of El Salvador. Samples were divided according to the region of origin within the
148 country (Centro, Occidente and Oriente). However, after applying the Bonferroni correction
149 for multiple testing taking into account the number of loci in each category (autosomal STRs,
150 X-STRs, Y-STRs, and SNPs), F_{ST} values among the regions were not statistically significantly
151 different from 0 ($p > 0.05$) for all but one locus (DYS390, $F_{ST} = 0.0514$) (see Supplementary Tables
152 2-5, and 9), and thus, for all subsequent analyses, the El Salvador sample has been treated as a
153 single entity.

154 Allele frequencies, heterozygosities, Hardy-Weinberg p-values and a priori informativeness
155 statistics for 27 autosomal STRs in a population sample of 391 individuals from El Salvador are
156 presented in Supplementary Tables 2 and 3, respectively for LB and RSB alleles. Sample sizes
157 for these analyses ranged from 488 to 782, with an average of 758 and a median of 758
158 chromosomes. We could generate genotypes for all individuals in the sample for 18 of the 27
159 STRs. In five additional loci, missing genotypes were less than 5% of the total population
160 sample. The most problematic loci were PentaE, with 37.6% missing genotypes, and PentaD,
161 with 24.0%. LB genotypes were in Hardy-Weinberg equilibrium ($p > 0.05$) at all but four STRs
162 (D22S1045, D5S818, PentaD, and PentaE) after Bonferroni correction. Random match
163 probability (RMP) was 6.79×10^{-31} , which was of the same order of magnitude of that in
164 Catalans [17] or US Hispanics [18]. The joint chance of paternity exclusion was $1 - (2.8 \times 10^{-11})$.
165 RSB variation was detected in 20 out of 27 autosomal STRs, for a total of 490 RSB alleles, up
166 from 293 LB alleles. Exactly the same four STRs mentioned above (namely, D22S1045, D5S818,
167 PentaD, and PentaE) also failed HWE after Bonferroni correction for RSB genotypes. RMP
168 decreased to 1.47×10^{-34} , that is 4,629 times lower than the LB-based RMP. Rare alleles (defined
169 here arbitrarily as those with a frequency $< 1\%$) can have an important contribution to solving
170 cases involving distant relatives [16]. We found 83 LB rare alleles with 178 (45.5%) individuals
171 carrying at least one, and up to four rare alleles across all autosomal STRs. These figures clearly
172 increased for RSB: there were a total of 237 RSB rare alleles, with 306 (78.3%) individuals
173 carrying at least one rare allele.

174 Allele frequencies, heterozygosities, Hardy-Weinberg test results and F_{ST} values for X-STRs are
175 presented in Supplementary Tables 4 and 5 for LB and RSB alleles. A volunteer was a
176 heterozygote for three of the seven X-STRs, yet genotypes were also recovered for 22 out of
177 24 Y-STRs. Additionally, AMELX and AMELY were sequenced in this individual at 116x and 68x
178 coverages, respectively. Contamination seems to be ruled out by the fact that, in 27 autosomal
179 STRs, 5 show no amplification, in 6 only one allele is detected, and for 16 STRs, two alleles are
180 detected, with no autosomal STR showing more than two alleles. Thus, these results fit the
181 expected pattern of a XXY karyotype, which also agrees with the fact that the volunteer self-
182 identifies as a male. Even though the sample was comprised of 142 women and 249 men, the
183 maximum number of X chromosomes in the population sample was 534 rather than 533
184 ($= 2 \times 142 + 249$). Allele frequencies, expected heterozygosities and F_{ST} values were estimated
185 from total chromosome samples ranging from 323 to 534, with an average 489 and a median
186 532 chromosomes; Hardy-Weinberg tests were performed on the female samples, which
187 ranged from 96 to 143, with an average of 132 and a median of 143. Genotypes could be
188 generated for all samples at 3 out of 7 loci, and in a fourth locus, only one sample had a
189 missing genotype. On the contrary, missing genotype rates were 6.6% for DXS10135, 12.3% for
190 DXS8378, and 41.9% for DXS10103. Missing rates were higher for men (47.2%) than for women
191 (41.9%) in DXS10103, while the reverse was true for DXS8378 (14.0% in women and 11.3% in
192 men) and DXS10135 (6.9% in women and 6.5% in men). Variation in the repeat sequence was
193 present in three X-STR loci, and, adding up across all seven loci, 82 different LB and 129 RSB
194 alleles were found. Considering that the STRs within the pairs DXS10135-DXS8378, DXS7132-
195 DXS10074, and DXS10103-HPRTB are in close proximity of each other, we also estimated
196 haplotype frequencies by direct counting in males and informative (i.e., heterozygote in at
197 most one locus within a particular pair) females (Supplementary Tables 6 and 7). Sample sizes

198 for these haplotypes were respectively 324, 346, and 237 chromosomes. The availability of
199 repeat sequences implied that the number of different haplotypes increased from 173 LB
200 haplotypes to 227 RSB haplotypes.

