

1 **Research Article**

2 **Palenque de San Basilio in Colombia: genetic data supports an oral history**

3 **of a paternal ancestry in Congo**

4 Naser Ansari-Pour,^{1*} Yves Moñino,² Constanza Duque,^{3,4} Natalia Gallego,^{3,5} Gabriel
5 Bedoya,³ Mark G Thomas⁶ and Neil Bradman⁷

6
7
8 ¹ Faculty of New Sciences and Technology, University of Tehran, Tehran, Iran

9 ²LLACAN, CNRS, Villejuif, Paris, France

10 ³ Universidad de Antioquia UdeA, Calle 70 No 52-21 Medellín Colombia

11 ⁴ Universidad Cooperativa de Colombia, Sede Medellín, Colombia

12 ⁵ Institución Universitaria Colegio Mayor de Antioquia, Medellín, Colombia

13 ⁶Department of Genetics, Evolution and Environment, University College London, London,
14 UK

15 ⁷Henry Stewart Group, 29/30 Little Russell Street, London, UK

16

17 **Corresponding Author:**

18

19 Dr Naser Ansari-Pour

20

21 E-mail: n.ansaripour@ut.ac.ir

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38 **Abstract**

39 The Palenque, a black community in rural Colombia, have an oral history of fugitive African
40 slaves founding a free village near Cartagena in the 17th Century. Recently linguists have
41 identified some 200 words in regular use that originate in a Kikongo language, with Yombe,
42 mainly spoken in the Congo region, being the most likely source. The non-recombining
43 portion of the Y chromosome (NRY) and mitochondrial DNA were analysed to establish
44 whether there was greater similarity between present day members of the Palenque and
45 Yombe than between the Palenque and 42 other African groups (for all individuals, n=2,799)
46 from which forced slaves might have been taken. NRY data are consistent with the linguistic
47 evidence that Yombe is the most likely group from which the original male settlers of
48 Palenque came. Mitochondrial DNA data suggested substantial maternal sub-Saharan African
49 ancestry and a strong founder effect but did not associate Palenque with any particular
50 African group. In addition, based on cultural data including inhabitants' claims of linguistic
51 differences, it has been hypothesized that the two districts of the village (Abajo and Arriba)
52 have different origins, with Arriba founded by men originating in Congo and Abajo by those
53 born in Colombia. Although significant genetic structuring distinguished the two from each
54 other, no supporting evidence for this hypothesis was found.

55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70

71

72

73 **Introduction**

74

75 In many locations throughout the Caribbean and Latin America during the Atlantic
76 slave trade, runaway slaves in pursuit of freedom established fortified villages. In Colombia,
77 these walled towns (known as palenques) were famed for their resistance to the Spanish
78 military conquest. This reputation is evident from colonial records which tell how inhabitants
79 successfully repulsed attacks by the authorities [1]. Despite their resistance, Palenque de San
80 Basilio is the only palenque to have survived to the present day [2].

81 Palenque de San Basilio (Palenque for short) is located some 70 km south east of the
82 regional capital of Bolivar, Cartagena, in North West Colombia (10.1°N, 75.2°W) (see Fig.
83 1). The residents comprise a community of about 3,500 individuals divided between two
84 major districts, Arriba and Abajo, although the reason for the division is not established [2,
85 3]. They have remained largely isolated from the prevailing Hispanic culture, living by
86 subsistence farming together with cattle husbandry [4]. Their oral history is one of descent
87 from a group of male slaves who escaped captivity early in the 17th century from nearby
88 Cartagena (then a major centre of the Latin American slave trade [5]).

89 Interestingly, Palenque is the only Colombian black community that speaks a creole
90 Spanish known as Palenquero [6]. Linguistic analysis of this creole led to the suggestion that
91 the language of the founding group originated in the area of present day Congo and/or
92 northern Angola [4, 7]. More recently, detailed lexical research has established that
93 Palenquero contains more than 200 words of African origin [8, 9] and that Kikongo is the
94 only demonstrable donor of the vocabulary [10]. The Kikongo group of languages
95 encompasses several extant tongues which are spoken by approximately one million people
96 in Republic of the Congo [11]. Although the recorded vocabulary in Palenquero does not
97 suggest a particular origin among them, the ritual vocabulary [10] and oral history [12]
98 suggest that Yombe is the most probable source. Today Yombe is spoken by the Yombe
99 people, an ethnic group living mainly in Pointe-Noire (Republic of the Congo). Furthermore,
100 many members of the Palenque community have claimed that a) Arriba residents are more
101 traditional and have better conserved Palenquero than their Abajo counterparts; and b) the
102 founding men of Arriba were born in Congo while Abajo would have been populated by
103 Maroons born in Colombia (Yves Moñino (YM), field work in Palenque and [3]).

104 To clarify these questions concerning Palenque history, we undertook a genetic
105 analysis of individuals from both Arriba and Abajo. DNA analysis has proved useful in
106 revealing origins of ethnic communities (for example see [13-15]). Sex specific genetic
107 systems (the non-recombining portion of the Y chromosome (NRY) and mitochondrial DNA
108 (mtDNA)) have been analysed to reveal connections between geographically separated
109 diaspora communities sharing a common identity [15-17] and to evaluate support for
110 alternative oral histories [18]. Recently, the geographic distribution of NRY haplotypes and
111 time to the most recent common ancestor (TMRCA) of paternal haplogroups were interpreted
112 as suggesting a late, exclusively eastern, expansion of the Bantu speaking peoples (EBSP)
113 [19].

114 The geographic origins of African diasporas, in particular those created by the
115 Atlantic slave trade, have been investigated using NRY and mtDNA. In studies of the
116 populations of Cape Verde Islands [20, 21] and Sao Tome Island [22, 23], sex specific
117 genetic systems were used to elucidate both maternal and paternal origins. In the case of the
118 Palenque, analysis of HLA autosomal markers and antigens [24, 25], and recently NRY
119 variation [26] has suggested a greater proportion of recent African descent (RAD) than other
120 Colombian groups.

