1 **Different historical generation intervals in human populations**

2 **inferred from Neanderthal fragment lengths and patterns of mutation**

3 **accumulation**

4

5 Moisès Coll Macià*#[1], Laurits Skov*[2], Benjamin Marco Peter[2] and Mikkel Heide Schierup#[1]

6

7    1. Bioinformatics Research Centre, Aarhus University, Aarhus C, Denmark

8    2. Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany

9

10 * Contributed equally

11 # Corresponding authors: moicoll@birc.au.dk, mheide@birc.au.dk

12

13 **Abstract**

14

15 After the main out-of-Africa event, humans interbred with Neanderthals leaving 1-2% of

16 Neanderthal DNA scattered in small fragments in all non-African genomes today[1,2]. Here we

17 investigate the size distribution of these fragments in non-African genomes[3]. We find

18 consistent differences in fragment length distributions across Eurasia with 11% longer

19 fragments in East Asians than in West Eurasians. By comparing extant populations and

20 ancient samples, we show that these differences are due to a different rate of decay in length

21 by recombination since the Neanderthal admixture. In line with this, we observe a strong

22 correlation between the average fragment length and the accumulation of derived mutations,

23 similar to what is expected by changing the ages at reproduction as estimated from trio

24 studies[4]. Altogether, our results suggest consistent differences in the generation interval

25 across Eurasia, by up to 20% (e.g. 25 versus 30 years), over the past 40,000 years. We use

26 sex-specific accumulations of derived alleles to infer how these changes in generation

27 intervals between geographical regions could have been mainly driven by shifts in either male

28 or female age of reproduction, or both. We also find that previously reported variation in the

29 mutational spectrum[5] may be largely explained by changes to the generation interval and not

30 by changes to the underlying mutational mechanism. We conclude that Neanderthal fragment

31 lengths provide unique insight into differences of a key demographic parameter among human

32 populations over the recent history.

**Introduction**

If Neanderthal sequences in all non-Africans stem from a single introgression event, then differences in Neandertal fragment length distribution across the world would be indicative of differences in the speed of the recombination clock. Assuming a constant number of recombinations per generation, this would then imply differences in the number of generations since the admixture event and consequently differences in generation times among populations. While recent studies point towards a single gene flow event[2], an additional admixture event private to Asians has also been proposed[6,7] because Asian genomes carry larger amounts of Neanderthal sequence compared to European genomes. However, Asian genomes will also have more archaic fragments if a single gene flow common to Eurasians was followed by dilution of Neanderthal content in Europeans due to subsequent admixture with a population without Neanderthal admixture[2,8].

An independent source of information for estimating differences in generation time is the rate and spectrum of derived alleles accumulating in genomes over a given amount of time[9,10]. Pedigree studies have shown that the yearly mutation rate slightly decreases when the generation time increases because the mutational burst in the germline before puberty represents a high proportion of new mutations in young parents[4]. Moreover, the relative proportion of different mutational types depends on both the paternal and maternal age at reproduction. This has been exploited to estimate differences in generation intervals for males and females between Neanderthals and humans[9].

Here we investigate archaic fragment length distributions among extant non-Africans genomes from the Simon Genome Diversity Project (SGDP)[3] and high coverage ancient genomes. We report strong evidence for a single Neanderthal admixture event shared by all Eurasian and American individuals, enabling us to make use of archaic fragment length distributions as a measure of generation intervals since admixture. Differences in estimated generation intervals are mirrored by concordant patterns of mutation accumulation, and suggest significant differences by up to 20% in the generation time interval experienced by different Eurasian regions since their splits.

64    **Results**

65

66    Neanderthal fragment length distributions differ across Eurasia

67

68    The average archaic fragment lengths in non-African individuals from the SGDP, inferred

69    using the approach of Skov et al[11], differs across Eurasia and America (Fig. 1a, S3,

70    Data1_archaicfragments.txt). It presents a clear west-east gradient with the lowest mean

71    fragment length in an individual from the Middle East (S_Jordanian-1, mean = 65.69 kb,  SE

72    = 2.49 kb, sd = 72.09 kb, S1) and the highest in an individual from China (S_Tujia-1, mean =

73    88.70 kb, SE = 3.29 kb, sd = 110.62 kb, S1). The pattern is qualitatively very similar when a)

74    median fragment length instead of mean lengths are used, b) restricting to fragments most

75    closely related to the Vindija Neanderthal genome, the sequenced Neanderthal that is most

76    closely related to the introgressing Neanderthal population[12] or c) only using high-confidence

77    fragments inferred by the model (Extended Figure 1a-c, S4). When individuals are grouped

78    into five main geographical regions, the average archaic fragment length distributions are

79    significantly different (P value < 1e-5, permutation test, S2) by up to 1.12 fold (Fig. 1b zoom

80    in, Table S1). These five regions also show significant differences in the number of archaic

81    fragments and in the amount of archaic sequence inferred per individual (P value < 1e-5 for

82    both, permutation test, S2, Fig. 1c and d, Table S1), mirroring the mean archaic fragment

83    length distribution patterns. In agreement with previous reports[2,7], we found, for example, that

84    East Asians have 1.32 fold more archaic sequence inferred per individual compared to West

85    Eurasians (P value < 1e-5, permutation test, S2, Fig. 1d, Table S1).

86

87    We next investigated whether the larger amount of archaic sequence in East Asians is

88    explained by having distinct archaic fragments due to a second Neanderthal admixture. We

89    did this by joining the fragments of the 45 East Asian individuals and comparing them to the

90    joined fragments of a subsample of 45 West Eurasian individuals (Extended Figure 2a and b,

91    S6, Extended Figure 3). A total of 916,369 kb of the genome is covered by archaic sequence

92    in East Asia and 866,945 kb in West Eurasia, with 485,255 kb (53%) of the archaic sequence

93    overlapping (Fig. 2a, Table S3). Thus, as a group, East Asia has 5% more genomic positions

94    with archaic introgression evidence. If we further remove fragments with the closest affinity to

95    the sequenced Denisovan, which East Asians are known to possess more of[13], the total

96    sequence covered by archaic fragments is almost identical (East Asia 853,065 kb, West

97    Eurasia 850,028 kb, Table S5). When we restrict to fragments with affinity to the Vindija or

98    Altai Neanderthal, East Asia has a ~7% higher proportion of the genome covered (East Asia

99    646,710 kb, West Eurasians 604,518 kb, Table S5). We ascribe this  latter difference to the

100   fact that shorter fragments in Western Eurasians both make them slightly harder to infer by

101 the Skov et al[11] approach and less likely to carry SNPs that directly classify them as closest

102 to the Vindija Neanderthal.

103 To compare shared fragments in terms of length, we only consider fragments in an East Asian

104 that overlap with regions in the genome of West Eurasians that contain archaic sequence and

105 vice versa (Extended Figure 1d, Extended Figure 2c, S6). We observe that shared fragments

106 in East Asian individuals are on the average 1.13 fold longer than in West Eurasians (P value

107 < 1e-5, permutation test, S2, Fig. 2b, Table S4) as also observed when all fragments were

108 used.

109

110 Based on these observations, we conclude that the vast majority and possibly all of

111 Neanderthal ancestry in East Asians and West Eurasians stems from the same Neanderthal

112 admixture event. The 32% higher total amount of archaic sequence in an East Asian compared

113 to a West Eurasian individual on average is primarily due to archaic fragments occurring at

114 higher frequency in East Asians (Fig. 2c, Extended Figure 2d, S6, Extended Figure 3). It is

115 unlikely that natural selection has acted much more strongly against archaic fragment

116 frequency in West Eurasia since the purging of Neanderthal introgression is expected to have

117 acted prior to the split of European and Asian populations[14,15]. We consider our observations

118 more compatible with Europeans mixing with a Basal Eurasian population with little or no

119 archaic content diluting the Neanderthal ancestry as has previously suggested from admixture

120 modelling using ancient samples[8]. Such a dilution process should have a negligible effect on

121 the Neanderthal fragment lengths observed today (Fig. 2b) but would shift the frequency

122 distribution of Neanderthal fragments as we observe (Fig. 2c). This leaves us with a difference

123 in the speed of the recombination clock and hence the number of generations since the

124 common admixture with Neanderthals as the major cause of differences in archaic fragment

125 length distributions.

126

127 Ancient genomes allow us to look at archaic fragment lengths back in time. We called archaic

128 fragments in three high-coverage ancient samples included in the SGDP data; Ust'-Ishim

129 (dated 45,000 BP, equally related to all Eurasians)[16], Stuttgart (dated 7,000 BP farmer, West

130 Eurasian ancestor)[17] and Loschbour (dated 8,000 BP hunter-gatherer, West Eurasian

131 ancestor)[17] (S3, Fig. S1, Table S2). As expected, the archaic fragments are much longer for

132 Ust'-Ishim compared to any of the ancient and extant individuals (Fig. 1b, Fig. S1, Table S2,

133 see also [16,18]). Loschbour's and Stuttgart's archaic fragments are on average longer than their

134 West Eurasian descendants. However, their mean fragment length are very similar to the other

135 extant populations, particularly East Asian populations (Fig. 1b zoom in) suggesting East

136 Asians archaic fragments have experienced the same amount of recombination as West

137 Eurasian ancestors, represented as Loschbour, 8,000 years ago. This corresponds to around

138    275 fewer generations in East Asia than in West Eurasian populations (assuming an average

139    generation time of 29 years) over the approximately 40,000 years since the split of European

140    and Asian populations[2,19–21]. Another way of stating this is that the 8,000 fewer years of

141    recombination over 40,000 years corresponds to a difference in generation time of about 20%

142    across Eurasia.

143

144    <u>Mutations accumulated differently across Eurasia</u>

145

146    The number of *de novo* mutations (DNM) transmitted to a child depends on the sex and the

147    age of the parents[4]. Thus, a change in generation time during recent human evolutionary

148    history, as suggested above, should leave a detectable pattern in the total number of

149    mutations accumulated. To test this, we estimated the number of derived alleles accumulated

150    in each individual's autosomes since the split of African and non-Afrcan populations (S7,

151    Data2_mutationspectrum.txt). This was done by first removing all derived alleles observed in

152    the Sub Saharan Africa outgroup, excluding those individuals with detectable West Eurasian

153    ancestry[3]. Furthermore, we masked all genomic regions with evidence of archaic introgression

154    in any individual since archaic variants would not be found in Sub-Saharan genomes and they

155    would affect our results because they accumulated under a different mutational process[9].

156    Masking those regions also ensures that this analysis is independent of the archaic fragment

157    length analysis above. After these procedures, we were left with ~20% of the callable genome

158    (S7).

159

160    Fig. 3a shows that the rate of accumulation of derived alleles is significantly different among

161    groups (P value = 2.8e-4, permutation test, S2, Table S6). West Eurasia has accumulated

162    1.09% more derived alleles than East Asia (P value = 1.18e-3, permutation test, S2) since the

163    Out-of-Africa event. However, this difference in the accumulation of derived alleles could only

164    have happened when West Eurasia and East Asia were separated, which is only a part of the

165    time since the Out-of-Africa (Fig. S3). If we assume >60,000 years for the out-of-Africa and a

166    West Eurasia/East Asia split of <40,000 years[2,19–21] the difference in the rate of derived allele

167    accumulation is at least 60,000/40,000*1.09%=1.64% while West-Eurasia and East Asia were

168    apart (SI8). Using the pedigree-based estimate of the relationships between mean parental

169    age and mutation rate per generation[4] (SI8), we estimate that this difference corresponds to a

170    2.68 or 3.39 years shorter generation interval in West Eurasia if East Asian mean generation

171    time was 28 or 32 years respectively (SI8). These are lower bounds of the inferred differences

172    in generation intervals since the difference between out-of-Africa and population split times is

173    minimized.

174

175    The age of parents at conception, and hence generation time, also impacts the frequency of
176    which types of single nucleotide mutations occur[4]. Thus, a shift in generation time is predicted
177    to change the spectrum of new mutations[9,10] and partially explain differences in mutation
178    spectrum described among human populations[5,22,23]. We calculated the relative frequencies
179    of the six different types of single nucleotide mutations depending on their ancestral and
180    derived allele (S7, Fig. S2, Table S7) and related that to the average Neanderthal fragment
181    length for each individual (Fig. 3b). We observe significant associations with average archaic
182    fragment lengths for all six types (Table S8). We further subdivide C>T mutations in three
183    types: CpG>TpG which present a distinct mutational process[24] (Fig. S2, Table S7), TCC>TTC,
184    which is in great excess in European genomes and has been studied as a population-specific
185    mutational signature[5,22] (Fig. S2, Table S7) and the rest, denoted as C>T' (Fig. 3b, Table S8).
186    We find that the frequency of CpG>TpG transitions depends the least on fragment length.

187

188    To investigate whether these correlations could be due to differences in generation time
189    between geographical regions, we reanalysed the proportion of DNM mutation types as a
190    function of mean parental ages in the deCODE trio data set[4,25] (Fig. 3b inserts, S9, Table S8).
191    Comparing the correlations from the SGDP data with the deCODE data we see a strong
192    correspondence for most mutational types: in all mutation types where correlations with either
193    dataset are significant, the direction of the effects are concordant (Figure 3b, Fig. S6). The
194    deCODE dataset has a slight bias towards probands having older fathers than mothers (mean
195    = 2.77 years, sd = 4.25, Fig. S5), and this could affect the response of mutation type fraction
196    depending on mean parental age. However, no major change in the correlation coefficients
197    was observed when only probands with similar parental ages were analysed (S9, Fig. S6).