201 Sample sizes for the 24 Y-STRs present in the Verogen ForenSeq™ primer mix A ranged from
202 134 to 249, with an average of 223 and a median of 241. For seven loci, all individuals could be
203 genotyped, and, in an additional four loci, missing values were <5%. On the contrary, DYS392
204 had 46.2% missing genotype calls, DYS389II reached 34.1%, and the value for DYS448 was
205 29.3%. It should also be noted that DYS438 produced genotypes in two individuals, at
206 coverages 52X and 71X, in which no other Y-STR could be genotyped and that were
207 heterozygotes for 7 and 5 out of 7 X-STRs. Presumably, these are XX persons with a spurious
208 amplification of DYS438. The average and median number of missing genotypes per individual
209 were 2.56 and 1, respectively, and, for 76 individuals, we could produce complete haplotypes.
210 We could detect RSB alleles in 11 Y-STRs; overall, the number of alleles increased from 244 LB
211 alleles to 367 RSB alleles. However, RSB variation did not imply an increase in the number of
212 haplotypes (Supplementary Table 8), since the 76 males (out of a population sample of 249
213 men) for which we could generate complete haplotypes all carried different LB haplotypes. The
214 average F_{ST} value among the three Salvadoran regions was low (0.0021), and it was not
215 significantly different from zero ($p>0.05$) in all but one locus, namely DYS390, with $F_{ST}=0.0631$
216 ($p=0.0001$). In Table 1, we report the frequencies of the predicted haplogroups for this subset.
217 In particular, we found haplogroup E1b1a (5.3%), which is much more frequent in African
218 populations than elsewhere, as well as haplogroup Q1a2-M346 (13.2%), which is found almost
219 exclusively in Native Americans. The rest of haplogroups and their frequencies, once the
220 frequencies of these two components are subtracted, are similar to those found in the Iberian
221 Peninsula [25,26]. Thus, these results can be interpreted as the paternal lineages of El Salvador
222 being admixed from African, Native American, and European sources, with the latter in a
223 greater proportion.

224 Allele frequencies, a priori statistics, HWE and F_{ST} values for the 94 autosomal identification
225 SNPs in Verogen ForenSeq™ in a population sample of El Salvador are given in Supplementary
226 Table 9. Sample sizes ranged from 392 to 782 chromosomes, with a mean of 754 and a median
227 of 782. For 53 loci, genotypes could be produced for all individuals, and overall 80 SNPs had a
228 fraction of missing genotypes <5%. On the contrary, four loci had >30% missing genotypes:
229 rs1736442 (49.87%), rs2920816 (47.83%), rs7041158 (39.9%), and rs1031825 (31.97%). Eight
230 SNPs failed HWE after Bonferroni correction. Average expected heterozygosity was 0.4406,
231 which is close to the maximum possible value of 0.5. RMP was 3.13×10^{-38} , and the chance of
232 excluding a false father is $1 - 2.28 \times 10^{-8}$. When combining the autosomal STRs and SNPs in
233 Verogen ForenSeq™, these a priori statistics take very low values: RMP becomes 4.60×10^{-72} ,
234 and the chance of excluding a false father is $1 - 1.06 \times 10^{-20}$.

235 We next computed F_{ST} distances between a set of reference populations: Roma and Catalans
236 from Spain [17] (the former is an ethnic minority, while the latter were sampled as individuals
237 with all four grandparents born in Catalonia), the main ethnic groups in the USA [18], and
238 world populations grouped by continental origins [19]. Table 2 shows the F_{ST} distance matrix
239 computed from the 27 autosomal STRs; in it, Salvadorans appear closest to US Hispanics ($F_{ST}=$
240 0.0042), while they are more distant from Catalans, Europeans and European Americans ($F_{ST} \approx$

241 0.0200). Intriguingly, the population of El Salvador is slightly closer to European Americans
242 than to Europeans or Catalans from the colonial power, Spain; a possible explanation might be
243 that European Americans are known to carry ~0.18% admixture from Native Americans [27].
244 Although their distance to Native Americans is larger ($F_{ST} = 0.0262$), it should be noted that
245 Salvadorans are the population Native Americans are closest to. We applied MDS to this
246 distance matrix, and the result can be seen in Figure 1a. In the plot, El Salvador groups with
247 Hispanics, as well as with East Asians and Asian Americans, even though their F_{ST} values with
248 the latter two populations are relatively high (0.0298-0.0375), and actually higher than the
249 distance between El Salvador and European populations (see above and Table 2). This cluster
250 is relatively close to a European set of populations. The patterns observed with X-STRs are
251 slightly different: while El Salvador is still closest to Hispanics ($F_{ST} = 0.0031$), it is much closer to
252 Native Americans ($F_{ST} = 0.0061$) than to populations of European descent ($F_{ST} = 0.0242 -$
253 0.0323); as in autosomal STRs, it is closer to European Americans than to populations from
254 Europe. The MDS plot (Figure 1b) reproduces the same general groupings of populations, but
255 Salvadorans, Hispanics and East Asians are closer to Native Americans than they were in the
256 autosomal STR-based plot (Figure 1a).