121 Although culturally and geographically isolated for most of its existence, during the
122 past few decades, the Palenque people have experienced more contact with those from
123 outside their group [10]. Therefore, in recent times, an increased level of gene flow may have
124 occurred. Given the substantial geographic structuring of NRY and mtDNA haplotypes at the
125 continental level, and assuming that the founding group was of RAD, genetic analysis can
126 provide evidence of geographic ancestry and potential gene flow from non-RAD groups. If
127 the male founders were of RAD and there has been little gene flow from Europeans and
128 Amerindians, it can then be expected that NRY haplotypes will match those common in sub-
129 Saharan Africa and will have low diversity respectively.

130 Palenque oral tradition provides a testable hypothesis for NRY but not mtDNA
131 variation. However, from colonial records, it appears that in the second half of the 18th
132 century 178 black families occupied Palenque [5]. Therefore, it can be hypothesised that the
133 majority of the females at that time had RAD. If there has been little female gene flow since
134 that time, then the expectation is that mtDNA haplotypes will match those commonly seen in
135 sub-Saharan Africa.

136 Sub-Saharan Africa is known for its relatively high human genetic diversity [27, 28],
137 and geographic structuring of mtDNA haplotypes has been recognized [29]. Furthermore, the

138 considerable increase in NRY polymorphic sites identified in recent years [30, 31] has
139 revealed geographic structuring of NRY haplotypes [19, 32]. These findings have made it
140 possible, in some cases, to reveal recent shared paternal descent of men with a RAD born
141 outside sub-Saharan Africa with men still living there [21, 23, 26].

142 To explore these questions about Palenque history based on anthropological and
143 linguistic studies we analysed NRY and mtDNA in the Palenque and 42 sub-Saharan groups.
144 We address the following three questions: a) Is there greater genetic similarity between the
145 inhabitants of Palenque and Yombe speakers than between Palenque and non-Yombe African
146 groups?; b) Is there a significant difference between the sex-specific genetic systems profiles
147 of present day residents of Abajo and Arriba?; and c) Are the NRY and mtDNA of the Arriba
148 inhabitants more similar to those of Republic of the Congo than are the NRY and mtDNA of
149 Abajo residents? Genetic data analysed in this paper support the prior hypotheses that (a)
150 Palenque have a paternal line founding origin in the Yombe and (b) there is significant
151 difference in NRY distribution between Abajo and Arriba, but not mtDNA. There is limited
152 NRY but not mtDNA support for an affirmative answer to (c).

153

154 **Materials and Methods**

155 **Sample collection**

156 In Palenque, buccal swabs were collected from males over eighteen years old
157 currently living in, or born in, the community. Donors were initially selected randomly but
158 after questioning, only one of each set of donors having a common paternal grandfather was
159 included in the study. Samples were collected from a very substantial proportion of
160 individuals satisfying the above criterion (estimated at >90%, n = 153: Abajo area n = 88,
161 Arriba area n = 52, others n = 13). Samples from eight groups in the Republic of the Congo
162 (n = 591) were collected at local gatherings in different areas of Brazzaville, Pointe Noire and
163 in the villages of Kakamoeka and Lovoulou, 90km and 70km inland from Pointe Noire
164 respectively.

165 Ethnographic data were gathered from each Palenque individual, adopting the
166 procedure reported in Ansari-Pour et al. [19]. Buccal swabs previously collected from 34 sub-
167 Saharan groups in West, Central West and South-East Africa, representing other potential
168 source populations for the Atlantic slave trade, were also analysed in this study (see
169 Supplementary Table S1; samples from all population groups other than Palenque were
170 included in Ansari-Pour et al. [19]). DNA from all Congolese and Palenque samples was

171 extracted using the Gentra protein precipitation method (Gentra Systems, Minneapolis) while
172 the standard phenol-chloroform method was used for all other samples

173 **DNA typing**

174 The battery of Y-chromosome presumed unique event polymorphisms (UEPs),
175 consisting of single nucleotide polymorphisms (SNPs) and insertion/deletion polymorphisms,
176 as well as a set of short tandem repeats (STRs) were typed in the Palenque samples as
177 described by Ansari-Pour et al. [19]. Briefly: (a) sixteen UEPs (see Figure 2) were used to
178 classify NRY into haplogroups, applying the nomenclature of the Y Chromosome
179 Consortium [31] with the ‘capital letter- mutation’ system, and within each haplogroup, (b)
180 six STRs (DYS19, DYS388, DYS390, DYS391, DYS392 and DYS393) were used to define
181 haplotypes. Equivalent Y chromosome data for all 42 sub-Saharan African population
182 samples, including the eight Congolese groups (see Table S1) were taken from Ansari-Pour et
183 al. [19].

184 The mtDNA HVR-1 region of all Congolese groups and Palenque was sequenced as
185 described by Veeramah et al. [18]. For all samples, HVR-1 variable site only (VSO)
186 haplotypes were determined by comparing sequences of nucleotide range 16020-16400 with
187 the revised Cambridge Reference Sequence [33]. Haplotypes were defined by substitutions,
188 insertions and deletions, and their corresponding nucleotide positions. Tentative mtDNA
189 haplogroup assignment, based on HVR-1 sequences, were inferred according to the scheme
190 of Salas et al. [34], although it should be noted that inferred haplogroups frequencies were
191 not used in our statistical analyses and are only presented for reference. To extend the
192 mtDNA dataset, HVR-1 haplotypes were also determined for 30 out of 34 non-Congolese
193 sub-Saharan population samples considered in the NRY analyses (i.e. all groups except Sena,
194 Tumbuka, Bantu speakers from Pretoria and Yao; unpublished data except for the Nigerian
195 groups [35]). To facilitate comparison of all population samples, the range of the HVR-1
196 region considered was reduced to 16023-16380.