198

199    Since there is no a priori reason to expect a relationship between archaic fragment lengths
200    and derived allele accumulation, we consider it likely that the same underlying factor has
201    affected both. The general correspondence of these correlations with those expected from
202    DNM studies supports our hypotheses that this causal element is a change in generation time.
203    More specifically, the matching decreasing correlation with parental age of TCC>TTC mutation
204    indicates that this mutation signature will increase when the mean parental age decreases.
205    Thus a considerable reduction in mean generation time in West Eurasians, as suggested in
206    this study, offers an alternative explanation to the excess of TCC>TTC mutations in that region
207    compared to the rest of the world[5,26].

208

209    An increase in the mean generation interval can be due to an increase in paternal or maternal
210    age, or both. Anthropological studies suggest that males have generally been older than
211    females at reproduction, but that the age gap is twice as large in hunter-gatherers compared

212  with sedentary populations[27]. To gain insight into sex-specific changes to generation time
213  intervals we first compared the accumulation of derived mutations between autosomes, which
214  spend the same amount of evolutionary time in both sexes, and X chromosomes, which spend
215  ⅔ of the time in females while ⅓ in males (S10). Thus, an increase of the male-to-female
216  generation interval is expected to increase the X chromosome to autosomes (X-to-A) mutation
217  accumulation ratio[28], although other factors such as reproductive variance and changes in
218  population size can also influence the ratio. Fig. 4a shows the X-to-A ratio of derived alleles
219  accumulated per base pair as a function of the mean archaic fragment length, as mean
220  generation time proxy, for the females in the SGDP data (S10). We observe that the X-to-A
221  ratio is significantly different among regions (P value = 3.6e-4, permutation test, S2). East
222  Asians have a higher X-to-A ratio compared to American and Central Asia and Siberia, with
223  similar Neanderthal fragment sizes, and higher than West Eurasians, with smaller Neanderthal
224  fragment sizes. This result is compatible with East Asians having a higher mean generation
225  time than West Eurasians primarily due to an increased paternal age at reproduction as
226  compared to Americans and Central Asia and Siberia where the age at reproduction of both
227  sexes are inferred to have increased similarly.

228  Another sex-specific mutation signature are C>G mutations in genomic regions with
229  clustered de novo mutations in old mothers[4,29]. This signature can be explored to compare
230  maternal ages among groups[9]. We estimated the proportion of derived C>G alleles to other
231  derived allele types in these genomic regions and contrasted it to the same ratio for the rest
232  of the genome, for each individual (S10). When samples are grouped in the 5 main regions,
233  the C>G ratio in DNM clusters differs significantly (P value = 2.77e-3, permutation test, S2),
234  increasing with increasing Neanderthal fragment length (Fig. 4b). Notably, America has a
235  higher ratio than Central Asia and Siberians for similar Neanderthal fragment lengths,
236  suggesting a relatively larger impact of old mothers to the overall mean generation time
237  throughout their history. This is in line with the X chromosome analysis in that longer
238  generation times in America were more driven by older mothers as compared to older fathers
239  in East Asia with an intermediate increase of both parental ages in Central Asia and Siberia.

240  Finally, the Y chromosome is also expected to accumulate more derived alleles in
241  populations with younger fathers, similarly to the autosomes, about 0.4-0.5% per year
242  difference in generation time between two populations. We observe a point estimate of 1.19%
243  larger accumulation between West Eurasia and East Asia (Fig. S7, Table S10) but this is not
244  significant with the limited data available for the Y chromosome (P value = 0.66, permutation
245  test, S2).

246   **Discussion**

247

248   We have shown that the length of Neanderthal fragments in modern human genomes can be

249   used to obtain meaningful information about a fundamental demographic parameter, the mean

250   generation interval. We estimate surprisingly large differences across eurasian and american

251   groups suggesting stable differences over tens of thousands of years. Our approach depends

252   on the assumption that archaic fragments trace back to a single Neanderthal admixture event

253   shared by all non-African populations, for which we provide further evidence. Consistent with

254   these results, the number of derived mutations accumulated in the geographic regions studied

255   here follow the expectations of the difference in generation time estimated from the fragment

256   lengths. The agreement between the recombination and the mutation clock signatures argues

257   against confounding factors. For example, a potential bias would be expected if the African

258   outgroup, here used to find archaic fragments in the other individuals, had experienced some

259   ancient gene flow from West Eurasia that we have not been able to detect. Such a scenario

260   would shorten and remove archaic fragments in West Eurasians, explaining the observed

261   gradient. However, it would also decrease the number of derived alleles in West Eurasia

262   compared to East Asia, which is the opposite to what we report.

263

264   Differences in generation intervals of the magnitude and duration that we estimate can account

265   for observed variation in the mutation spectrum of human populations without an underlying

266   change to the mutational repair system. An example of this is the increased frequency of the

267   TCC>TTC mutation in West Eurasians. The differences in generation time, inferred here from

268   archaic fragment lengths, explain more than half of the total variation among individuals

269   (adjusted $R^2$ = 55.53%).

270

271   Our results have direct implications for previous investigations of demographic human

272   parameters, which have typically assumed that the generation interval was shared and

273   constant for distinct human populations. Thus, future investigations should take variation in

274   the generation time under consideration. We do not have an explanation for the underlying

275   causes of large generation interval differences, but it is plausible that both low population

276   densities and harsh environmental conditions increase generation time, whereas agriculture

277   decreases mean generation times and reduces generation time differences between sexes.

278   With an increasing number of sequenced ancient and modern genomes we anticipate that the

279   approach we present here can be used to obtain a fine-grained picture of shifts in generation

280   interval during the last 40,000 years that can be directly related to changes in population

281   densities, climate and culture.

282 **Methods**

283

284 A description of all analyses performed in this study is detailed in the Supplementary
285 Information.

286

287 **Data availability**

288

289 The archaic fragments and their basic statistics are provided in Data1_archaicfragments.txt;
290 the counts of the 96 mutation types per individual per chromosome are provided in
291 Data2_mutationspectrum.txt (S11).

292

293 **Code availability**

294 The scripts coded to produce data and tables, perform statistical analysis and plot figures for
295 this manuscript are accessible on Github
296 (https://github.com/MoiColl/TheGenerationTimeProject).

297

298 **Acknowledgements**

299

307

308 **Contributions**

309

310 M.C.M., L.S. and M.H.S designed the study. M.C.M. and L.S. created the methods to assess
311 the data and, with M.H.S., analysed the results with input from B.M.P. M.C.M., L.S. and M.H.S.
312 wrote the manuscript with comments from B.M.P.

313

314 **Competing interests**

315

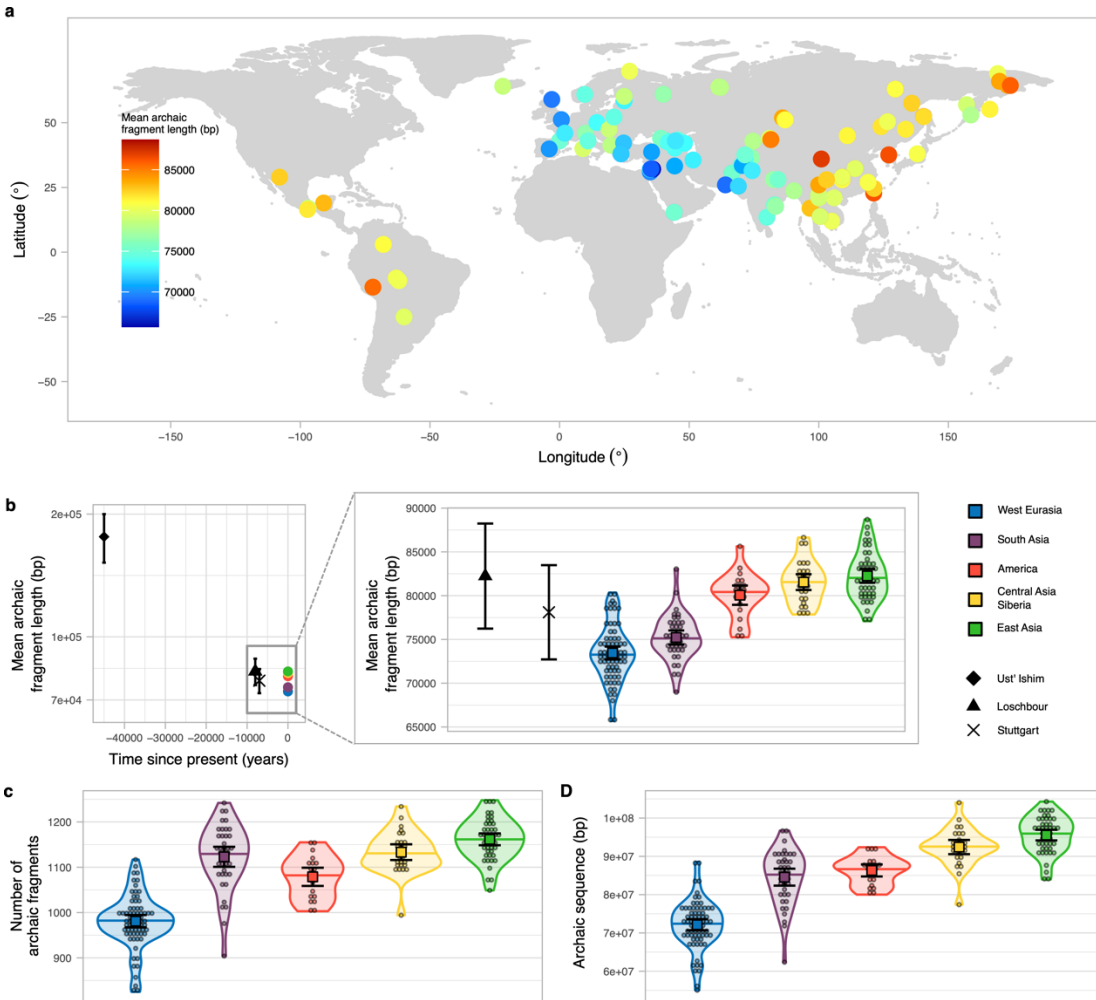316 The authors declare no competing interests.