257

258 Discussion

259 We have typed the 58 STRs and 94 SNPs contained in the Verogen Forenseq™ primer Mix A in
260 391 samples from El Salvador, and provide allele frequencies for the general population of El
261 Salvador, which will be of invaluable help in the quest for the thousands of people that were
262 disappeared during the El Salvador Civil War (1980-1992), and for identifying human remains
263 of putative migrants to the USA. In particular, the large number of loci contained and the
264 degree of informativity added by sequencing rather than sizing alleles make it particularly
265 adequate when DNA degradation results in a high proportion of missing genotypes. Out of 27
266 autosomal STRs, we found four (D22S1045, D5S818, PentaD, and PentaE) that failed HWE after
267 Bonferroni correction, in all cases due to a homozygote excess over the expected values under
268 HWE. Technical and/or population cases could have caused this homozygote excess. However,
269 Novroski et al. [18] found 1-3 loci failing HWE in the main ethnic groups in the USA, which are
270 presumed to be large and relatively homogeneous populations. In particular, D5S818 also
271 failed HWE in African Americans, Hispanics and European Americans (although only below the
272 Bonferroni threshold in the latter), as did PentaD in African Americans. This would point to
273 technical causes in heterozygote detection in at least these two STRs. In particular, to date,
274 five sequence variations at the primer binding region of D5S818 have been reported to cause
275 discrepancies in paternity testing since they cause null alleles [28]. And although rarer, null
276 alleles have also been reported for PentaD [29].

277 As expected given the cultural and ethnic homogeneity of this relatively small country (21,041
278 Km², roughly the size of Slovenia, Israel or New Jersey), we found that allele frequencies are
279 not significantly different among the three main regions in El Salvador (Centro, Occidente, and
280 Oriente). In El Salvador, population mobility has always been high, both seasonally (of laborers
281 to the agricultural areas) and permanently (to the capital or abroad) [30]. The war increased
282 migration to the capital, and, while some migrants returned to their home towns, many settled

283 in San Salvador [31]. Still currently, internal migration remains high due to economic or social
284 reasons, such as escaping high-crime areas; the small size of the country implies that it can be
285 covered in a single, countrywide migration network.

286 Also as expected given that ~90% of the population of El Salvador are Mestizos (the product of
287 historic admixture mostly between Native Americans and Europeans), allele frequencies in El
288 Salvador are closest to those in Hispanics. We should note that we sampled blindly with
289 respect to ethnicity and did not record it; we expect that our sample reflects the average
290 genetic composition of El Salvador. It is noteworthy that, in the case of X-STRs, Salvadorans,
291 Hispanics and other Central Americans [10] are especially close to Native Americans. This is
292 likely the result of the well-known sex-specific admixture patterns in the Americas, where
293 Native American ancestry was contributed mostly by females, while most European migrants
294 were male [32–34]. In the case of El Salvador, Native American mtDNA sequences were found
295 at a ~95% frequency [15], while the Native American Q haplogroup was found at a 31%
296 frequency [12], slightly higher than our 13% STR-based estimate. Since we estimated
297 haplogroup frequencies from Y-STRs rather than from the Y-SNPs that define them, our
298 haplogroup frequencies should be taken with caution[35]. Still, the presence of inferred Q and
299 E1b1a Y chromosomes in our population sample highlights the internal diversity of Y-
300 chromosome haplotypes in El Salvador, which must be taken into account in casework.

301 It should be noted that we based our population genetics analysis on the standard F_{ST} metric,
302 which does not take into account the mutational distance between alleles. Alternatively, R_{ST}
303 [36] is based on a quantitative variance apportionment of allele size, and it reflects better the
304 evolutionary history of size-based alleles. One would expect that, besides the effects of genetic
305 drift, founder effects and gene flow that can be expected in an admixed population such that
306 of El Salvador, and that can be captured by F_{ST} , R_{ST} would add the mutational history, that is,
307 the fact that presumably two populations with smaller allele size differences are more closely
308 related to each other than populations with larger allele size differences. However, massive
309 parallel sequencing of STR alleles has uncovered different layers of complexity in STR structure,
310 comprising different repeat arrays, single nucleotide mutations or repeat conversions, which
311 imply that isometric alleles cannot be considered as single evolutionary entities. Thus, it would
312 be desirable to implement some ad hoc metric of evolutionary distance among RSB alleles,
313 which would probably need to be locus-specific. Instead, as a basic phenetic distance, we have
314 used F_{ST} , which, as detailed in the paragraph above, has allowed us to retrieve the expected
315 population genetic patterns of Salvadorans, and which may imply that the time scale of
316 differentiation by mutation is deeper than that caused by drift and admixture in mixed
317 American populations.

318

319 ACKNOWLEDGEMENTS

320 This article is dedicated to the memory of the late Cristián Orrego Benavente, who had this
321 initiative, and for his general contribution to the defense of human rights in El Salvador.

322

323 We would like to particularly thank all the volunteers participating in this study. We are
324 particularly grateful to (now deceased) for this initiative and for. This work was supported by
325 the Spanish Ministry of Economy and Competitiveness and Agencia Estatal de Investigación
326 (grant numbers CGL2016-75389-P (MINEICO/FEDER, UE), PID2019-106485GB-
327 I00/AEI/10.13039/501100011033 (MINEICO), and “Unidad María de Maeztu” (MDM-2014-
328 0370) to FCal; Agència de Gestió d’Ajuts Universitaris i de la Recerca (Generalitat de Catalunya,
329 grant 2017SGR00702); Agència Catalana de Cooperació al Desenvolupament (
330 ACCD004/17/00019 and ACCD016/18/00031); Fundación Panamericana para el Desarrollo
331 (PADF, No. PRDHD-RFA-R-2017-009). We thank also the Ministry of Health of El Salvador,
332 which, in 2018, allowed us to take samples at their facilities.