197

198 **Statistical analysis**

199 Pairwise genetic differences between population samples were assessed using the
200 exact test of population differentiation (ETPD) [36] which is analogous to Fisher’s exact test
201 extended to an $m \times n$ matrix, where m is the number of groups and n is the number of distinct
202 haplotypes. Gene diversity and its standard error were estimated using the unbiased formula

203 of Nei [37]. Genetic distances calculated were: F_{ST} [38] based on UEP haplogroups, STR
204 haplotypes (respecting their classification within haplogroups, i.e. UEP+STR haplotypes),
205 and mtDNA HVR-1 haplotypes and imputed haplogroups, R_{ST} [39] based on six STRs on the
206 NRY, and Kimura's two-parameter model with gamma value of 0.47 [40] for mtDNA HVR-
207 1 sequences. It should be noted that F_{ST} is used here as in Thomas et al. [17] as a convenient
208 statistic summarizing multidimensional differences in allele frequencies. No further
209 assumptions regarding the underlying population genetic model were applied in its
210 interpretation, other than a monotonic relationship between F_{ST} and genetic differences.

211 Analyses based on a selection of UEPs may suffer from biases in their ascertainment.
212 However, given the geographic structuring of the NRY variation, the choice of UEPs is
213 appropriate for the comparisons of sub-Saharan and RAD individuals. To test if genetic
214 distances differed significantly from zero, haplogroups/haplotypes were permuted among
215 samples; 1,000 permutations were performed to generate a null distribution of pairwise
216 genetic distances.

217 All of the above analyses were performed using Arlequin software version 3.0 [41].
218 Principal Component Analysis (PCA) and the nonparametric 'Sign Test' were performed
219 using the 'R' statistical programming language (www.R-project.org) [42], using 'princomp'
220 and 'binom.test' functions respectively. PCA plots were used to visualize relationships
221 among population samples based on NRY haplogroup frequencies.

222

223 **Results**

224 **Frequencies of NRY haplogroups and NRY based genetic distances**

225 The frequencies of all observed NRY haplogroups in the Palenque and the 42 sub-
226 Saharan groups analysed in this study are included in Table 1. The phylogenetic relationships
227 of the haplogroups can be seen in Figure 2. Thirteen NRY haplogroups were observed, of
228 which ten were present in the Palenque. The modal haplogroup in the Palenque was E-U175
229 (27%), but the two districts of the village had different modal haplogroups (details below).
230 Notably, there were only three haplogroups present in sub-Saharan African groups not
231 observed in the Palenque dataset: DE-YAP which is found at very low frequency in Nigeria
232 [43]; A-M13 which forms a very basal clade in the NRY phylogeny and has a wide
233 distribution at low frequency in Africa [44-46]; and E-U181 which has been proposed as a
234 signature of an exclusively eastern expansion of the Bantu speaking peoples [19]. P-92R7 and
235 R1a1, both widely considered to be 'non-African origin haplogroups' [47], were observed at

236 18% and 2.7% respectively in Palenque while observed as a singleton or at low frequencies
237 and completely absent in sub-Saharan African groups respectively. STR haplotypes within
238 each haplogroup were then analysed (see Supplementary Table 2). Of note, the two most
239 common STR haplotypes within P-92R7 in Palenque were haplotype 14-12-24-11-13-13
240 (N=5) and its one-step neighbour (14-12-24-12-13-13) (N=5). Both were absent from the
241 sub-Saharan African dataset. The former has been designated the Atlantic Modal Haplotype
242 (AMH) due to its high frequency in Western European populations [48, 49].

243

244 Gene diversity in the Palenque based on all haplogroups and E-sY81 component
245 haplogroups was 0.830 ± 0.013 and 0.638 ± 0.035 respectively, while the equivalent statistics
246 in the sub-Saharan African dataset were 0.753 ± 0.007 and 0.679 ± 0.008 respectively (for
247 gene diversity in each individual group see Supplementary Table S3).

248

249 The genetic distinctiveness of Palenque compared with each of the sub-Saharan
250 African groups was apparent as assessed by ETPD ($P < 0.001$). Also, all F_{ST} values between
251 Palenque and the sub-Saharan African groups were significant as assessed by random
252 permutation (see Methods) ($P < 0.00001$) with only two below 0.05 (Chewa, an East African
253 group from Malawi ($F_{ST} = 0.027$) and Yombe ($F_{ST} = 0.035$)) (see Supplementary Table S4). F_{ST}
254 between the Chewa and Yombe was not significant. This pattern was also consistently
255 observed based on R_{ST} (see Supplementary Table S5). Comparison of haplogroup profiles in
256 Palenque, Chewa and Yombe, revealed twelve NRY haplogroups present in at least one of
257 the groups. Six were observed in all three groups (see Fig 3). Of the remaining six, four were
258 observed in Palenque and Yombe, one in Palenque and Chewa, and one was observed only in
259 the Palenque (see Fig. 3). Most notably, all the haplogroups observed in the Yombe were also
260 observed in the Palenque, while the proposed signature haplogroup of the eastern EBSP
261 (E1b1a8a1a; E-U181) [19] was absent in the Palenque and the Yombe.

262

263 Similar to the approach taken by Di Giacomo et al. [50], we compared the distribution
264 of NRY variation within E-sY81 (E1b1a; the signature haplogroup of EBSP [19, 35, 51, 52]),
265 a clade which was present in all population samples including Palenque, and observed only in
266 men of RAD. F_{ST} between the Palenque and the other groups, based on the frequencies of the
267 E-sY81 component haplogroups, revealed only two groups with a non-significant F_{ST} (Chewa
268 and Yombe) with the Yombe-Palenque $F_{ST} < 0.001$ (Supplementary Table S6). Based on the
same dataset, pairwise differentiation between Palenque and all sub-Saharan African

269 population samples were also assessed using ETPD. Interestingly, all were significant at the
270 5% level except the Yombe ($P=0.507$).