**Fig. 1. Archaic fragment statistics distributions around the world and in ancient samples. a)** World map showing the samples from SGDP used in this study coloured according to the mean archaic fragment length. **b)** Mean archaic fragment length of extant geographical regions and ancient samples. Ust'-Ishim, Loschbour and Stuttgart mean archaic fragment length are shown as black points with specific shapes with their corresponding 95%CI as error bars. The average of the mean archaic fragment length among all individuals in each of the 5 main regions are shown as points (colour-coded). The zoom-in shows the mean archaic fragment length distribution per region (colour coded) as a violin plot. Individual values are shown as dots. The median is shown as a horizontal line in each violin plot. The mean and its 95%CI of each distribution is shown as a coloured square with their corresponding error bars. Loschbour and Stuttgart mean length are also shown for comparison. **c) and d)** the number of archaic fragments and the archaic sequence distributions respectively per region (colour coded) as violin plot. Individual values are shown as dots. The median is shown as a horizontal line in each violin plot. The mean and its 95%CI of each distribution is shown as a coloured square with their corresponding error bars. (width = 18cm)
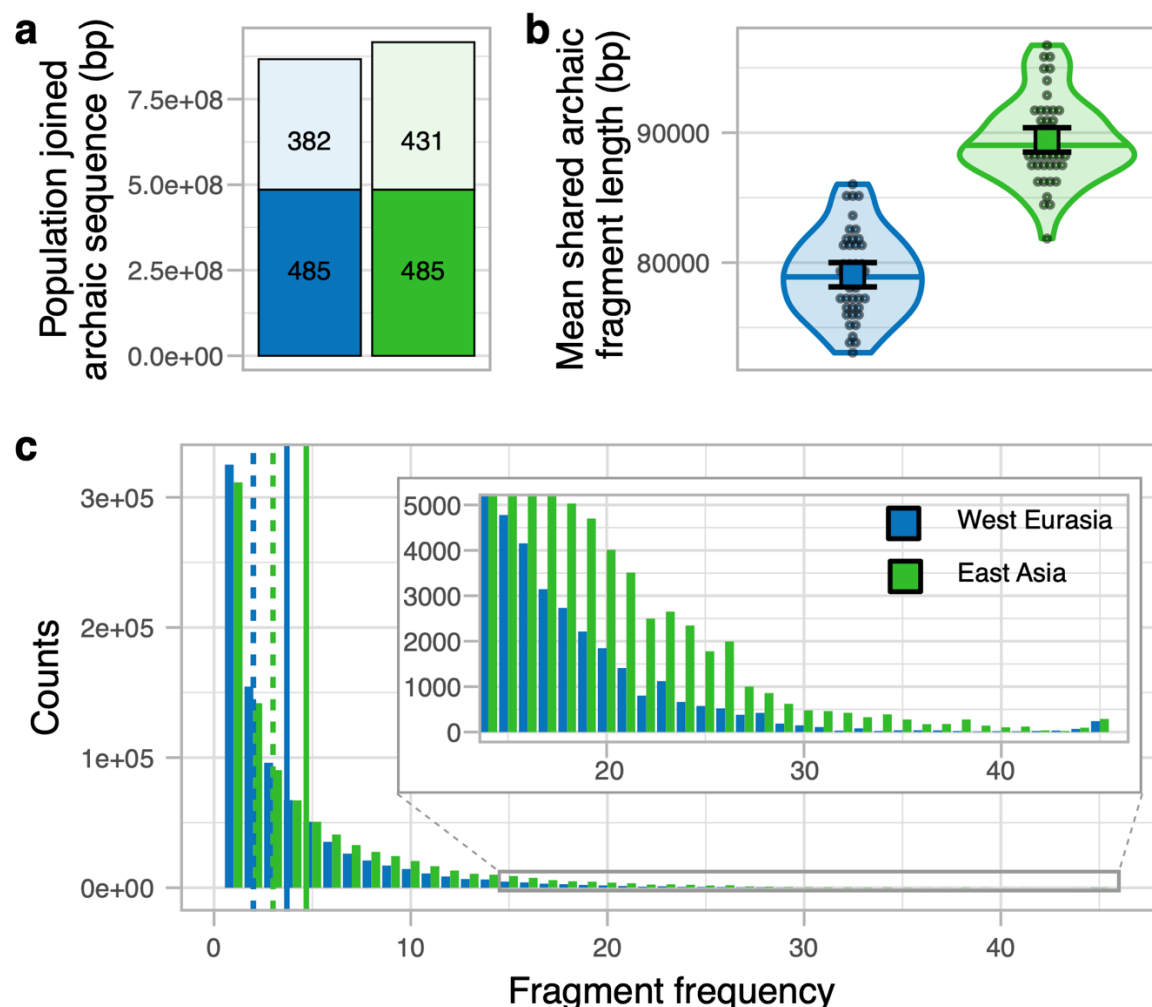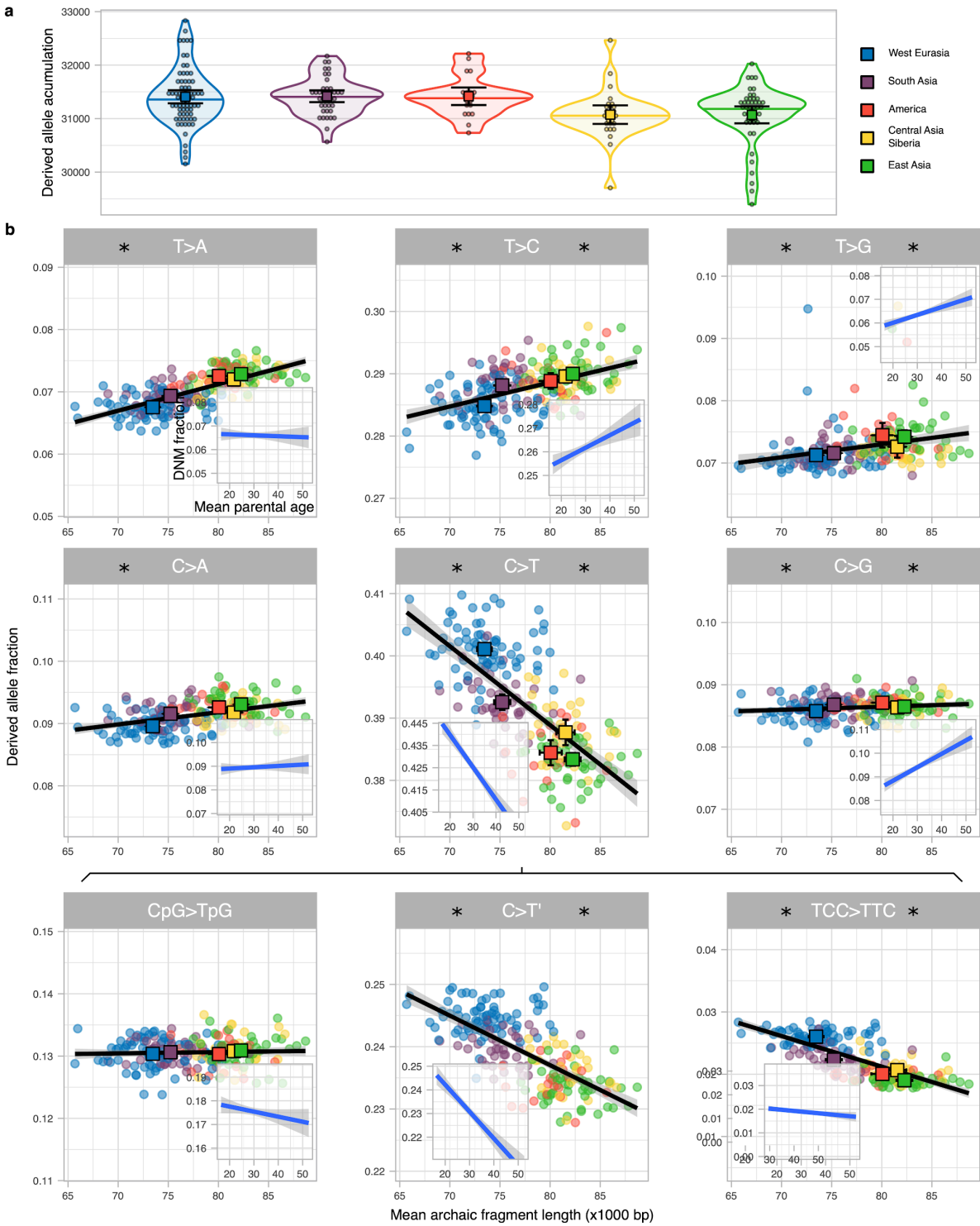
**Fig. 2. West Eurasia and East Asia archaic fragments comparison. a)** Joined archaic sequence in both geographic regions (colour coded). The portion of the bar painted in plain colour shows the shared amount between the regions. The rest of the column shows the sequence private of each region. The numbers in each section denote the corresponding archaic sequence in Mb. **b)** The mean archaic fragment length distributions of individual shared fragments among regions per region (colour coded) as violin plot. Individual values are shown as dots. The median is shown as a horizontal line in each violin plot. The mean and its 95%CI of each distribution is shown as a coloured square with their corresponding error bars. **c)** The number of 1 kb genomic windows (y-axis) in which an archaic fragment has been found in a certain amount of individuals (x-axis) for each region. The insert shows the high-frequency bins. Vertical lines show the mean (plain lines) and median (dashed lines) for each region. (width = 9cm)

346 **Fig. 3. Derived allele accumulation distributions and their mutation spectrum. a)**
347 Distribution of the derived allele accumulation (y-axis) per region (colour coded) as violin plot.
348 Individual values are shown as dots. The median is shown as a horizontal line in each violin
349 plot. The mean and its 95%CI of each distribution is shown as a coloured square with their
350 corresponding error bars. **b)** Correlation between the derived allele proportion (y-axis) with the
351 mean archaic fragment length (x-axis) for 9 mutation types. Each dot represents an individual
352 coloured according to the region they belong to. For each region, The mean and its 95%CI of
353 both axes is shown as a coloured square with their corresponding error bars. Linear
354 regressions (black lines) are shown with their corresponding SE (shaded area). For each
355 mutation, the linear regression and corresponding SE between the fraction of DNM and mean
356 parental age per proband of the deCODE data (**S9**) is shown as an insert. Note that the total
357 span of the y-axis is the same for all panels and inserts but centred at the mean value
358 specifically in each panel and insert. Asterisk on the left and right side of each mutation type
359 indicates that the slope of the linear regression is significantly different from 0 for the SGDP
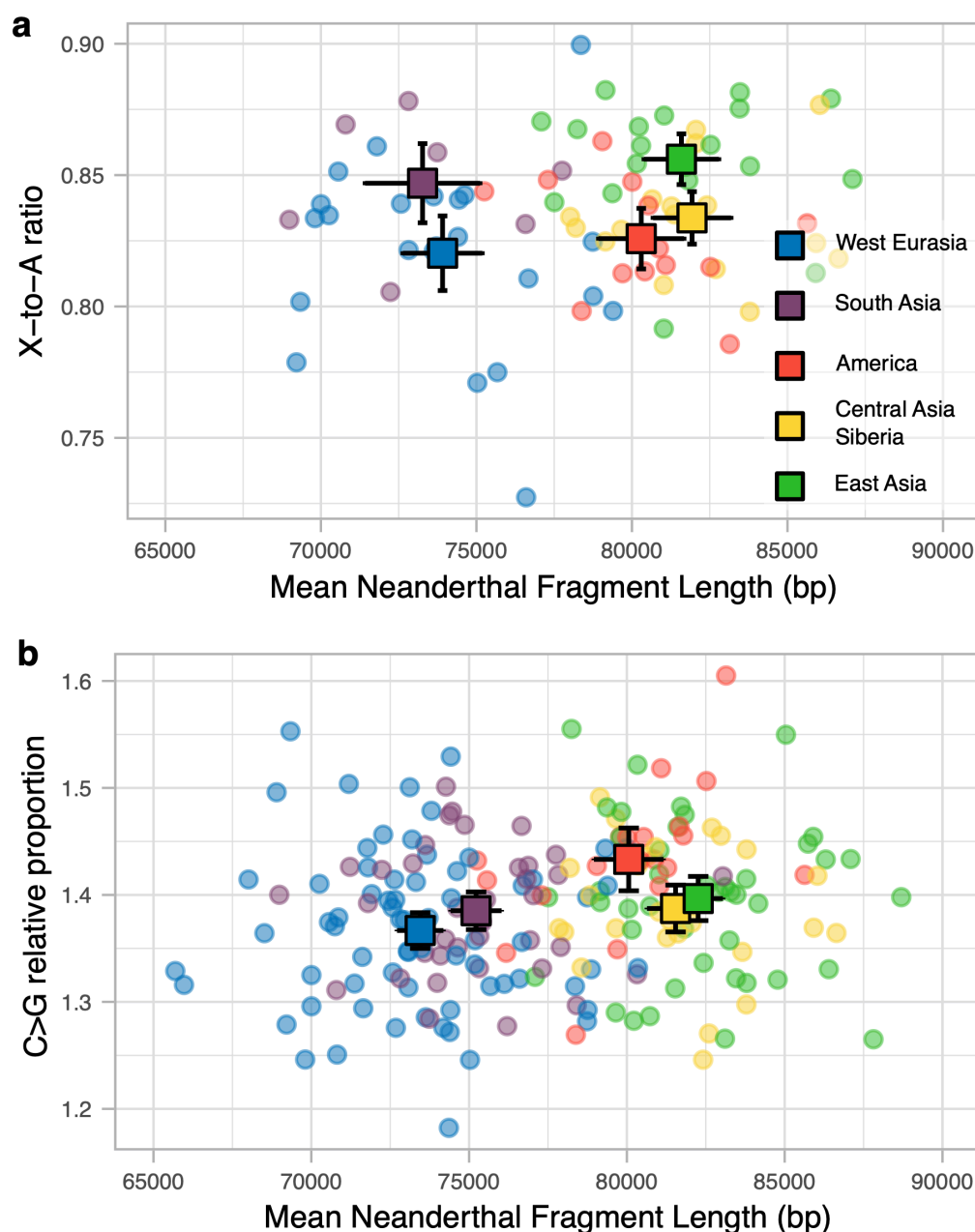360 and the deCODE data respectively. (width = 18cm)

13

**Fig. 4. Sex-specific mutation patterns. a)** Scatterplot of the X chromosome vs Autosome derived allele accumulation ratio (y-axis) and the mean Neanderthal fragment length (x-axis) for each region (colour coded). Each dot represents an individual in the corresponding population. The mean for each population for each axis is shown in squares and their 95%CI are denoted by the error bars. Only females were used to produce this plot. **b)** The same as in **a**, but the ratio between the proportions of C>G derived alleles in cDNM and the rest of the genome (**S10**). All samples were used to produce this plot. (width = 9cm)

**a**

Median Archaic Fragment Length



**b**

Vindija–like fragments



**c**

High–confidence fragments



**d**

East Asia and West Eurasia shared fragments



369

15

370   **Extended Figure 1. Archaic fragment length distribution around the world with specific**

371   **filters.** World map with samples from SGDP used in this study coloured according to the mean

372   or median average archaic fragment length applying filters to the data. **a)** Median archaic

373   fragment length is plotted instead of the mean. **b)** Only fragments with more SNPs shared with

374   the Vindjia genome than the Denisova or the Altai genomes are used **c)** Only high confidence

375   archaic fragments (posterior probability >= 90%) are used. **d)** Only shared individual fragments

376   (Extended Figure 2, **S6**) between East Asians and West Eurasians.

**Extended Figure 2. West Eurasia and East Asia fragment comparison methods.** Diagram showing the different methods to compare archaic fragments between West Eurasians and East Asians (**S6**). Each horizontal line represents a genome. Wide bands on each genome represent archaic sequences. East Asia is represented in green colours and West Eurasia in blue. Grey colours are used when sequences are shared by both. Plain colours denote joined sequence and transparent colours show individual sequences. Vertical dashed lines are mainly used to point to genomic windows of interest.  **a)** Joined region fragments. **b)** Shared and private joined region sequence. **c)** Shared and private individual fragments. **d)** Archaic frequency in 10 kb windows represented as the vertical grey lines intervals (note that in the main text, 1 kb windows are used instead).