333

334

335

336 REFERENCES

- 337 [1] DIGESTYC, Encuesta de Hogares de Hogares Múltiples, Ministerio de Economía, San
338 Salvador, 2020.
- 339 [2] J. Lemus, Sociolinguistic Atlas of Indigenous Peoples in Latin America, UNICEF and
340 FUNPROEIB, Cochabamba, Bolivia, 2009.
- 341 [3] Centre for the Autonomy and Development of Indigenous Peoples updated by IFAD,
342 Country technical note on indigenous peoples' issues, Republic of El Salvador, San
343 Salvador, 2017.
- 344 [4] P. Muñoz, E.L. Pinto de Erazo, C. Baeza, E. Arroyo-Pardo, A.M. López-Parra, Genetic
345 polymorphism of 15 STR loci in El Salvador, *Int. J. Legal Med.* 129 (2015) 991–993.
346 <https://doi.org/10.1007/s00414-015-1148-8>.
- 347 [5] J.C. Monterrosa, J. Morales, I. Yurrebaso, O. García, Population genetic data for 16 STR
348 loci (PowerPlex ESX-17 kit) in El Salvador, *Forensic Sci. Int. Genet.* 6 (2012).
349 <https://doi.org/10.1016/j.fsigen.2011.12.004>.
- 350 [6] J. Lovo-Gómez, A. Salas, Á. Carracedo, Microsatellite autosomal genotyping data in four
351 indigenous populations from El Salvador, *Forensic Sci. Int.* 170 (2007) 86–91.
352 <https://doi.org/10.1016/j.forsciint.2006.05.031>.
- 353 [7] J.C. Monterrosa, J.A. Morales, O. García, Genetic variation for 15 short tandem repeat
354 loci in an El Salvadoran (Central America) population, *J. Forensic Sci.* 51 (2006) 451–452.
355 <https://doi.org/10.1111/j.1556-4029.2006.00097.x>.
- 356 [8] B. Martínez-Jarreta, P. Vásquez, E. Abecia, M. Garde, I. de Blás, B. Budowle, Autosomic
357 STR Loci (HUMTPOX, HUMTH01, HUMVWA, D18S535, D1S1656 and D12S391) in San
358 Salvador (El Salvador, Central America), *J. Forensic Sci.* 49 (2004) 1–2.
359 <https://doi.org/10.1520/jfs2003395>.
- 360 [9] J.A. Morales, J.C. Monterrosa, J.C. Alvarez, C. Entrala, J.A. Lorente, M. Lorente, B.
361 Budowle, E. Villanueva, Population Data on Nine STR Loci in an El Salvadoran (Central
362 American) Sample Population, *J. Forensic Sci.* 47 (2002) 15461J.
363 <https://doi.org/10.1520/jfs15461j>.
- 364 [10] M. Baeta, E. Prieto-Fernández, C. Núñez, T. Kleinbielen, P. Villaescusa, L. Palencia-
365 Madrid, O. Alvarez-Gila, B. Martínez-Jarreta, M.M. de Pancorbo, Study of 17 X-STRs in
366 Native American and Mestizo populations of Central America for forensic and
367 population purposes, *Int. J. Legal Med.* (2021). <https://doi.org/10.1007/s00414-021-02536-9>.
- 369 [11] J.C. Monterrosa, J.A. Morales, I. Yurrebaso, L. Gusmão, O. García, Population data for 12
370 Y-chromosome STR loci in a sample from El Salvador, *Leg. Med.* 12 (2010) 46–51.
371 <https://doi.org/10.1016/j.legalmed.2009.10.003>.
- 372 [12] J. Lovo-Gómez, A. Blanco-Verea, M. V. Lareu, M. Brión, A. Carracedo, The genetic male
373 legacy from El Salvador, *Forensic Sci. Int.* 171 (2007) 198–203.
374 <https://doi.org/10.1016/j.forsciint.2006.07.005>.
- 375 [13] B. Martínez-Jarreta, P. Vásquez, E. Abecia, B. Budowle, A. Luna, F. Peiró,
376 Characterization of 17 Y-STR Loci in a Population from El Salvador (San Salvador, Central
377 America) and Their Potential for DNA Profiling, *J. Forensic Sci.* 50 (2005) 1–4.