271 A PCA plot using only E-sY81 component haplogroup frequencies showed Palenque
272 as an outlier. While a mixed collection of Bantu speaking groups are nearer than other Niger-
273 Congo groups to Palenque, strikingly it is the Yombe who are the closest of all (Figure 4).

274 **The distribution of mtDNA variation**

275 The Palenque sample contained 26 mtDNA HVR-1 haplotypes. The modal haplotype
276 was at a frequency of 0.166 and five common haplotypes together accounted for 66.2% of the
277 total. The imputed haplogroups were almost all sub-lineages of L (the sub-Saharan African
278 modal haplogroup) with over 70% within either L1 or L3 (see Supplementary Table S7).
279 Nei's gene diversity, based on mtDNA HVR-1 haplotypes, was 0.903 ± 0.010 . Across all
280 sub-Saharan African groups (38 groups), 723 mtDNA HVR-1 haplotypes (427 singletons)
281 were observed. Gene diversity for the combined set was 0.993 ± 0.004 and in individual
282 groups ranged from 0.968 to 0.997 (see Supplementary Table S8) with mean of 0.987 and SD
283 of 0.007.

284 ETPD based on HVR-1 haplotypes was significant between Palenque and all sub-
285 Saharan groups. This was also the case based on imputed haplogroups (except Sundi ($N=25$)
286 with borderline P-value of 0.057). F_{ST} and $K2P$ between Palenque and the sub-Saharan
287 African groups were also all significant, and all were in the range of 0.037-0.066 and 0.039-
288 0.126 respectively (see Supplementary Tables S9 and S10 respectively).

289 **Intra-village analysis of Palenque**

290 Summary statistics were calculated for both districts (Abajo and Arriba) in Palenque
291 (see Supplementary Tables S11 and S12 for NRY and mtDNA raw data respectively). The
292 modal NRY haplogroups in Abajo and Arriba were E-U175 (34.1%) and E-U290 (23.1%),
293 respectively. The proportion of the E-sY81 clade in Abajo and Arriba was 55.3% and 51.9%
294 respectively and not significantly different ($P=0.727$). Gene diversity based on all NRY UEP
295 in Abajo and Arriba was 0.800 ± 0.024 and 0.853 ± 0.019 respectively, and 0.5597 ± 0.0632
296 and 0.6809 ± 0.0495 respectively when restricting analysis to E-sY81 NRY types. At the
297 UEP+STR level, the modal haplotype was one STR mutation different from the EBSP modal
298 haplotype (i.e. E-sY81-15-12-21-10-11-13) [19] in Abajo (E-sY81-16-12-21-10-11-13;
299 10.6%) and in Arriba (E-sY81-15-12-21-10-12-13; 19.2%), while the EBSP modal haplotype
300 was at a frequency of 5.9% and 5.8% respectively.

301 Analysing P-92R7 NRY, the frequency was not significantly different in the two
302 districts (15.3% in Abajo and 19.2% in Arriba, $P=0.639$). The distribution of constituent
303 haplotypes was also not significantly different as measured by ETPD ($P=0.738$).

304 To investigate the hypothesis derived from observed cultural differences and local
305 oral history (YM, personal field notes) that the two village districts can be distinguished from
306 each other, several statistical tests were performed using the NRY and mtDNA data (see
307 Supplementary Table S13). Strikingly, genetic comparisons based on mtDNA were not
308 significant based on both haplotypes and imputed haplogroups, while all comparisons using
309 NRY markers were significant at the 5% level. Notably F_{ST} between the two village districts,
310 calculated using E-sY81 component haplogroups only, was both significant ($P<0.05$) and
311 greater than between either of them and Yombe (Supplementary Table S14). The
312 distinctiveness of Abajo and Arriba NRY was also confirmed by ETPD (UEP, $P=0.006$;
313 UEP-E-sY81, $P=0.012$; UEP+STR, $P=0.017$). We then examined whether the ETPD was
314 significant because of haplogroups introduced into Palenque probably through non-RAD
315 introgression. Y chromosomes were divided into a) those collectively belonging to
316 haplogroups K, P and R (all with inferred origins outside sub-Saharan Africa) and b)
317 Y(xK,P,R). The difference was driven by the African Y(xK,P,R) NRY at both UEP
318 ($P=0.008$) and UEP+STR ($P=0.008$) levels and not by the non-African (Y(K,P,R) NRY
319 ($P=0.114$ and 0.271 at UEP and UEP+STR levels respectively). No significant difference was
320 observed between the two districts based on mtDNA HVR-1 haplotypes and imputed
321 haplogroups as assessed by ETPD ($P=0.985$ and $P=0.77$). When applying the same test both
322 districts differed significantly from all sub-Saharan African groups ($P<0.00001$) at the
323 haplotype level. Analysis of imputed haplogroups also showed a consistent pattern. This may
324 be due to the presence of non-African mtDNA haplotypes including those defined by 16290T
325 and 16319A (possibly belonging to the Amerindian A2 haplogroup) and high frequency of
326 founder haplotypes such as that bearing 16294T and 16309G (possibly haplogroup L2a1).

327 **Abajo and Arriba in the context of sub-Saharan Africa**

328 F_{ST} between Abajo and Arriba, treated as separate samples, and 42 sub-Saharan
329 African populations were estimated based on all haplogroups and E-sY81 component
330 haplogroups (Supplementary Table S14) but not mtDNA HVR-1 haplotypes, since no
331 significant genetic distance (F_{ST} and $K2P$) was observed between the two districts.