**Extended Figure 3. The archaic landscape across the West Eurasian and East Asian genomes.** Each horizontal rectangle represents a chromosome (hg19). In each chromosome, it is shown the joined region fragments for West Eurasia (blue upper bands) and East Asia (green lower bands). The shared joined region fragments are shown as black bands in the middle of each chromosome. For each region, the number of individuals that have an archaic fragment in a particular 1kb window are represented as lines (maximum number of individuals is 45 for each region). Grey bands on the chromosomes show the non-callable portions of the genome (hg19).

# Supplementary Information

**1 - Confidence Interval calculation**

The mean and its confidence intervals (CI) for any statistic are calculated using the mean and standard deviations of the 100,000 bootstrap sampling distribution of the observed statistic. The code to compute them is provided on the GitHub page.

**2 - Statistical significance assessment by permutation test**

The statistical significance of a statistic to compare different groups is assessed by contrasting the observed statistic with a non-parametric null distribution. The null distribution is generated by permuting 100,000 the original data and calculating the statistic in each permutation. P values are then calculated as the fraction of permutations which yield a value as extreme or more extreme than what is observed in the data. If no such event is observed in all permutations, we considered the fraction to be < 1/100,000 = 1e-5. The significance level ($\alpha$) in all tests is considered to be 0.05.

To test if there are differences between two groups for a statistic (for example, average archaic fragment length), we subtract the means of each group. In this case, since this test is a two-tailed hypothesis test, we multiply the obtained P value by two. When we test differences for multiple populations, we compute the F statistic.

The code to compute the statistical significance is provided on the GitHub page.

421 **3 - Identification of archaic fragments in non-African individuals and ancient samples**

422

423 We called archaic fragments in individuals of the Simon Genome Diversity Project (SGDP)

424 from West Eurasia, South Asia, America, Central Asia Siberia and East Asia regions as

425 described in [9,11] - a step by step tutorial is also available at

426 https://github.com/LauritsSkov/Introgression-detection.

427

428 In short, the method first removes a set of variants (SNPs) which are present in an outgroup

429 with no presumed archaic admixture (Sub-Saharan African populations) from the samples in

430 which we want to detect archaic fragments (non-Africans). Then, taking into account window-

431 specific mutation rate and callability, the method classifies non-overlapping windows into

432 archaic ancestry and non-archaic ancestry depending on the derived allele density.

433

434 <u>Outgroup variants set, window mutation rate and callability and derived allele polarization</u>

435

436 To generate the set of variants in the outgroup, we merged all variants from the following

437 populations:

438

439   1. All Sub-Saharan Africans (populations: YRI, MSL, ESN) from the 1000 Genomes

440      Project [30] and

441   2. All Sub-Saharan African populations from SGDP (this excludes Sharawi and Mozabite

442      populations from the African supergroup) [3] except individuals from the Masai and

443      Somali populations because they are reported to have some West Eurasian genetic

444      component.

445

446 We determine the background mutation rate as the SNP density in the outgroup samples in

447 windows of 100 kb.

448

449 To generate the callability regions, we merged the following files:

450

451   1. 1000 Genomes Project Callability file (hg19)

452

453 ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/accessible_genome_

454 masks/StrictMask/

455

456   2. Repeatmask file (hg19)

457

20

458    hgdownload.cse.ucsc.edu/goldenpath/hg19/bigZips/chromFaMasked.tar.gz

459

460    To polarize alleles into ancestral and derived alleles we used the following file:

461

462    http://web.corral.tacc.utexas.edu/WGSAdownload/resources/human_ancestor_GRCh37_e7

463    1/

464

465    Training the Hidden Markov model and decoding archaic fragments in each sample

466

467    For each extant non-African individual and 3 ancient samples (Stuttgart, Loschbour and Ust-

468    ishim) from the SGDP, we filtered out all sites where the derived variant is found in our

469    outgroup population and sites that are not in our callable regions.

470

471    Then we trained the HMM and found the best fitting emission and transition values. Finally we

472    identified tracks of archaic introgression in the whole genome of each individual

473    (Data1_archaicfragments.txt). The archaic fragments in Stuttgart, Loschbour and Ust-ishim

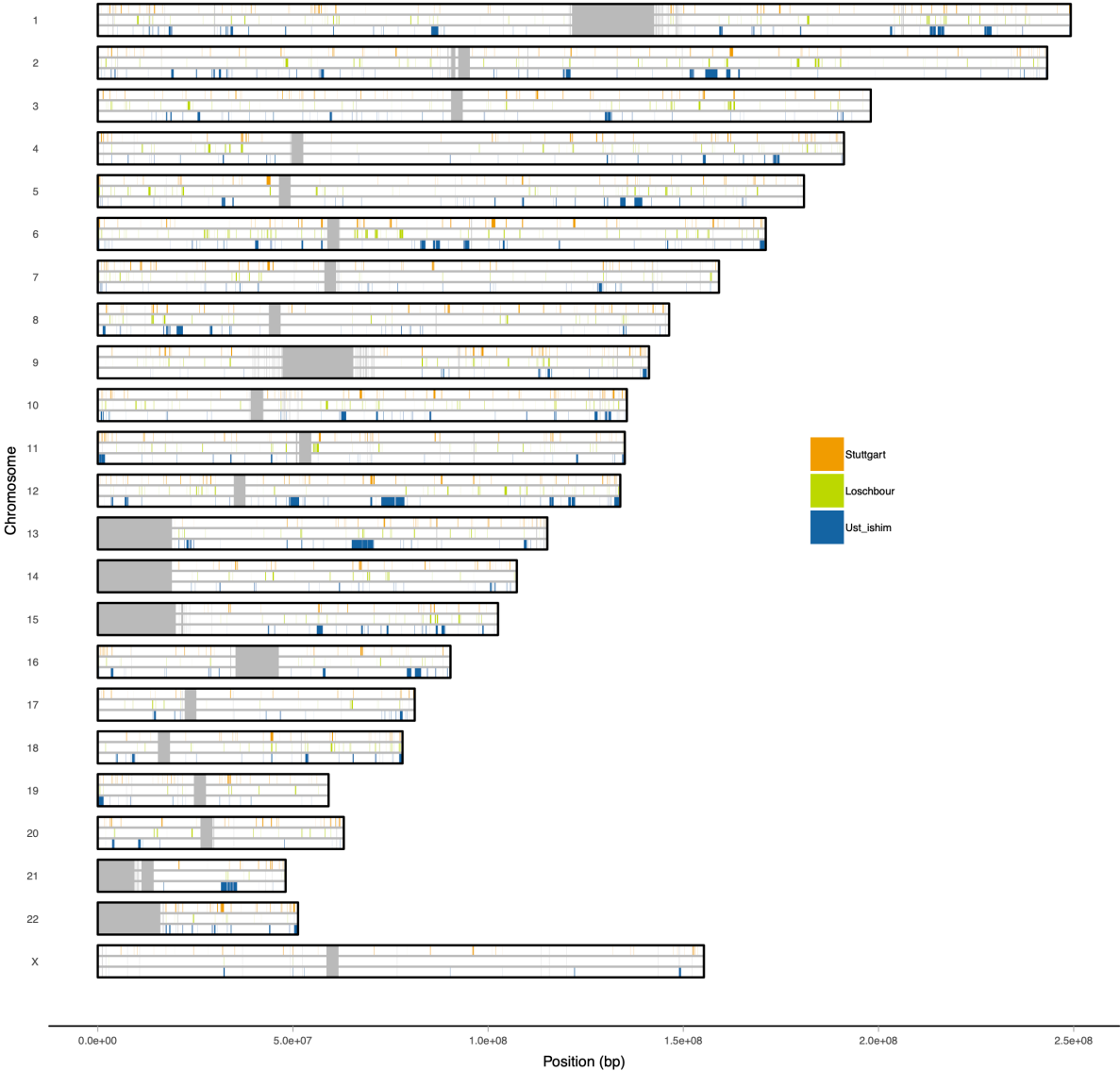474    are visualized in Fig. S1.

**Fig. S1**. **Archaic fragments in Ust 'Ishim, Loschbour and Stuttgart ancient samples.** Each horizontal rectangle represents a chromosome (hg19). In each chromosome, it is shown the archaic fragments found in Ust 'Ishim, Loschbour and Stuttgart ancient samples (colour coded). Wide grey bands on the chromosomes show the non-callable portions of the genome (hg19).

22

**4 - Archaic fragment length gradient around the world is consistent to multiple filters**

We studied the robustness of the difference in mean archaic fragment length among the 5 geographical groups studied applying multiple filters.

    1) Median instead of mean

The mean is very sensitive to outliers. In our case, very long archaic fragments, for example in East Asians, could increase the mean and thus show an unrealistic pattern among regions. To avoid that, we use median instead because it is more robust to outliers.

    2) Vindija genome-like fragments

The method used in this study is able to find archaic fragments whose variation is not fully captured by the sequenced archaic individuals [11]. The difference in archaic fragment length can potentially be affected if there is a distinct archaic content among the extant populations studied here - for example, a greater and more recent Denisova component in Asia [31].

It is known that the majority of the archaic component in Eurasia and America is from a Neanderthal population closely related to the Vindija genome [12]. Thus, we restrict fragments used in this analysis to share more variation with the Vindija Neanderthal genome than the Altai Neanderthal genome or the Denisovan genome.

    3) High confidence archaic fragments

The method used in this study, returns the archaic fragments found in a genome with an associated mean posterior probability. We restricted archaic fragments compared to be of a high confidence (mean posterior probability >= 0.9).

When we study the archaic fragment difference among individuals in Eurasia and America applying the multiple filters explained above, we can see that the pattern observed using all fragments holds (Extended Figure 1). We conclude that the difference in archaic fragment length is genuine and not depending on the factors exposed above.

514 **5 - Archaic fragment summary statistics per individual per region in extant populations**

515 **and ancient samples**

516

| Region | Number of samples | Number archaic fragments | | Archaic seq (bp) | | Mean archaic fragment length (bp) | |
|---|---|---|---|---|---|---|---|
| | | mean | SE | mean | SE | mean | SE |
| West Eurasia | 71 | 980.92 | 6.93 | 72,134,504.74 | 747,098.22 | 73,449.52 | 375.37 |
| South Asia | 39 | 1,122.99 | 11.29 | 84,561,986.37 | 1,127,788.19 | 75,221.52 | 410.92 |
| America | 20 | 1,078.87 | 10.25 | 86,322,668.28 | 803,716.12 | 80,057.38 | 561.94 |
| Central Asia Siberia | 27 | 1,133.28 | 8.86 | 92,424,440.59 | 946,306.77 | 81,548.03 | 460.74 |
| EastAsia | 45 | 1,161.58 | 6.65 | 95,552,691.46 | 712,217.98 | 82,258.79 | 402.10 |

517

518 **Table S1**. **Archaic fragment summary statistics per individual per region.** Summary

519 statistics of the fragments found among the individuals of the 5 main regions. For each statistic,

520 the mean and the of SE (S1) is provided.

521

| Ancient samples | Number archaic fragments | Archaic seq (bp) | Archaic fragment length (bp) | |
|---|---|---|---|---|
| | | | mean | SE |
| Ust 'Ishim | 763 | 134,360,000 | 176,100.09 | 12,464.19 |
| Loschbour | 921 | 75,757,000 | 82,190.02 | 3,087.97 |
| Stuttgart | 1,101 | 85,950,000 | 78,070.12 | 2,705.90 |

522

523 **Table S2**. **Archaic fragment summary statistics per ancient sample.** Summary statistics

524 of the fragments found in the three ancient samples. For the archaic fragment length, the mean

525 and the of SE (S1) is provided.

**6 - West Eurasia and East Asia fragment comparison**

In this study, we compare fragments in West Eurasians and East Asians. The more individuals used to recover archaic fragments, the more undiscovered fragments can be found [9]. Thus, the imbalance in the number of individuals in each region in the SGDP data (71 West Eurasians and 45 East Asians) can potentially affect any comparison between the two regions. Therefore, we downsample the number of individuals used in West Eurasians to 45 randomly chosen individuals to make comparisons fair.

First, we join all overlapping fragments for each region, hereby "joined region fragments" (Extended Figure 2). To do that, we used bedtools software [32] with the following command:

```
bedtools merge -i ind1_regx.bed ind2_regx.bed … indN_regx.bed > joined_regx.bed
```

where x denotes either West Eurasia or East Asia regions and N denotes the number of individuals in the corresponding region.

Then, we compared how much archaic sequence the two regions share (Extended Figure 2). For that, we call the intercept between the two joined sets of fragments. We refer to it as the "shared joined region sequence". We use the following command:

```
bedtools intercept -a joined_regx.bed -b joined_regy.bed > shared_joined.bed
```

where x denotes either West Eurasia or East Asia and y denotes the other region different than x.

It follows that the rest of the fragments not included in this set are the "private joined region sequence".

The amount of sequence for shared, private and total joined region fragments are provided in Table S3.

For each individual, we classified the fragments as shared depending upon if there was an overlapping fragment in the other joined region fragments (Extended Figure 2). We name these fragments as "shared individual fragments". To get them, we ran the following command:

```
beedtools intercept -u -a indn_regx.bed -b joined_regy.bed > shared_indn_regx.bed
```

It follows that the rest of the fragments not included in this set are the "private individual fragments".

Summary statistics for shared and private individual fragments are provided in Table S4.