- 378 <https://doi.org/10.1520/jfs2005173>.
- 379 [14] J. Saul, M. Fondevila, A. Salas, M. Brión, M.V. Lareu, Á. Carracedo, Y-chromosome STR-
380 haplotype typing in El Salvador, *Forensic Sci. Int.* 142 (2004) 45–49.
381 <https://doi.org/10.1016/j.forsciint.2004.02.004>.
- 382 [15] A. Salas, J. Lovo-Gómez, V. Álvarez-Iglesias, M. Cerezo, M.V. Lareu, V. Macaulay, M.B.
383 Richards, Á. Carracedo, Mitochondrial echoes of first settlement and genetic continuity
384 in El Salvador, *PLoS One.* 4 (2009) e6882.
385 <https://doi.org/10.1371/journal.pone.0006882>.
- 386 [16] F. Calafell, R. Anglada, N. Bonet, M. González-Ruiz, G. Prats-Muñoz, R. Rasal, C. Lalueza-
387 Fox, J. Bertranpetit, A. Malgosa, F. Casals, An assessment of a massively parallel
388 sequencing approach for the identification of individuals from mass graves of the
389 Spanish Civil War (1936–1939), *Electrophoresis.* 37 (2016).
390 <https://doi.org/10.1002/elps.201600180>.
- 391 [17] F. Casals, R. Anglada, N. Bonet, R. Rasal, K.J. van der Gaag, J. Hoogenboom, N. Solé-
392 Morata, D. Comas, F. Calafell, Length and repeat-sequence variation in 58 STRs and 94
393 SNPs in two Spanish populations, *Forensic Sci. Int. Genet.* 30 (2017).
394 <https://doi.org/10.1016/j.fsigen.2017.06.006>.
- 395 [18] N.M.M. Novroski, J.L. King, J.D. Churchill, L.H. Seah, B. Budowle, Characterization of
396 genetic sequence variation of 58 STR loci in four major population groups, *Forensic Sci.*
397 *Int. Genet.* 25 (2016) 214–226. <https://doi.org/10.1016/j.fsigen.2016.09.007>.
- 398 [19] C. Phillips, L. Devesse, D. Ballard, L. van Weert, M. de la Puente, S. Melis, V. Álvarez
399 Iglesias, A. Freire-Aradas, N. Oldroyd, C. Holt, D. Syndercombe Court, Á. Carracedo,
400 M.V. Lareu, Global patterns of STR sequence variation: Sequencing the CEPH human
401 genome diversity panel for 58 forensic STRs using the Illumina ForenSeq DNA Signature
402 Prep Kit, *Electrophoresis.* 39 (2018) 2708–2724.
403 <https://doi.org/10.1002/elps.201800117>.
- 404 [20] F.R. Wendt, J.D. Churchill, N.M.M. Novroski, J.L. King, J. Ng, R.F. Oldt, K.L. McCulloh, J.A.
405 Weise, D.G. Smith, S. Kanthaswamy, B. Budowle, Genetic analysis of the Yavapai Native
406 Americans from West-Central Arizona using the Illumina MiSeq FGx™ forensic
407 genomics system, *Forensic Sci. Int. Genet.* 24 (2016) 18–23.
408 <https://doi.org/10.1016/j.fsigen.2016.05.008>.
- 409 [21] L. Excoffier, H.E.L. Lischer, Arlequin suite ver 3.5: a new series of programs to perform
410 population genetics analyses under Linux and Windows, *Mol. Ecol. Resour.* 10 (2010)
411 564–567.
- 412 [22] R.A. Fisher, Standard calculations for evaluating a blood system, *Heredity (Edinb.)* 5
413 (1951) 95–102.
- 414 [23] P.E. Smouse, R. Chakraborty, The use of restriction fragment length polymorphisms in
415 paternity analysis, *Am. J. Hum. Genet.* 38 (1986) 918–939.
416 <https://pubmed.ncbi.nlm.nih.gov/3014872/> (accessed November 5, 2021).
- 417 [24] M. van Oven, A. Van Geystelen, M. Kayser, R. Decorte, M.H.D. Larmuseau, Seeing the
418 wood for the trees: a minimal reference phylogeny for the human Y chromosome.,
419 *Hum. Mutat.* 35 (2014) 187–91. <https://doi.org/10.1002/humu.22468>.
- 420 [25] S.M.M. Adams, E. Bosch, P.L.L. Balaesque, S.J.J. Ballereau, A.C.C. Lee, E. Arroyo, A.M.