332

333 At the NRY-UEP level, all F_{ST} estimates were significant ($P < 0.05$) with the exception
334 of Abajo and Yombe. When considering only E-sY81 component haplogroups, F_{ST} was
335 similarly not significant between Abajo and (a)Yombe and (b) Chewa. In addition, Arriba
336 had a non-significant F_{ST} with Bembe and Yombe. Based on the magnitude of F_{ST} , at the
337 NRY-UEP level, the Congolese groups were split in half; four were closer to Arriba than
338 Abajo, and four were closer to Abajo than Arriba. However, at the E-sY81 level, six of seven
339 (excluding Yombe which had a non-significant F_{ST}) were closer to Arriba. Comparisons with
340 the complete African dataset presented a clearer difference. At both NRY-UEP and E-sY81
341 levels, Arriba was closer to 33 (Sign Test $P = 0.0001$) and 36 (Sign Test $P < 0.0001$) out of 42
342 sub-Saharan groups, respectively.

343
344

345

346 **Discussion**

347 The evolutionary processes of mutation and genetic drift, including founder-effect, as
348 well as the possible influence of natural selection, make direct inference of population history
349 from genetic data challenging. However, such challenges become more tractable when clear
350 hypotheses can be formulated from existing anthropological, linguistic and ethnographic
351 research. In such circumstances, genetic data can, in some cases, be analysed to test those
352 hypotheses. Even though the NRY and mtDNA are effectively single loci (because they are
353 non-recombining regions) they can be appropriate systems for testing such hypotheses,
354 particularly where the prior hypotheses concern only patrilineal or matrilineal history. In the
355 current study, there were three prior hypotheses which we address in turn in the following
356 sections.

357

358 **The founding fathers of the Palenque community were primarily Yombe**

359 Sex-specific genetic systems are particularly susceptible to genetic drift [53, 54]. The
360 larger the population size and the fewer the generations since a postulated event, the less the
361 effect of genetic drift should be. Because forced slaves are recorded to be mainly from Niger-
362 Congo speaking groups, we analysed a set of haplogroups (E-sY81) that are collectively in
363 high frequency in Niger-Congo speaking peoples but are at only low frequencies or absent in
364 other groups. This should, at least to some extent, have the added benefit of reducing the
365 effect of any recent contribution from Amerindian and European males. We repeated the

366 analysis including haplogroups within Y(xP,K,R) (see Figure 2), to which the NRY
367 haplotypes of the great majority of residents in sub-Saharan Africa belong.

368 Notably, in the PCA visualization, the Yombe are the closest to the present day
369 Palenque out of all the 42 sub-Saharan groups (see Fig. 4). In addition, Yombe was the only
370 group from the Republic of Congo for which there was not a significant F_{ST} value (Yombe
371 $P=0.378$, other seven groups $P < 0.001$). We also calculated F_{ST} distances after including ten
372 West African groups from Montano et al. [52] with E-sY81 chromosome set equal to or
373 above the minimum set in this study (Sundi, $N(E-sY81)=22$). Yombe remained the closest
374 group to Palenque. Analysing the Y(xP,K,R) set of haplotypes produced a similar outcome
375 but with the Chewa marginally closer to the Palenque than were the Yombe (both $F_{ST} <$
376 0.02). Interestingly, in both the E-sY81 and Y(xP,K,R) analyses, Yombe and Chewa had an
377 $F_{ST} < 0.001$.

378 Even though there is considerable genetic similarity among the many widely
379 distributed groups having an origin in the rapid EBS [19], the small genetic distance
380 between Chewa (a group from Malawi) and Palenque is so similar to that between the Yombe
381 and Palenque that it would be surprising, were it not for an oral history of the Chewa that
382 records an origin in the “Luba country of the southern Congo basin” [55]. This description
383 could place their origin only about 400 miles east of the region where Yombe is currently
384 spoken and may even reflect a migration from a more western location, passing through Luba
385 country rather than commencing within it. The date of this migration is uncertain with the
386 earliest record of the group as ‘Chévas’ only appearing in 1831-2 [55]. Marwick [55] also
387 records that the Chewa have an equally prevalent alternative origin story that places their
388 genesis south west of Lake Malawi. Marwick agrees with Hamilton [56] that the two
389 traditions can be reconciled by the migration from the north being by “chiefly invaders” who
390 gained control over “long-established autochthones”. Our results are more consistent with
391 this interpretation of the oral accounts, as the Chewa-Yombe genetic distances were non-
392 significant at the NRY level but highly significant at the mtDNA level ($P < 0.001$).

393 Additional support for the Yombe origin of the Palenque comes from the absence of
394 NRY E-U181 chromosomes in both the Yombe and Palenque and their presence in the
395 Chewa. The presence of the E-U181 – previously reported as characteristic of East African
396 populations – in the Chewa can be explained by post-migration male gene flow following
397 their arrival in Malawi [56]. Nevertheless, given that a prior hypothesis exists – based on
398 linguistic evidence – for a Yombe origin, and that no such evidence has yet been advanced to

399 support a Chewa origin, it is reasonable to conclude that the genetic analysis of NRY
400 haplotypes supports a Yombe origin.

401

402 **Is there a significant difference between the sex-specific genetic systems profiles of**
403 **residents of Abajo and Arriba?**

404 Differences in the paternal demographic histories of the two areas of the village are
405 clearly supported by the presence of different modal haplotypes, a slightly higher haplotype
406 diversity in Arriba compared with Abajo, and a significant ETPD between the two. The
407 summed frequencies of P-92R7 haplotypes and the distribution of STR haplotypes within this
408 haplogroup were, however, similar suggesting a similar extent of non-African genetic
409 introgression into the two districts. These results, in general, contrasted with comparisons
410 using mtDNA where no statistical differences in diversity were observed and ETPD was not
411 significant ($P = 0.985$). The similarity in mtDNA profiles but not NRY supports field
412 observations of YM that patrilocality is practiced in Palenque with men choosing to live close
413 to their fathers and grandfathers, and women marrying men either from their own district or
414 another, with the latter being common.