Finally, we calculated the number of individuals that have an overlapping archaic fragment in a certain 1kb window in the genome. This way, we calculate the **archaic frequency**. For that,

25

571    we first divided each fragment in the joined region fragments into 1 kb segments

572    (`joined_regx_1kb.bed`). Then, we counted the number of individuals with an overlapping

573    archaic fragment for each 1kb segment with the following command:

574

575     `bedtools intersect -c -a joined_regx_1kb.bed -b ind1_regx.bed ind2_regx.bed … indN_regx.bed >`
576     `freq_regx.bed`

577

578    Extended Figure 3 shows a summary of the joined region fragments, shared joined region

579    sequence and the archaic frequency for each region.

580

581

582    <u>Shared joined region fragments filtering by archaic affinity</u>

583

584    The collapsed East Asian archaic sequence (916,369,000 bp) is 1,06 fold greater than the

585    collapsed West Eurasian archaic sequence (866,945,000 bp) and more than half of the

586    sequence is shared between the two (485,255,000 bp, Table S5). We partially attribute this

587    difference to the fact that East Asians have a higher Denisova component than West

588    Eurasians [31]. To study that we repeated the analysis above filtering archaic fragments in each

589    individual (before collapsing) depending on which of the three archaic genomes (Vindija

590    Neanderthal genome [12], Altai Neanderthal genome [33], Denisova genome [34]) share the most

591    variants to (below), following the methods in [9]. Some fragments do not share variants with any

592    of the 3 sequenced archaic genomes, and thus we classify them as unknown. There are also

593    instances in which an archaic fragment does not share more SNPs with one of the archaic

594    genomes but multiple, so we can't classify the affinity of the fragments; these fragments are

595    called ambiguous fragments.

596

597      1)  Denisova fragments
598

599    We only include archaic fragments which share more variants to Denisova genome than any

600    of the two Neanderthal genomes.

601

602      2)  nonDenisova fragments
603

604    In this analysis we exclude fragments used above from all the fragments. Thus, we include

605    Vindija-like, Altai-like, ambiguous and unknown.

606

607      3)  Neanderthal fragments

608

609   We only include archaic fragments that share more variants with either the Altai Neanderthal
610   or the Vindija Neanderthal genomes than the Denisova genome. Neanderthal ambiguous
611   fragments, fragments that share the same number of SNPs with Vindija or Altai but this number
612   is higher than what is shared with the Denisova, are also included.

613

614   All results for the different filters are shown in Table S5. The Denisova content is 3 times
615   greater in East Asia than in West Eurasia (Denisova fragments filter). When this unequal
616   component is removed (non-Denisova fragments filter), we can see that the collapsed archaic
617   sequence is very similar between the two regions.

618

619   The analysis was repeated with fragments that share more variation with Neanderthal than
620   with Denisova (Neanderthal fragments). In this case, we observe a 1.07 fold higher
621   Neanderthal content in the East Asian group. We attribute this to the fact that since West
622   Eurasia archaic fragments tend to be shorter, they do not contain enough SNPs to classify
623   them to the category that they belong to. Thus, they are going to be more often classified as
624   unknown compared to fragments in East Asia. Furthermore, the [11] method has higher false
625   negative rate with short fragments, which will artificially decrease the total number of
626   fragments in that region.

| Region | Number of samples | Type | Archaic Sequence (kb) |
|---|---|---|---|
| West Eurasia | 45 | Shared | 485,255 (55,97%) |
| | | Private | 381,690 (44,03%) |
| | | All | 866,945 (100%) |
| East Asia | 45 | Shared | 485,255 (52,95%) |
| | | Private | 431,114 (47,05%) |
| | | All | 916,369 (100%) |

627

628 **Table S3**. Summary table of shared, private and total joined archaic sequence of West Eurasia
629 and EastAsia regions. Percent in respect of the total are shown in parenthesis.

630

| Region | Number of samples | Type | Number archaic fragments | | Archaic seq (bp) | | Archaic fragment length (bp) | |
|---|---|---|---|---|---|---|---|---|
| | | | mean | SE | mean | SE | mean | SE |
| West Eurasia | 45 | Shared | 756.17 | 9.43 | 59,867,254.64 | 945,255.36 | 79,060.32 | 473.90 |
| | | Private | 221.48 | 2.83 | 11,477,725.85 | 239,923.71 | 51,722.32 | 737.11 |
| | | All | 977.85 | 9.83 | 71,357,270.60 | 988,585.16 | 72,879.08 | 447.28 |
| East Asia | 45 | Shared | 913.80 | 5.04 | 81,720,490.56 | 586,860.89 | 89,448.78 | 473.32 |
| | | Private | 247.80 | 3.72 | 13,824,416.74 | 249,255.72 | 55,785.00 | 573.58 |
| | | All | 1161.57 | 6.76 | 95,555,705.32 | 704,807.56 | 82,258.99 | 400.50 |

631

632 **Table S4**. Summary statistics of the shared, private and total individual archaic fragments of
633 West Eurasians and East Asians. For each statistic, the mean and the of SE (S1) is provided.

634

635

| | Joined East Asia archaic sequence (kb) | Joined West Eurasia archaic sequence (kp) | Fold diff | Shared joined archaic sequence (kb) | East Asia shared (%) | West Eurasia shared (%) |
|---|---|---|---|---|---|---|
| All fragments | 916,369 | 866,945 | 1.06 | 485,255 | 52.95 | 55.97 |
| Denisova fragments | 107,695 | 36,850 | 2.92 | 16,004 | 14.86 | 43.43 |
| nonDenisova fragments | 853,065 | 850,028 | 1.003 | 460,490 | 53.98 | 54.17 |
| Neanderthal fragments | 646,710 | 604,518 | 1.07 | 309,043 | 47.79 | 51.12 |

636

637 **Table S5**. Joined archaic sequence in East Asia and West Eurasia and comparative statistics
638 for different subsamples of archaic fragments (S6).

639 **7** - Derived alleles call outside regions with evidence of archaic introgression and acquired

640 after the Out of Africa in SGDP samples

641

642 We retrieved the genotypes of all polymorphic loci for each individual in the 5 main regions

643 and African samples using the cpoly script from the Ctools software[3] for chromosomes 1 - 22.

644 In the parameter file, we specified the minimum quality to be 1 (as recommended by[3]) and

645 alleles to be polarized with the chimpanzee reference genome (PanTro2) provided with the

646 SGDP data.

647

648 Next, we masked repetitive regions and regions of the genome in which there is some

649 evidence of archaic introgression.

650

651     1) Neandertal introgressed regions

652

653 Neanderthals had a different mutation profile than modern humans[9]. Thus, differences in

654 Neanderthal content per individual could influence those analyses that explore the mutation

655 spectrum differences among populations. Also, by removing these regions, we will base the

656 mutation analysis on regions of the genome that we haven't explored in the archaic fragment

657 length part of the study. Thus, the tests are going to be independent of each other.

658

659 To do that, we disregarded any polymorphism localized in a region with evidence of archaic

660 introgression in any of the individuals analyzed in this study (S3). For that, we joined all archaic

661 fragments called in any individual included in this study using this command:

662
663 ```
bedtools merge -i ind1.bed ind2.bed … indN.bed > joined.bed
```

664
665 where `N` denotes the total number of individuals.

666
667 In total, the joined archaic region adds up to 1,632,776,000 bp.

668

669     2) Repeats

670

671 We also excluded repetitive regions in which sequencing errors are expected to be more

672 prevalent. For that, we downloaded the human reference genome by using the following

673 command:

674
675 ```
for chr in `seq 1 22` X Y;
```
676 ```
do
```
677 ```
rsync -avzP
```

678  `rsync://hgdownload.cse.ucsc.edu/goldenPath/hg19/chromosomes/chr${chr}.fa.gz .;`
679  `done`

680

681  from which we created a bed file with the coordinates of the repeats from RepeatMasker and

682  Tandem Repeats Finder (represented in the reference genomes fastas as lowercase letters

683  in the fasta file).

684

685  These regions add up to 1,431,504,380 bp in total.

686

687  The intersection between the repetitive regions and the archaic regions correspond to

688  806,042,777 bp, which corresponds to 56.31% of the total repetitive regions sequence and

689  49.37% of the archaic sequence. Together, these regions add up to 2,258,237,603 bp. If we

690  consider only the callable fraction - instead of the total genomic length of 3,036,303,846 bp -

691  of the human genome (2,835,673,565 bp), 577,435,962 bp remain after masking by archaic
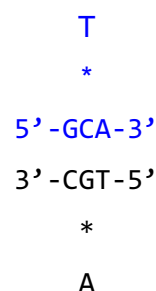
692  and repetitive regions (20.36%).

693

694  Other filters on the SNP level were imposed for each polymorphism:

695

696  1) The SNP must be biallelic

697  2) The contiguous 5' and 3' base pairs of the focal SNP (context) must be called in the
698  human reference genome (hg19)

699  3) 20% of the individuals have to be called

700  4) The chimpanzee reference genome in human coordinates must have the homologous
701  base pair called for that position

702  5) No Sub Saharan African (which excludes S_Mozabite-1, S_Mozabite-2, S_Saharawi-
703  1 and S_Saharawi-2 samples from the African supergroup) samples can have the
704  derived allele

705

706  The latter filter ensures that the polymorphisms investigated most probably arose after the Out

707  of Africa expansion. S_Masai-1, S_Masai-2 and S_Somali-1 samples are not included in the

708  Sub Saharan African group because they are reported to have some West Eurasian genetic

709  component in [3], which would affect our results. If African genomes with West Eurasian

710  components are included in the African set, then, by the 5) filter, we are going to more likely

711  remove derived alleles private to West Eurasia than other regions.

712  Homozygous locus for the derived allele count as 2 mutations and individuals heterozygous

713  count as 1 for a given individual. The distribution of derived allele accumulation per region is

714     shown in Fig. 3 and the mean derived allele accumulation counts per region are provided in

715     Table S6.

716

717     Finally, we classified loci in different mutation types depending on the derived allele nucleotide,

718     the ancestral allele nucleotide and their 5' and 3' nucleotide context. For example, as shown

719     by the diagram below, a derived allele T that had an ancestral allele C with the context G and

720     A (5' and 3' respectively) would be denoted as GCA>T. Because we do not make distinction

721     of the strand in which the mutation occurred, we collapsed strand-symmetric mutations. This

722     is the same as saying that GCA>T is equivalent to TGC>A. This way, we end up with 96

723     mutation types.

724

725                                        T

726                                      *

727                              `5'-GCA-3'`

728                              `3'-CGT-5'`

729                                      *

730                                      A

731

732     Data2_mutationspectrum.txt provides the resulting counts of each individual for each mutation

733     type in each chromosome.

734

735     The mutation types investigated in this study are the 9:

736       -  6 mutation types in which only the ancestral and derived allele nucleotides were taken

737           into account and C and T were used as ancestral (T>A, T>C, T>G, C>A, C>T, C>G)

738       -  C>T mutations were further divided into 3 mutation types:

739           -  CpG>TpG mutations which are shown to evolve in a more clock like manner[24].

740           -  TCC>TTC mutations which are in excess in Europeans compared to other

741               human populations[5,22].

742           -  C>T' mutations which contain the rest of C>T mutations not included in the

743               previous 2 types.