- 421 López-Parra, M. Aler, M.S.G.S. Grifo, M. Brion, A. Carracedo, J. Lavinha, B. Martínez-
422 Jarreta, L. Quintana-Murci, A. Picornell, M. Ramon, K. Skorecki, D.M.M. Behar, F.
423 Calafell, M.A.A. Jobling, A.M. Lopez-Parra, M. Aler, M.S.G.S. Grifo, M. Brion, A.
424 Carracedo, J. Lavinha, B. Martinez-Jarreta, L. Quintana-Murci, A. Picornell, M. Ramon, K.
425 Skorecki, D.M.M. Behar, F. Calafell, M.A.A. Jobling, The genetic legacy of religious
426 diversity and intolerance: paternal lineages of Christians, Jews, and Muslims in the
427 Iberian Peninsula, *Am. J. Hum. Genet.* 83 (2008) 725–736.
428 <https://doi.org/10.1016/j.ajhg.2008.11.007>.
- 429 [26] N. Solé-Morata, J. Bertranpetit, D. Comas, F. Calafell, Y-chromosome diversity in Catalan
430 surname samples: insights into surname origin and frequency., *Eur. J. Hum. Genet.* 23
431 (2015) 1549–57. <https://doi.org/10.1038/ejhg.2015.14>.
- 432 [27] K. Bryc, E.Y. Durand, J.M. Macpherson, D. Reich, J.L. Mountain, The genetic ancestry of
433 African Americans, Latinos, and European Americans across the United States, *Am. J.*
434 *Hum. Genet.* 96 (2015) 37–53. <https://doi.org/10.1016/J.AJHG.2014.11.010>.
- 435 [28] C. Shao, Y. Yao, X. Pan, M. Wu, B. Zhang, H. Xu, J. Xie, K. Sun, Variants in linkage status
436 at D5S818 detected by multiple STR kits comparison and Sanger sequencing, *Mol.*
437 *Genet. Genomic Med.* 9 (2021). <https://doi.org/10.1002/MGG3.1765>.
- 438 [29] K.J. Van Der Gaag, R.H. De Leeuw, J. Hoogenboom, J. Patel, D.R. Storts, J.F.J. Laros, P. De
439 Knijff, Massively parallel sequencing of short tandem repeats—Population data and
440 mixture analysis results for the PowerSeq™ system, *Forensic Sci. Int. Genet.* 24 (2016)
441 86–96. <https://doi.org/10.1016/J.FSIGEN.2016.05.016>.
- 442 [30] S. Montes, Displaced persons and Salvadoran refugees, *Int. Relations. Natl. Univ. Costa*
443 *Rica.* 13 (1985) 11–21.
- 444 [31] J.D. Morán Mendoza, Guerra y migración interna en El Salvador, in: P.C. de Población
445 (Ed.), *Semin. Int. Población Del Istmo Al Final Del Milen.*, San José, Costa Rica, 1999: pp.
446 307–333.
- 447 [32] I. Mendizabal, K. Sandoval, G. Berniell-Lee, F. Calafell, A. Salas, A. Martinez-Fuentes, D.
448 Comas, A. Martínez-Fuentes, D. Comas, Genetic origin, admixture, and asymmetry in
449 maternal and paternal human lineages in Cuba, *BMC Evol Biol.* 8 (2008) 213.
450 <https://doi.org/10.1186/1471-2148-8-213>.
- 451 [33] L. Ongaro, M.O. Scliar, R. Flores, A. Raveane, D. Marnetto, S. Sarno, G.A. Gneccchi-
452 Ruscone, M.E. Alarcón-Riquelme, E. Patin, P. Wangkumhang, G. Hellenthal, M.
453 Gonzalez-Santos, R.J. King, A. Kouvatsi, O. Balanovsky, E. Balanovska, L. Atramentova, S.
454 Turdikulova, S. Mastana, D. Marjanovic, L. Mulahasanovic, A. Leskovac, M.F. Lima-
455 Costa, A.C. Pereira, M.L. Barreto, B.L. Horta, N. Mabunda, C.A. May, A. Moreno-Estrada,
456 A. Achilli, A. Olivieri, O. Semino, K. Tambets, T. Kivisild, D. Luiselli, A. Torroni, C. Capelli,
457 E. Tarazona-Santos, M. Metspalu, L. Pagani, F. Montinaro, The Genomic Impact of
458 European Colonization of the Americas, *Curr. Biol.* 29 (2019) 3974–3986.e4.
459 <https://doi.org/10.1016/j.cub.2019.09.076>.
- 460 [34] A. Moreno-Estrada, S. Gravel, F. Zakharia, J.L. McCauley, J.K. Byrnes, C.R. Gignoux, P.A.
461 Ortiz-Tello, R.J. Martínez, D.J. Hedges, R.W. Morris, C. Eng, K. Sandoval, S. Acevedo-
462 Acevedo, P.J. Norman, Z. Layrisse, P. Parham, J.C. Martínez-Cruzado, E.G. Burchard,
463 M.L. Cuccaro, E.R. Martin, C.D. Bustamante, Reconstructing the Population Genetic
464 History of the Caribbean, *PLoS Genet.* 9 (2013).
465 <https://doi.org/10.1371/journal.pgen.1003925>.

466 [35] J. Jannuzzi, J. Ribeiro, C. Alho, G. de Oliveira Lázaro e Arão, R. Cicarelli, H. Simões Dutra
467 Corrêa, S. Ferreira, C. Fridman, V. Gomes, S. Loiola, M.F. da Mota, Â. Ribeiro-dos-
468 Santos, C.A. de Souza, R.S. de Sousa Azulay, E.F. Carvalho, L. Gusmão, Male lineages in
469 Brazilian populations and performance of haplogroup prediction tools, *Forensic Sci. Int.*
470 *Genet.* 44 (2020). <https://doi.org/10.1016/J.FSIGEN.2019.102163>.

471 [36] M. Slatkin, A measure of population subdivision based on microsatellite allele
472 frequencies, *Genetics* 139 (1995) 457–462.