415

416 **Are the NRY and mtDNA profiles of the Arriba inhabitants compared with Abajo**
417 **residents more similar to those of residents of Republic of the Congo?**

418 Since no significant difference was observed between Arriba and Abajo residents in
419 mtDNA haplotype distribution, the answer to the question posed is “no”. With respect to
420 paternal ancestry alone, the proposition has some limited support from the results of NRY
421 when restricting analysis to E-sY81 component haplogroups. Here, where F_{ST} was
422 significant, six out of seven of Congo groups had a smaller F_{ST} with the Arriba. More striking
423 is that when compared with all 42 sub-Saharan groups, Arriba had lower genetic distances
424 ($p=0.0001$). Although genetic drift cannot be discounted as the cause, one possible
425 explanation is that practices associated with Africa such as matrilinearity (involving
426 inheritance from a maternal uncle to his nephew, as seen in the Congo) were retained longer
427 in Arriba than Abajo. There might therefore be an association between cultural practice and
428 patterns of genetic diversity but not necessarily a causative relationship.

429

430

431

432 This study has explored an important aspect of the genetic ancestry of a fugitive
433 African slave community in Colombia and contributed to a fuller understanding of their
434 history. Further analysis of DNA of the Palenque alongside that of Colombian, European and
435 sub-Saharan African groups using genome-wide markers and a more detailed characterisation
436 of NRY and mtDNA should reveal more of the genetic history of the Palenque including
437 contributions made by other communities in Colombia.

438

439

440

441

442

443 **Ethics**

444

445 All samples were collected anonymously with informed consent. This study received ethical
446 approval from the Ministry of Health of the Republic of Congo (741/MSP/DGS/S), the
447 Scientific Committee of the Academic Corporation for the Studies of Tropical Pathologies of
448 the Universidad de Antioquia in Colombia (CPT-8840-03-054), the village council of
449 Palenque de San Basilio and the Joint UCL/UCLH Committees
450 on the Ethics of Human Research Committee A (99/0196).

451

452 **Data accessibility**

453 Supplementary Tables S1–S14 have been uploaded as part of the electronic supplementary
454 material.

455 **Competing interests**

456 The authors declare no conflict of interest.

457 **Author contributions**

458 N.A., N.B. and Y.M. conceived and designed the study, C.D., N.G., G.B. and N.B. collected
459 DNA samples. N.A. and M.G.T genotyped and sequenced the samples, Y.M. analysed the
460 anthropological data, N.A., M.G.T. and N.B. analysed the genetic data, and N.A., N.B. and
461 M.G.T wrote the paper.

462 **Acknowledgments**

463 We thank all DNA donors and those assisting in sample collection, and David Balding for
464 advice on statistical analysis. N.A was supported by NERC CASE award.

465

466

467 **References**

468

- 469 [1] Wade P. 1995 The cultural politics of Blackness in Colombia. *Am Ethnol* **22**, 341-357.
470 [2] Friedemann N, Patiño Rosselli C. 1983 *Lengua y sociedad en El Palenque de San Basilio*.
471 Bogotá, Instituto Caro y Cuervo.
472 [3] Moñino Y. 2012 Pasado, presente y futuro de la lengua de Palenque. In *Palenque*
473 *(Colombia): Oralidad, identidad y resistencia* (eds. G. Maglia & A. Schwegler), pp. 221-256.
474 Bogotá, Editorial Pontificia Universidad Javeriana.
475 [4] Bickerton D, Escalante A. 1970 Palenquero: A Spanish-based creole of northern
476 Colombia. *Lingua* **24**, 254-267.
477 [5] Del Castillo N. 1984 El lexico negro-africano de San Basilio de Palenque. *Thesaurus* **39**,
478 80-169.
479 [6] Schwegler A. 2006 Palenquero. In *Concise encyclopedia of languages of the world* (eds.
480 K. Brown & S. Oligvie). Oxford, UK, Elsevier Ltd.
481 [7] Granda G. 1971 Sobre la procedencia africana del habla "criolla" de San Basilio de
482 Palenque. *Thesaurus* **26**, 84-94.
483 [8] Schwegler A. 2000 The African vocabulary of Palenque (Colombia). Part 1: Introduction
484 and corpus of previously undocumented Afro-Palenquerisms. *J Pidgin Creole Lang* **15**, 241-
485 312.
486 [9] Schwegler A. 2002 El vocabulario africano de Palenque (Colombia). Segunda Parte:
487 compendio de palabras (con etimologías). In *Palenque, Cartagena y Afro-Caribe: historia y*
488 *lengua* (eds. Y. Moñino & A. Schwegler). Tübingen, Max Niemeyer.
489 [10] Schwegler A. 2011 Palenque(ro): the search for its African substrate. In *Creoles, their*
490 *substrates, and language typology* (ed. C. Lefebvre), pp. 225-249. Amsterdam, John
491 Benjamins.
492 [11] Lewis MP. 2009 *Ethnologue: Languages of the World*. Dallas, Texas, SIL International.
493 [12] Moñino Y. 2007 Les rôles du substrat dans les créoles et dans les langues secrètes: le cas
494 du palenquero, créole espagnol de Colombie. In *Grammaires créoles et grammaire*
495 *comparative* (eds. K. Gadelii & A. Zribi-Hertz), pp. 49-72. Saint-Denis, Vincennes
496 University Press.
497 [13] Helgason A, Sigureth ardottir S, Gulcher JR, Ward R, Stefansson K. 2000 mtDNA and
498 the origin of the Icelanders: deciphering signals of recent population history. *Am J Hum*
499 *Genet* **66**, 999-1016. (doi:S0002-9297(07)64026-9 [pii]).
500 [14] Hill EW, Jobling MA, Bradley DG. 2000 Y-chromosome variation and Irish origins.
501 *Nature* **404**, 351-352. (doi:10.1038/35006158).
502 [15] Thomas MG, Skorecki K, Ben-Ami H, Parfitt T, Bradman N, Goldstein DB. 1998
503 Origins of Old Testament priests. *Nature* **394**, 138-140. (doi:10.1038/28083).
504 [16] Thomas MG, Parfitt T, Weiss DA, Skorecki K, Wilson JF, le Roux M, Bradman N,
505 Goldstein DB. 2000 Y chromosomes traveling south: the cohen modal haplotype and the
506 origins of the Lemba--the "Black Jews of Southern Africa". *Am J Hum Genet* **66**, 674-686.
507 (doi:10.1086/302749).