744

745     The distribution of derived allele accumulation per region is shown in Fig. S2 and the mean

746     derived allele accumulation counts per region are provided in Table S7.
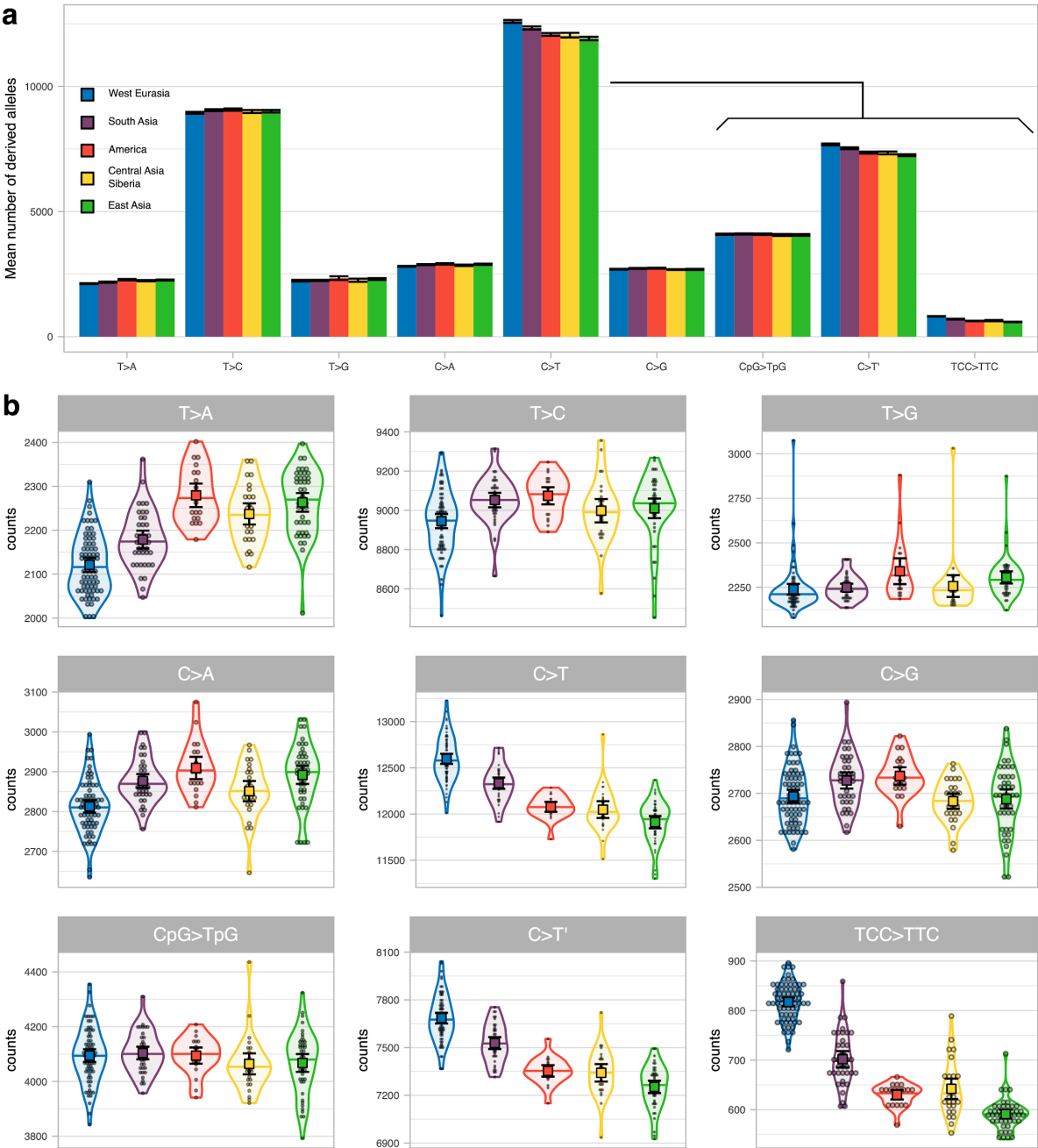
**Fig. S2**. **Mean derived allele accumulation of the 9-mutation types per region. a)** The mean number of derived alleles of each mutation type accumulated among individuals of the 5 regions (colour coded). The 95%CI of each mean is shown as error bars. **b)** The number of derived alleles of each mutation type per region (colour coded) as violin plot. Individual values are shown as dots. The median is shown as a horizontal line in each violin plot. The mean and its 95%CI of each distribution is shown as a coloured square with their corresponding error bars.

| Region | Number of samples | Derived allele accumulation | |
|---|---|---|---|
| | | mean | SE |
| West Eurasia | 71 | 31,408.54 | 62.61 |
| South Asia | 39 | 31,418.28 | 56.07 |
| America | 20 | 31,417.64 | 83.31 |
| Central Asia Siberia | 27 | 31,074.21 | 88.53 |
| EastAsia | 45 | 31,069.77 | 80.99 |

755

756 **Table S6. Derived allele accumulation per region.** Summary statistics of the derived allele
757 accumulation per region (S7). For each region, the mean and the of SE (S1) is provided.

| Region | T | | | | | |
|--------|----|----|----|----|----|----|
| | T>A | | T>C | | T>G | |
| | mean | SE | mean | SE | mean | SE |
| West Eurasia | 2,121.02 | 8.11 | 8,944.94 | 18.57 | 2,238.91 | 15.64 |
| South Asia | 2,178.71 | 10.37 | 9,052.32 | 18.91 | 2,249.05 | 10.88 |
| America | 2,279.45 | 13.67 | 9,073.85 | 22.37 | 2,340.34 | 37.02 |
| Central Asia Siberia | 2,237.11 | 12.33 | 8,997.90 | 30.43 | 2,257.22 | 31.10 |
| EastAsia | 2,263.41 | 10.86 | 9,009.86 | 25.46 | 2,305.85 | 17.52 |

758

| Region | C | | | | | |
|--------|----|----|----|----|----|----|
| | C>A | | C>T | | C>G | |
| | mean | SE | mean | SE | mean | SE |
| West Eurasia | 2,812.88 | 7.93 | 12,597.01 | 27.81 | 2,693.51 | 6.96 |
| South Asia | 2,876.76 | 8.80 | 12,333.87 | 31.10 | 2,727.42 | 8.86 |
| America | 2,909.30 | 14.36 | 12,077.68 | 27.39 | 2,737.04 | 9.47 |
| Central Asia Siberia | 2,851.11 | 13.00 | 12,047.72 | 46.23 | 2,683.31 | 8.42 |
| EastAsia | 2,891.90 | 11.58 | 11,911.38 | 34.26 | 2,687.80 | 10.43 |

759

760

761

762

763

764

765

766

767

| Region | C | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | CpG>TpG | | T>C' | | TCC>TTC | |
| | mean | SE | mean | SE | mean | SE |
| West Eurasia | 4,094.57 | 11.82 | 7,684.86 | 17.14 | 817.35 | 4.42 |
| South Asia | 4,103.80 | 11.75 | 7,528.27 | 19.06 | 701.81 | 8.40 |
| America | 4,094.27 | 14.91 | 7,353.29 | 18.45 | 629.95 | 4.81 |
| Central Asia Siberia | 4,064.23 | 19.55 | 7,341.42 | 28.09 | 642.05 | 10.48 |
| EastAsia | 4,067.34 | 16.60 | 7,252.96 | 19.80 | 591.16 | 4.48 |

768

769 **Table S7. Derived allele accumulation per region stratified per mutation type.** Summary
770 statistics of the derived allele accumulation per region for each mutation type (S7). For each
771 region and mutation type, the mean and the of SE (S1) is provided.

**8 - Estimation of the different parental generation time in West Eurasia and East Asia**

As described in the main text, West Eurasia individuals have accumulated 1.09% more derived alleles than East Asians since the split with Africans (Out of Africa). Because we are only interested in the proportion of derived alleles accumulated after the split of West Eurasians and East Asians, we need to correct for the span of time since the Out of Africa event until the split of the two Eurasian populations (Fig. S3). Thus, we need to assume dates for the split between Africans and non-Africans and the split between Eurasians.

We note that in the literature dating the Out of Africa is widely discussed and controversial, since it was not a clean split between non-Africans and Africans. Instead, from MCMC results and cros coalescence rate analysis in [2,19] the authors note that there might have been a gradual separation among African populations and between Africans and non-Africans. They suggest that this process created population structure between 200,000 - 100,000 years ago within Africa and that the non-African group had more gene flow with certain African groups (i.e., Yorubans) than others (i.e., San). After that, the rate increased, indicative of an accelerated split between Africans and non-Africans which has the median divergence point between 80,000 - 60,000 years ago. Similarly, the split among Eurasians was not clean either. All splits started around 70,000 years ago with a median divergence point between 40,000 and 20,000 years ago for East Asians and West Eurasians. Nonetheless, studies of ancient DNA show that around 40,000 years ago East Asians and West Eurasians were already diverging: the ancient human sample of Kostenki (36,000 year old sample) presents higher affinity to present day West Eurasians [21] and Tianyuan (40,000 year old sample) to East Asians [20].

In this analysis, we assume that the split between Africans and non-Africans happened 60,000 years ago and that the split between West Eurasians and East Asians happened 40,000 years ago. This is because if the proportion of time the West Eurasians and East Asians were apart decreases in respect of the time since both splited from Africans (i.e., out of Africa happening 80,000 instead of 60,000 years ago), the rate at which mutations should have accumulated would have been higher. Thus, a conservative measurement will be assuming a lower bound for the out of Africa.

In consequence, the excess of derived alleles accumulated in West Eurasians compared to East Asia is:

808
$$1.09\% \cdot \frac{60{,}000}{40{,}000} = 1.64\%$$

809

810    In [4] a poisson regression is derived for the number of mutations transmitted in each generation

811    from trio data for each parental lineage depending on their age at reproduction:

$$\hat{\mu}_{f,g} = 6.05 + 1.51a_f \tag{1}$$

$$\hat{\mu}_{m,g} = 3.61 + 0.37a_m \tag{2}$$

812    Where subscripts $f$ and $m$ denote paternal and maternal respectively, $\hat{\mu}$ is the estimation of

813    the mean mutation rate per generation $(g)$ and $a$ is the mean parental age. Thus, assuming

814    the same mean parental age for both progenitors $(a_f = a_m = a)$ we get that the total mutation

815    rate per generation is calculated by the equation 3 and the yearly $(y)$ rate by equation 4.

$$\hat{\mu}_g = \hat{\mu}_{f,g} + \hat{\mu}_{m,g} = 9.66 + 1.88a \tag{3}$$

$$\hat{\mu}_y = \hat{\mu}_g/a \tag{4}$$

816

817    Then, to compare the mutation rate per year in two different populations ($x$ and $z$) with different

818    mean parental ages, we get that

819

$$\frac{\hat{\mu}_{yx}}{\hat{\mu}_{yz}} = \frac{\frac{9.66}{a_x} + 1.88}{\frac{9.66}{a_z} + 1.88} \tag{5}$$

820    The number of derived alleles accumulated in a genome during a period of time $(d)$ depends

821    on the mutation rate per year and the time span $(T)$

$$d = \mu_y T \tag{6}$$

822    However, the ratio of $d$ between two populations, will only depend on their mutation rate

823    because $T$ has been the same for both

824

$$\frac{d_x}{d_y} = \frac{\hat{\mu}_{gx}}{\hat{\mu}_{gy}} \tag{7}$$

825

826    Thus, we can estimate the $a_x$ if $a_z$ and the $d_x/d_z$ are known
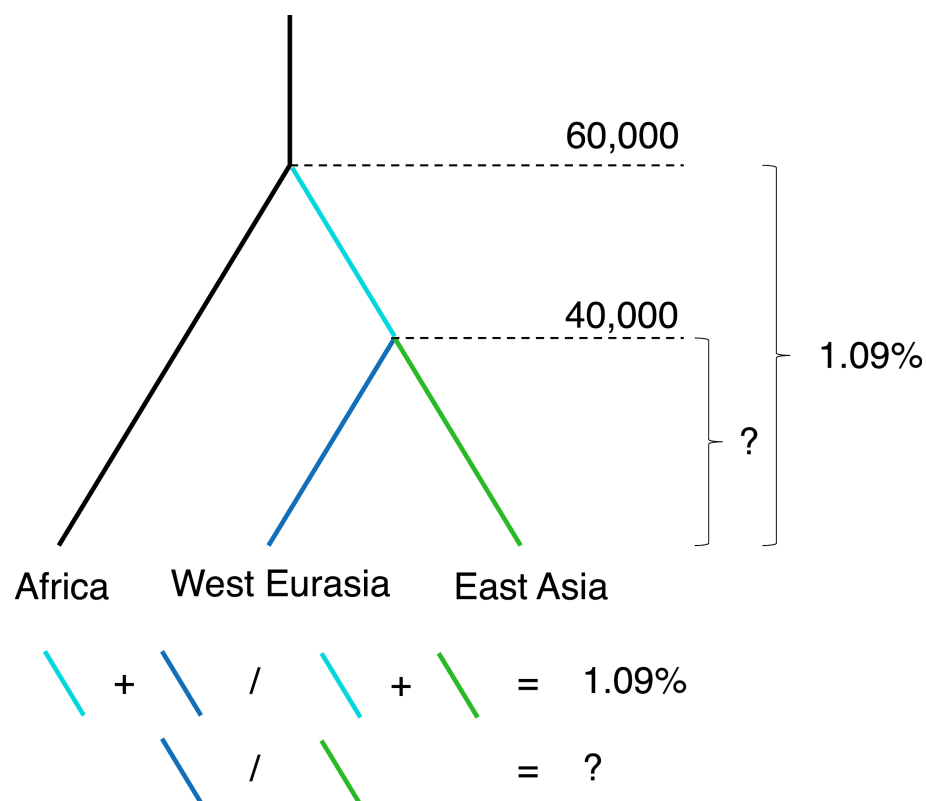
827

$$\frac{\hat{\mu}_{gx}}{\hat{\mu}_{gz}} = \frac{d_x}{d_z} = \frac{\frac{9.66}{a_x} + 1.88}{\frac{9.66}{a_z} + 1.88}$$

$$a_x = \frac{9.66}{\frac{d_x}{d_z}\left(\frac{9.66}{a_z} + 1.88\right) - 1.88} \tag{8}$$

828

829    In this study, we find that the ratio of the mean derived allele accumulation in West Eurasia

830    ($WE$) vs East Asia ($EA$, $d_{WE}/d_{EA}$) is 1.0164 (1.64%). With formula (8), we check for reasonable

831    $a_{EA}$ values between 28 and 32 years and found that the values of $a_{WE}$ ranged between 25.32

832    and 28.59 respectively. Thus, we estimate that generation times in East Asians have been

833    2.68 to 3.39 years longer than in West Eurasians since the split of the two populations. This

834    corresponds to West Eurasians having had approx. 150 generations more than East Asians.

835



836
837
838 **Fig. S3**. **Mutation rate difference between West Eurasia and East Asia.** This diagram
839 shows conceptually that the mutation rate could only be different after the split between East
840 Asians and West Eurasians (blue and green terminal branches). However, the difference in
841 derived allele accumulation is calculated since the split with Africans for each group (cyan and
842 blue, cyan and green).