473

474

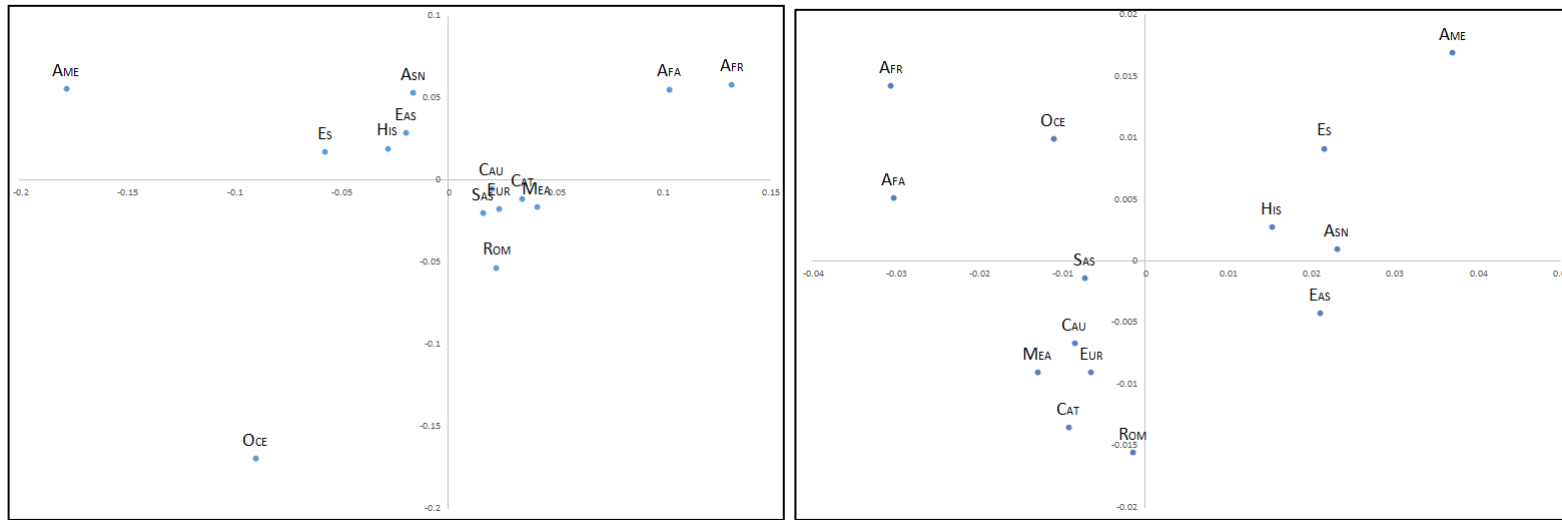
Haplogroup	N	Relative freq. (%)
E1b1a	4	5.3
E1b1b-M35	1	1.3
E1b1b-M78	9	11.8
E1b1b-V22	1	1.3
E1b1b-M123	2	2.6
E1b1b-M81	2	2.6
G2a2b-M406	1	1.3
I1-M253	1	1.3
I2-M26	3	3.9
I2-L460	1	1.3
I2-L596	1	1.3
J1-P58	1	1.3
J2a-L26	1	1.3
J2a-M319	1	1.3
J2a-L25	1	1.3
J2a-L581	1	1.3
Q1a2-M346	10	13.2
R1a	2	2.6
R1b	30	39.5
T	2	2.6
Unknown	1	1.3

475

476 Table 1. Relative haplogroup frequencies (*Relative freq.*) predicted from Y-STRs using the
477 Bayesian predictor nevgen [www.nevgen.org] for the 76 complete Y-STR haplotypes in samples
478 from El Salvador. In the case labelled “Unknown”, no predicted haplotype reached a posterior
479 probability > 80%

	ES	CAT	ROM	AFA	AFR	ASN	EAS	CAU	EUR	MEA	HIS	AME	OCE	SAS
ES	0	0.0323	0.0278	0.0389	0.0388	0.0109	0.0121	0.0242	0.0265	0.0310	0.0031	0.0061	0.0315	0.0230
CAT	0.0212	0	0.0102	0.0234	0.0319	0.0322	0.0271	0.0027	0.0022	0.0032	0.0202	0.0500	0.0188	0.0075
ROM	0.0250	0.0146	0	0.0271	0.0327	0.0279	0.0213	0.0090	0.0094	0.0136	0.0173	0.0471	0.0126	0.0111
AFA	0.0338	0.0253	0.0313	0	0.0044	0.0449	0.0435	0.0175	0.0213	0.0195	0.0337	0.0605	0.0191	0.0176
AFR	0.0509	0.0399	0.0455	0.0083	0	0.0454	0.0447	0.0197	0.0260	0.0255	0.0383	0.0619	0.0187	0.0213
ASN	0.0298	0.0268	0.0288	0.0287	0.0423	0	0.0007	0.0251	0.0261	0.0313	0.0100	0.0147	0.0239	0.0215
EAS	0.0375	0.0314	0.0351	0.0369	0.0388	0.0090	0	0.0222	0.0197	0.0260	0.0081	0.0200	0.0237	0.0196
CAU	0.0197	0.0012	0.0141	0.0242	0.0372	0.0231	0.0277	0	0.0024	0.0060	0.0150	0.0423	0.0122	0.0032
EUR	0.0219	0.0021	0.0151	0.0264	0.0355	0.0275	0.0256	0.0029	0	0.0028	0.0163	0.0438	0.0128	0.0059
MEA	0.0267	0.0079	0.0142	0.0208	0.0279	0.0244	0.0244	0.0085	0.0047	0	0.0210	0.0489	0.0175	0.0050
HIS	0.0042	0.0156	0.0200	0.0274	0.0433	0.0223	0.0320	0.0119	0.0155	0.0198	0	0.0121	0.0293	0.0163
AME	0.0262	0.0644	0.0600	0.0710	0.0759	0.0613	0.0529	0.0580	0.0573	0.0618	0.0317	0	0.0489	0.0377
OCE	0.0499	0.0639	0.0611	0.0670	0.0678	0.0639	0.0570	0.0590	0.0523	0.0562	0.0552	0.0675	0	0.0095
SAS	0.0248	0.0122	0.0145	0.0244	0.0311	0.0189	0.0156	0.0107	0.0072	0.0076	0.0197	0.0544	0.0459	0