508 [17] Thomas MG *et al.* 2002 Founding mothers of Jewish communities: geographically
509 separated Jewish groups were independently founded by very few female ancestors. *Am J*
510 *Hum Genet* **70**, 1411-1420. (doi:10.1086/340609).

511 [18] Veeramah KR, Zeitlyn D, Fanso VG, Mendell NR, Connell BA, Weale ME, Bradman N,
512 Thomas MG. 2008 Sex-Specific Genetic Data Support One of Two Alternative Versions of
513 the Foundation of the Ruling Dynasty of the Nso' in Cameroon. *Curr Anthropol* **49**, 707-714.
514 (doi:10.1086/590119).

515 [19] Ansari Pour N, Plaster CA, Bradman N. 2013 Evidence from Y-chromosome analysis
516 for a late exclusively eastern expansion of the Bantu-speaking people. *Eur J Hum Genet* **21**,
517 423-429. (doi:10.1038/ejhg.2012.176).

518 [20] Brehm A, Pereira L, Bandelt HJ, Prata MJ, Amorim A. 2002 Mitochondrial portrait of
519 the Cabo Verde archipelago: the Senegambian outpost of Atlantic slave trade. *Ann Hum*
520 *Genet* **66**, 49-60. (doi:doi:10.1017/S0003480001001002).

521 [21] Goncalves R, Rosa A, Freitas A, Fernandes A, Kivisild T, Villems R, Brehm A. 2003 Y-
522 chromosome lineages in Cabo Verde Islands witness the diverse geographic origin of its first
523 male settlers. *Hum Genet* **113**, 467-472. (doi:10.1007/s00439-003-1007-4).

524 [22] Trovoadá MJ, Pereira L, Gusmao L, Abade A, Amorim A, Prata MJ. 2004 Pattern of
525 mtDNA variation in three populations from Sao Tome e Principe. *Ann Hum Genet* **68**, 40-54.

526 [23] Goncalves R, Spinola H, Brehm A. 2007 Y-chromosome lineages in Sao Tome e
527 Principe islands: evidence of European influence. *Am J Hum Biol* **19**, 422-428.
528 (doi:10.1002/ajhb.20604).

529 [24] Jimenez S, Martinez B, Hernandez M. 1996 HLA antigens and gene distribution in San
530 Basilio (SB) an isolated black population of Colombia [Abstract]. *Hum Immunol* **47**, 62.

531 [25] Arnaiz-Villena A, Reguera R, Parga-Lozano C. 2009 HLA Genes in Afro-American
532 Colombians (San Basilio de Palenque): The First Free Africans in America. *Open Immunol J*
533 **2**, 59-66.

534 [26] Noguera MC *et al.* 2013 Colombia's racial crucible: Y chromosome evidence from six
535 admixed communities in the Department of Bolivar. *Ann Hum Biol*.
536 (doi:10.3109/03014460.2013.852244).

537 [27] Bowcock AM, Ruiz-Linares A, Tomfohrde J, Minch E, Kidd JR, Cavalli-Sforza LL.
538 1994 High resolution of human evolutionary trees with polymorphic microsatellites. *Nature*
539 **368**, 455-457. (doi:10.1038/368455a0).

540 [28] Yu N *et al.* 2002 Larger genetic differences within africans than between Africans and
541 Eurasians. *Genetics* **161**, 269-274.

542 [29] Rosa A, Brehm A. 2011 African human mtDNA phylogeography at-a-glance. *J*
543 *Anthropol Sci* **89**, 25-58. (doi:10.4436/jass.89006).

544 [30] Sims LM, Garvey D, Ballantyne J. 2007 Sub-populations within the major European and
545 African derived haplogroups R1b3 and E3a are differentiated by previously phylogenetically
546 undefined Y-SNPs. *Hum Mutat* **28**, 97. (doi:10.1002/humu.9469).

547 [31] Karafet TM, Mendez FL, Meilerman MB, Underhill PA, Zegura SL, Hammer MF. 2008
548 New binary polymorphisms reshape and increase resolution of the human Y chromosomal
549 haplogroup tree. *Genome Res* **18**, 830-838. (doi:10.1101/gr.7172008).

550 [32] de Filippo C *et al.* 2011 Y-chromosomal variation in sub-Saharan Africa: insights into
551 the history of Niger-Congo groups. *Mol Biol Evol* **28**, 1255-1269.
552 (doi:10.1093/molbev/msq312).

553 [33] Andrews RM, Kubacka I, Chinnery PF, Lightowlers RN, Turnbull DM, Howell N. 1999
554 Reanalysis and revision of the Cambridge reference sequence for human mitochondrial DNA.
555 *Nat Genet* **23**, 147. (doi:10.1038/13779).

556 [34] Salas A, Richards M, Lareu MV, Scozzari R, Coppa A, Torroni A, Macaulay V,
557 Carracedo A. 2004 The African diaspora: mitochondrial DNA and the Atlantic slave trade.
558 *Am J Hum Genet* **74**, 454-465. (doi:10.1086/382194).