**9 - Mutation spectrum correlation with mean parental age**

The germline mutation spectrum is dependent on the parental sex and age at conception [4]. In this study, we observe differences in the abundance of derived alleles accumulated after the out of Africa event when stratified by mutation type (Fig. S2, Table S7). Here, we study to which extent these differences can be explained by changes in generation time in the 5 regions. For that, we compare the mutational patterns of *de novo* mutations (DNM) depending on parental age in trio studies [4,25] (deCODE dataset) with the differences in mutation spectrum of extant populations with the mean archaic fragment length as a proxy of mean generation time (SGDP dataset).

SGDP dataset

We classified the derived alleles found in the autosomes of each individual into 6 mutation types depending on the ancestral and derived allele as explained in S7. C>T mutations were also classified in 3 subtypes: TCC>TTC, CpG>TpG and the rest (C>T'). In total, we divide all mutations into 9 types. In order to obtain the fraction of each mutation type per individual, we divided the number of each mutation type by the total amount of derived alleles. C>T mutations are duplicated since we subdivide them into 3 extra categories (TCC>TTC, CpG>TpG and C>T'). Thus, the total amount of derived alleles do not consider these 3 types. We correlated the fraction of derived alleles of each type with the mean archaic fragment length as a proxy of mean generation time (Fig. 3b). We obtained the linear model of such correlation for each mutation type using the following R function (Table S8).

```
lm(mutation_fraction~mean_fragment_length)
```

deCODE dataset

We downloaded the set of DNM called in [25] and the additional proband information from the supplementary data provided in the publication. We join both in order to compute the mean parental age for each DNM for each proband. Indels are filtered out. Following the methodology in a similar test in [4], we aggregate all mutation counts for each of the 9 types of all probands with the same mean parental age. We then compute the fraction of each mutation type. In other words, for each mutation type and mean parental age we have a single mutation fraction value. Those data points that were obtained aggregating information from less than 2 probands were discarded. We obtained linear models for each mutation type using the following R function (Table S8).

880

881   `lm(mutation_fraction~mean_parental_age, weights = n_probands)`

882

883   The correlations between the slopes of both datasets is shown in Fig. S4.

884

885   The probands of the deCODE dataset have a bias towards fathers being older than mothers,

886   with a mean of 2.77 years and the largest difference of more than 40 years (Fig. S5a). To

887   study if the correlation of mutation spectrum with the mean parental age is affected by the

888   mentioned bias, we rerun the correlation test with the deCODE dataset with only probands

889   that have parents with an age difference of less than 4 years.  This way, we retaining more

890   than 50% of the data (Fig. S5b) and reduce the bias (mean = 0.94 differences in years,

891   Fig. S5c). We then compared the slopes of the linear models calculated in the original

892   deCODE dataset and when we impose the parental age difference filter explained above

893   (Fig. S6). We don't observe qualitative changes in the slopes when comparing the two and

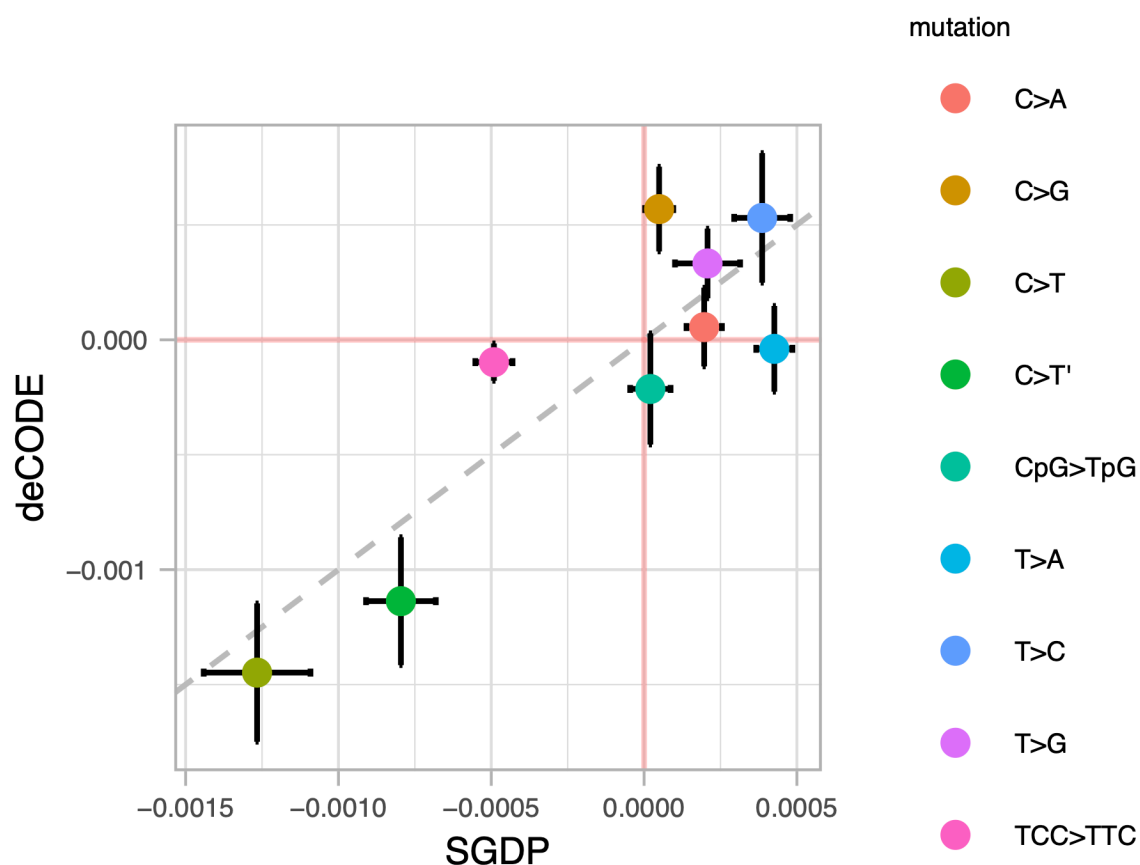894   thus, we used all probands for our analysis.

895

896

**Fig. S4**. **Slope coefficient correlation between SGDP data and deCODE data linear models.** Dot plot graph illustrating the correlation between linear model slope coefficients derived from the SGDP data (x-axis) and deCODE data (y-axis) for each mutation type (color code). 95%CI for each estimate are shown as error bars. The 1-to-1 correspondence is denoted by the gray dashed diagonal line.
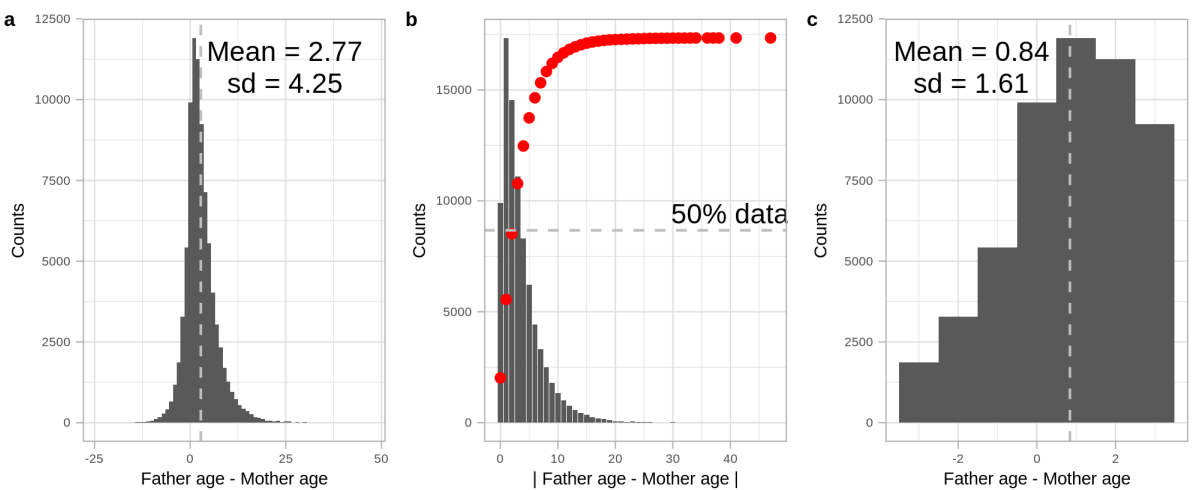
**Fig. S5**. **Parental age difference in the deCODE data. a)** Histogram of the number of probands with a certain parental age difference. The mean is shown as a vertical gray line and annotated as a numeric figure. **b)** Histogram of the number of probands with a certain absolute parental age difference. The cumulative distribution of provands is denoted by red dots. The horizontal gray line shows the 50% data threshold. **c)** Histogram of the number of probands with a certain parental age difference with less than 4 years difference. The mean is shown as a vertical gray line and annotated as a numeric figure.
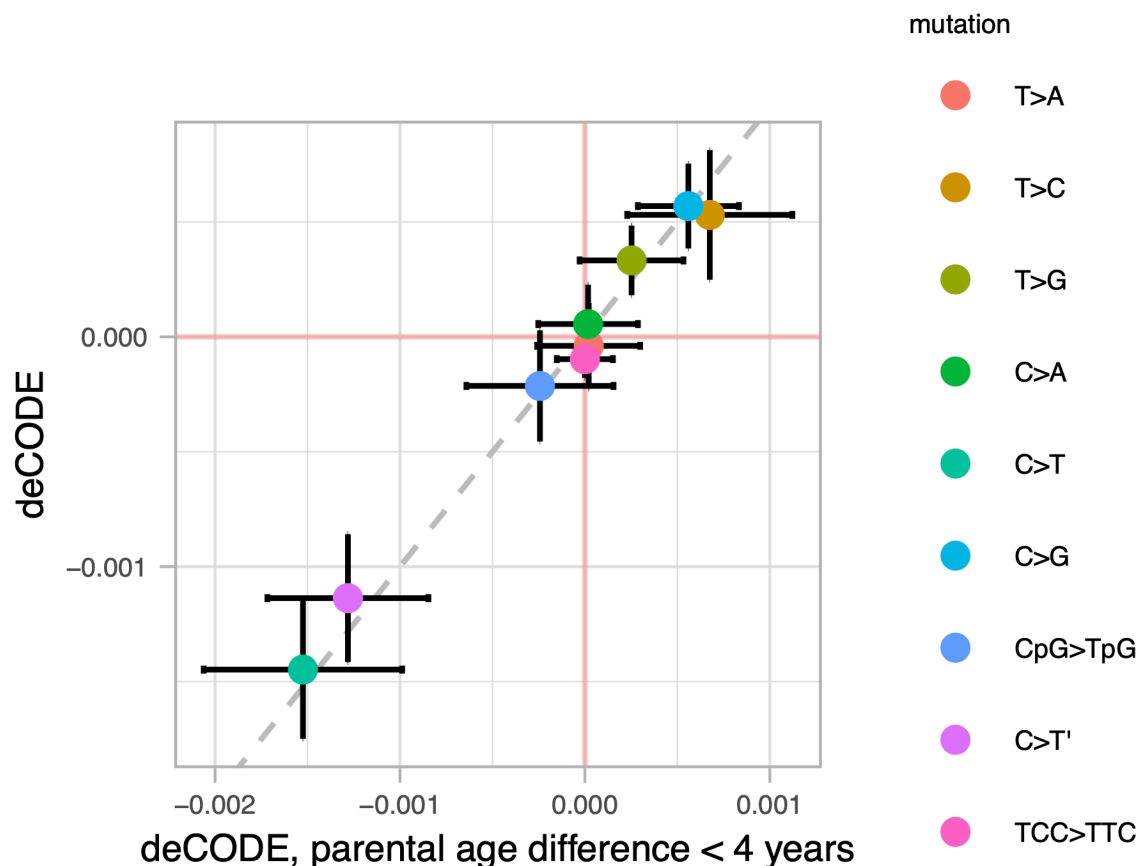
44

**Fig. S6**. **Slope coefficient correlation between linear models of deCODE data and deCODE data when only using probands with parental age difference less than 4 years.** Dot plot graph illustrating the correlation between linear model slope coefficients derived from the deCODE data (y-axis) and the deCODE data when only using probands with parental age difference less than 4 (x-axis) for each mutation type (color code). 95%CI for each estimate are shown as error bars. The 1-to-1 correspondence is denoted by the gray dashed diagonal line.

| Mutation | Dataset | Intercept | SE | t vauel | P vauel |
|---|---|---|---|---|---|
| T>A | deCODE | 6.72e-2 | 2.86e-3 | 23.47 | 9.70e-34 |
| | SGDP | 3.71e-2 | 2.34e-3 | 15.85 | 3.40e-37 |
| T>C | deCODE | 2.46e-1 | 4.33e-3 | 56.79 | 9.27e-58 |
| | SGDP | 2.58e-1 | 3.63e-3 | 71.00 | 8.39e-144 |
| T>G | deCODE | 5.34e-2 | 2.33e-3 | 22.98 | 3.34e-33 |
| | SGDP | 5.64e-2 | 4.22e-3 | 13.36 | 1.66e-29 |
| C>A | deCODE | 8.79e-2 | 2.63e-3 | 33.43 | 4.26e-43 |
| | SGDP | 7.61e-2 | 2.29e-3 | 33.19 | 2.77e-83 |
| C>T | deCODE | 4.69e-1 | 4.62e-3 | 101.48 | 3.35e-74 |
| | SGDP | 4.90e-1 | 6.92e-3 | 70.86 | 1.22e-143 |
| C>G | deCODE | 7.70e-2 | 2.83e-3 | 27.21 | 1.37e-37 |
| | SGDP | 8.25e-2 | 1.83e-3 | 45.16 | 7.19e-107 |
| CpG>TpG | deCODE | 1.82e-1 | 3.71e-3 | 48.96 | 1.32e-53 |
| | SGDP | 1.29e-1 | 2.62e-3 | 49.30 | 7.50e-114 |
| C>T' | deCODE | 2.65e-1 | 4.27e-3 | 62.02 | 3.08e-60 |
| | SGDP | 3.01e-1 | 4.52e-3 | 66.54 | 2.14e-138 |
| TCC>TTC | deCODE | 2.19e-2 | 1.25e-3 | 17.54 | 1.55e-26 |
| | SGDP | 6.05e-2 | 2.41e-3 | 25.13 | 8.48e-64 |