Table 2. F_{ST} distance matrices based on 27 autosomal STRs (below diagonal) or 7 X-STRs (above diagonal). ES, El Salvador (present data). CAT, Catalans; ROM, Spanish Roma [17]. AFA, African Americans; ASN, Asian Americans; CAU, *Caucasians*; HIS, Hispanics [18]. AFR, Africans; EAS, East Asians; EUR, Europeans; MEA, Middle Easterners and North Africans; AME, Native Americans; OCE, Oceanians; SAS, South Asians [19].



a)

b)

Figure 1. Multidimensional scaling (MDS) plots based on F_{ST} distances from a) 27 autosomal STRs and b) 7 X-STRs. as in Table 2. Stress values are a) 12.1%, b) 7.0%. In MDS, stress values measure the degree of concordance between the original distance (F_{ST} in this case) and the distances between points as they appear in the plot; stress values closer to 0 signify a better concordance. Population abbreviations: ES, El Salvador; CAT, Catalans; ROM, Spanish Roma; AFA, African Americans; ASN, Asian Americans; CAU, *Caucasians*; HIS, Hispanics; AFR, Africans; EAS, East Asians; EUR, Europeans; MEA, Middle Easterners and North Africans; AME, Native Americans; OCE, Oceanians; SAS, South Asians. Note the position of ES closer to Hispanics, and in b), to Native Americans as well.

Supplementary Table 1. Repeat sequence based allele structure and nomenclature for the STRs in the Verogen Forenseq™ Primer Mix A. The number and sequence of the repetitive units in each allele are indicated. LB allele: length-based allele names; RSB allele: repeat-sequence based allele names; sequence: repeat-region sequence region, as provided by the Verogen Forenseq™ UAS.

Supplementary Table 2. Length-based absolute (N) and relative (rel) allele frequencies, observed heterozygosity (Obs. Het.), expected heterozygosity (Exp. Het.), HWE p-value (p HWE), power of discrimination (POD), chance of exclusion in a paternity trio (CE), F_{ST} and the p-value for F_{ST} (p Fst) in 27 autosomal STRs in a population sample of 391 individuals from El Salvador. The number of individuals for which we could obtain a genotype is indicated next to each locus name.

Supplementary Table 3. Repeat sequence-based absolute (N) and relative (rel) allele frequencies, , observed heterozygosity (Obs. Het.), expected heterozygosity (Exp. Het.), HWE p-value (p HWE), power of discrimination (POD), chance of exclusion in a paternity trio (CE), F_{ST} and the p-value for F_{ST} (p Fst) in 27 autosomal STRs in a population sample of 391 individuals from El Salvador.. See Supplementary Table 1 for allele sequences. The number of individuals for which we could obtain a genotype is indicated next to each locus name.

Supplementary Table 4. Length-based allele absolute (N) and relative (rel) frequencies, observed heterozygosity (Obs. Het.) , expected heterozygosity (Exp. Het.) (both in females), HWE p-value (p HWE), F_{ST} and the p-value for F_{ST} (p Fst) in 7 X-STRs, in a population sample of 534 X chromosomes from El Salvador The number of individuals for which we could obtain a genotype is indicated next to each locus name.

Supplementary Table 5. Repeat sequence-based absolute (N) and relative (rel) allele frequencies, observed heterozygosity (Obs. Het.) , expected heterozygosity (Exp. Het.) (both in females), HWE p-value (p HWE), F_{ST} and the p-value for F_{ST} (p Fst) in 7 X-STRs, in a population sample of 534 X chromosomes from El Salvador. See Supplementary Table 1 for allele sequences. The number of individuals for which we could obtain a genotype is indicated next to each locus name.

Supplementary Table 6. Length-based haplotype frequencies for the DXS10135-DXS8378, DXS7132-DXS10074, and DXS10103-HPRTB pairs. Haplotype frequencies were estimated by direct counting in males and informative (i.e, heterozygote in at most one locus within a particular pair) females. Exp. Het.: expected heterozygosity.

Supplementary Table 7. Repeat sequence-based haplotype frequencies for the DXS10135-DXS8378, DXS7132-DXS10074, and DXS10103-HPRTB pairs. See Supplementary Table 1 for allele sequences. Haplotype frequencies were estimated by direct counting in males and informative (i.e, heterozygote in at most one locus within a particular pair) females. Exp. Het.: expected heterozygosity.

Supplementary Table 8. Repeat sequence-based Y-STR haplotype frequencies, for the 76 Y chromosomes from El Salvador for which we could genotype the complete set of Y-STRs in the

Verogen Forenseq™ Primer Mix A . See Supplementary Table 1 for allele sequences. Abs. freq.: absolute frequency; Rel. freq.: relative frequency.

Supplementary Table 9. Allele frequencies, observed heterozygosity (Obs. Het.), expected heterozygosity (Exp. Het.), HWE p-value (p HWE), power of discrimination (PD), chance of exclusion in a paternity trio (CE), F_{ST} and the p-value for F_{ST} (p Fst) in the 94 SNPs in the Verogen Forenseq™ Primer Mix A in a population sample of 391 individuals from El Salvador. N: number of individuals for which we could obtain a genotype for each locus.