559 [35] Veeramah KR *et al.* 2010 Little genetic differentiation as assessed by uniparental
560 markers in the presence of substantial language variation in peoples of the Cross River region
561 of Nigeria. *BMC Evol Biol* **10**, 92. (doi:10.1186/1471-2148-10-92).

562 [36] Raymond M, Rousset F. 1995 An Exact Test for Population Differentiation. *Evolution*
563 **49**, 1280-1283.

564 [37] Nei M. 1987 *Molecular Evolutionary Genetics*. New York, Columbia University Press.

565 [38] Reynolds J, Weir BS, Cockerham CC. 1983 Estimation of the coancestry coefficient:
566 basis for a short-term genetic distance. *Genetics* **105**, 767-779.

567 [39] Slatkin M. 1995 A measure of population subdivision based on microsatellite allele
568 frequencies. *Genetics* **139**, 457-462.

569 [40] Kimura M. 1980 A simple method for estimating evolutionary rates of base substitutions
570 through comparative studies of nucleotide sequences. *J Mol Evol* **16**, 111-120.

571 [41] Excoffier L, Laval G, Schneider S. 2005 Arlequin (version 3.0): an integrated software
572 package for population genetics data analysis. *Evol Bioinform Online* **1**, 47-50.

573 [42] Team RDC. 2011 *R: a language and environment for statistical computing*. Vienna,
574 Austria, R Foundation for Statistical Computing.

575 [43] Weale ME *et al.* 2003 Rare deep-rooting Y chromosome lineages in humans: lessons for
576 phylogeography. *Genetics* **165**, 229-234.

577 [44] Underhill PA *et al.* 2000 Y chromosome sequence variation and the history of human
578 populations. *Nat Genet* **26**, 358-361. (doi:10.1038/81685).

579 [45] Semino O, Santachiara-Benerecetti AS, Falaschi F, Cavalli-Sforza LL, Underhill PA.
580 2002 Ethiopians and Khoisan share the deepest clades of the human Y-chromosome
581 phylogeny. *Am J Hum Genet* **70**, 265-268. (doi:10.1086/338306).

582 [46] Knight A, Underhill PA, Mortensen HM, Zhitovovsky LA, Lin AA, Henn BM, Louis D,
583 Ruhlen M, Mountain JL. 2003 African Y chromosome and mtDNA divergence provides
584 insight into the history of click languages. *Curr Biol* **13**, 464-473. (doi:S0960982203001301
585 [pii]).

586 [47] Wells RS *et al.* 2001 The Eurasian heartland: a continental perspective on Y-
587 chromosome diversity. *Proc Natl Acad Sci U S A* **98**, 10244-10249.
588 (doi:10.1073/pnas.171305098).

589 [48] Moore LT, McEvoy B, Cape E, Simms K, Bradley DG. 2006 A Y-chromosome
590 signature of hegemony in Gaelic Ireland. *Am J Hum Genet* **78**, 334-338.
591 (doi:10.1086/500055).

592 [49] Wilson JF, Weiss DA, Richards M, Thomas MG, Bradman N, Goldstein DB. 2001
593 Genetic evidence for different male and female roles during cultural transitions in the British
594 Isles. *Proc Natl Acad Sci U S A* **98**, 5078-5083. (doi:10.1073/pnas.071036898).

595 [50] Di Giacomo F *et al.* 2004 Y chromosomal haplogroup J as a signature of the post-
596 neolithic colonization of Europe. *Hum Genet* **115**, 357-371. (doi:10.1007/s00439-004-1168-
597 9).

598 [51] Berniell-Lee G *et al.* 2009 Genetic and demographic implications of the Bantu
599 expansion: insights from human paternal lineages. *Mol Biol Evol* **26**, 1581-1589.
600 (doi:10.1093/molbev/msp069).

601 [52] Montano V, Ferri G, Marcari V, Batini C, Anyaele O, Destro-Bisol G, Comas D. 2011
602 The Bantu expansion revisited: a new analysis of Y chromosome variation in Central
603 Western Africa. *Mol Ecol* **20**, 2693-2708. (doi:10.1111/j.1365-294X.2011.05130.x).

- 604 [53] Scheinfeldt LB, Soi S, Tishkoff SA. 2010 Working toward a synthesis of archaeological,
605 linguistic, and genetic data for inferring African population history. *Proc Natl Acad Sci U S A*
606 **107**, 8931-8938.
- 607 [54] Jobling MA, Tyler-Smith C. 2003 The human Y chromosome: an evolutionary marker
608 comes of age. *Nat Rev Genet* **4**, 598-612.
- 609 [55] Marwick MG. 1963 History and tradition in East Central Africa through the eyes of the
610 Northern Rhodesian Cewa. *J Afr Hist* **4**, 375-390.
- 611 [56] Hamilton RA. 1955 Oral Tradition: Central Africa. In *History and Archaeology in Africa*
612 (ed. R.A. Hamilton). London, London: School of Oriental and African Studies.

613
614

615

616

617 **Figure Legends**

618

619 Figure 1. Geographic location of the village of Palenque de San Basilio in Colombia.

620

621 Picture adapted from Google maps (www.maps.google.com).

622

623 Figure 2. Phylogenetic relationships of UEP markers used to define NRY haplogroups.

624

625 The box identifies the E-sY81 (E1b1a) clade, exclusively observed in population groups with recent African ancestry.

626

627 Figure 3. NRY haplogroup profiles in Chewa, Palenque and Yombe.

628

629 Note that all haplogroups present in Yombe are observed in Palenque. Haplogroups unobserved in a group are shown as
630 blank. For definition of abbreviations of population names see Table 1.

631

632 Figure 4. Visual representation of genetic relationships among all groups using PCA based on
633 NRY UEP within E-sY81 (E1b1a).

634

635 SE represents population groups in South-East Africa, namely CH, YA, TU, SE and BN. Percentages in parentheses are the
636 amount of variation explained by each component. For definition of abbreviations of population names see Table 1.

637