918

919

920

921

922

923

924

925

| Mutation | Dataset | Slope | SE | t vaue | P vaue |
|---|---|---|---|---|---|
| T>A | deCODE | -3.92e-5 | 9.52e-5 | -0.41 | 6.82e-1 |
|  | SGDP | 4.26e-4 | 3.02e-5 | 14.13 | 7.04e-32 |
| T>C | deCODE | 5.30e-4 | 1.44e-4 | 3.69 | 4.61e-4 |
|  | SGDP | 3.86e-4 | 4.68e-5 | 8.27 | 1.91e-14 |
| T>G | deCODE | 3.32e-4 | 7.73e-5 | 4.30 | 5.77e-5 |
|  | SGDP | 2.08e-4 | 5.44e-5 | 3.82 | 1.78e-4 |
| C>A | deCODE | 5.55e-5 | 8.74e-5 | 0.63 | 5.28e-1 |
|  | SGDP | 1.97e-4 | 2.95e-5 | 6.66 | 2.60e-10 |
| C>T | deCODE | -1.45e-3 | 1.54e-4 | -9.43 | 7.52e-14 |
|  | SGDP | -1.27e-3 | 8.91e-5 | -14.21 | 3.85e-32 |
| C>G | deCODE | 5.69e-4 | 9.41e-5 | 6.05 | 7.66e-8 |
|  | SGDP | 4.92e-5 | 2.35e-5 | 2.09 | 3.77e-2 |
| CpG>TpG | deCODE | -2.14e-4 | 1.23e-4 | -1.73 | 8.82e-2 |
|  | SGDP | 2.07e-5 | 3.37e-5 | 0.61 | 5.41e-1 |
| C>T' | deCODE | -1.14e-3 | 1.42e-4 | -8.01 | 2.58e-11 |
|  | SGDP | -7.96e-4 | 5.82e-5 | -13.67 | 1.81e-30 |
| TCC>TTC | deCODE | -9.72e-5 | 4.15e-5 | -2.34 | 2.22e-2 |
|  | SGDP | -4.91e-4 | 3.10e-5 | -15.84 | 3.79e-37 |

926

927 **Table S8**. **Linear models between mutation type fraction and mean generation time**
928 **estimate in the SGDP and deCODE data sets.** Two separate tables are given for the
929 intercept and the slope of the linear models. For each mutation type and data set, the
930 coefficients estimate the SE, t value and the associated P value are provided.

931    **10 - Sex Specific mutational patterns**

932

933    X-to-A ratio

934

935    Due to the inheritance pattern of the X chromosome - 2 copies transmitted in females while

936    only 1 in males - compared to autosomes - 2 copies in both females and males -, it is expected

937    that the X chromosome has ¾ the diversity of the autosomes. However, this can be altered if

938    the mutation rate changes disproportionately between females and males due to shifts in

939    generation time between sexes. For example, an increase in the male mean generation time

940    will decrease the yearly mutation rate in males and thus, proportionally less mutations are

941    going to be accumulated in autosomes compared the X chromosomes [28]. Therefore, the ratio

942    of derived allele accumulation between the X chromosome and the autosomes will reflect

943    variation on the generation time between males and females: higher values of the X-to-A ratio

944    will be indicative of longer generation times in males compared to females and vice versa.

945    Although here we only consider generation time differences to affect the ratio, there are other

946    factors that can perturbe this ratio such as reproductive variance between sexes [35],

947    demographic events [36] or differences in selection [37].

948

949    To investigate that, we obtained the number of derived alleles in the autosomes and X

950    chromosomes of the females of the SGDP data (Table S9), as described in S7 (included in

951    Data2_mutationspectrum.txt), and computed the X-to-A ratio as:

952

953    $$\frac{d_X}{L_X} \Big/ \frac{d_A}{L_A}$$

954

955    where $d$ denotes the number of derived alleles, $L$ the number of callable base pairs in either

956    $X$(X chromosome) or $A$(autosomes). We then correlated the ratio with the mean archaic

957    fragment length for each individual obtained in S3 (Fig. 4a).

958

959    C>G maternaly enriched regions

960

961    As described in [4], there are regions of the genome in which DNM are clustered (cDNM). Those

962    regions appear to be enriched in C>G mutations which originate in the maternal lineage. They

963    also show that these clusters increase in number more rapidly with maternal than paternal age

964    at conception.

965

966 Here we explore if there is a difference on the number of C>G segregating sites in cDNM
967 genomic windows among the 5 regions.

968

969 For that we compute the number of derived alleles that are C>G and non-C>G along the
970 genome in windows of 1Mb. We join this information with the annotation of 1Mb-window of the
971 genome as cDNM or non-cDNM provided in [4]. Then, for each individual we compute the
972 following:

973

$$p = \frac{d_{C>G}}{d_{non-C>G}}$$

974

975

976 where $d$ denotes the number of derived alleles of C>G or non-C>G. Thus, p is the ratio
977 between the two quantities. Then, to compare this ratio between cDNM and non-cDNM
978 regions we compute the mean $p$ ($\underline{p}$) over all regions and compute the following ratio

979

$$r = \frac{\bar{p}_{C>G}}{\bar{p}_{non-C>G}}$$

980

981

982 If $r = 1$, it shows that there are a similar number of C>G mutations in cDNM regions compared
983 to the rest of the genome. If $r > 1$, then there is an excess and if $r < 1$, then there is a depletion.
984 Nonetheless, we are not interested in the actual ratio, but the comparison among regions on
985 this quantity. We then correlated the ratio with the mean archaic fragment length for each
986 individual obtained in S3 (Fig. 4b).

987

988 <u>Y chromosome</u>

989

990 Male individuals with shorter generation time are predicted to increase the mutation rate per
991 year. Thus, Y chromosomes are expected to accumulate more derived alleles in individuals
992 with a historically shorter mean generation time compared to others with longer ones.

993

994 To investigate that, we followed a similar procedure as in S7, changing certain steps and filters
995 listed below:

996

997     1. We only used males in SGDP data
998     2. Alleles were polarized using the Chimp sequence in human coordinates. Since the
999        chimpanzee Y chromosome is not provided with the SGDP data, this was achieved by

1000    taking the chimpanzee sequence from the hg19-panTro6 alignment into a fasta file

1001    with the human coordinates. The alignment can be downloaded from the following link:

1002

1003    http://hgdownload.cse.ucsc.edu/goldenpath/hg19/vsPanTro6/reciprocalBest/axtRBestNet/hg

1004    19.panTro6.rbest.axt

1005

1006    3.  No archaic regions were masked since there is no evidence of archaic sequence in

1007        the modern human Y chromosome

1008    4.  Only polymorphisms in the X degenerate regions are considered (coordinates from [38])

1009        and no further filters regarding repetitive regions were imposed

1010    5.  Individuals S_Finnish-2, S_Finnish-3, S_Palestinian-2, S_Mansi-1 and S_Masai-2 were

1011        discarded from the analysis because they didn't yield any callable polymorphism

1012    6.   For each individual, all heterozygous sites were classified as non callable sites

1013    7.  Only African individuals with Y haplogroups A and B (metadata provided in [3], A:

1014        S_Ju_hoan_North-2, S_Dinka-2; B: S_Biaka-1, S_Biaka-2, S_Mbuti-3, S_Ju_hoan_North-

1015        3, S_Ju_hoan_North-1) were used as the outgroup. If polymorphisms were found to be

1016        segregating in these individuals, they were filtered out from this analysis

1017    8.  We didn't require the 5' and 3' contiguous base pairs (context) of a polymorphic site to be

1018        callable

1019

1020    The accumulation of derived alleles in the Y chromosome per geographical region is shown in

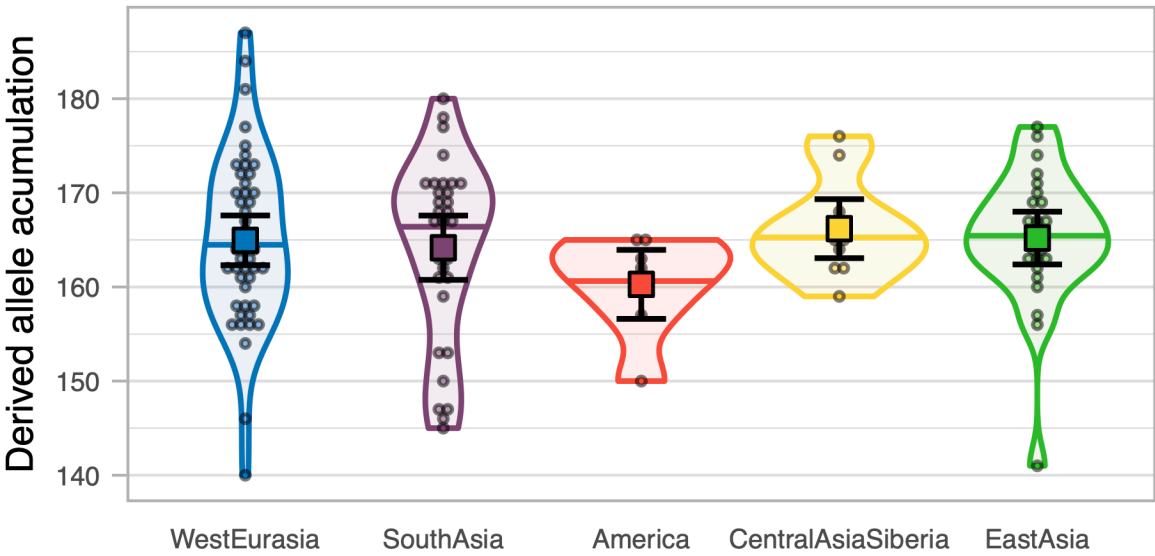1021    Fig. S7 (included in Data2_mutationspectrum.txt) and in Table S10.

1022



1023

1024

1025 **Fig. S7**. **Mean derived allele accumulation of the 7-mutation types per region in the Y**
1026 **Chromosome. a)** The mean number of derived alleles of each mutation type accumulated
1027 among individuals of the 5 regions (colour coded). The 95%CI of each mean is shown as error
1028 bars. **b)** The number of derived alleles of each mutation type per region (colour coded) as
1029 violin plot. Individual values are shown as dots. The median is shown as a horizontal line in
1030 each violin plot. The mean and its 95%CI of each distribution is shown as a coloured square
1031 with their corresponding error bars.

| Region | Number of samples | Derived allele accumulation (X chromosome) | |
| --- | --- | --- | --- |
| | | mean | SE |
| West Eurasia | 23 | 2,820.21 | 21.17 |
| South Asia | 8 | 2,911.41 | 26.69 |
| America | 13 | 2,839.21 | 21.73 |
| Central Asia Siberia | 16 | 2,818.77 | 18.56 |
| East Asia | 20 | 2,900.35 | 17.46 |

**Table S9. Derived allele accumulation per region for the X chromosome in female individuals.** Summary statistics of the derived allele accumulation per region on the X chromosome of females. For each region, the mean and the of SE (S1) is provided.

| Region | Number of samples | Derived allele accumulation (Y chromosome) | |
| --- | --- | --- | --- |
| | | mean | SE |
| West Eurasia | 45 | 164.95 | 1.35 |
| South Asia | 31 | 164.17 | 1.74 |
| America | 7 | 160.28 | 1.87 |
| Central Asia Siberia | 10 | 166.20 | 1.60 |
| East Asia | 25 | 165.20 | 1.43 |

**Table S10. Derived allele accumulation per region for the Y chromosome in male individuals.** Summary statistics of the derived allele accumulation per region on the X chromosome of males. For each region, the mean and the of SE (S1) is provided.

**11 - Datasets**

Data1_archaicfragments.txt: Archaic fragments found in individuals from the 5 main geographical regions and ancient samples in the SGDP investigated in this study. Each line is a fragment with the following attributes:

1. name: individual the fragment belongs to.
2. region: region that the individual belongs to as defined by [3].
3. chrom: chromosome in which the fragment is located.
4. start: starting fragment position in hg19 coordinates.
5. end: ending fragment position in hg19 coordinates.
6. length: fragment length (end - start).
7. MeanProb: mean posterior probability for the fragment outputted by the [11] method.
8. snps: number of SNPs found in the fragment that are not segregating in any of the Sub Saharan African genomes (S3).
9. Altai: number of SNPs found in the fragment that are shared with the Altai Neanderthal [33].
10. Denisova: number of SNPs found in the fragment that are shared with the Denisova [34].
11. Vindija: number of SNPs found in the fragment that are shared with Vindija Neanderthal [12].

Data2_mutationspectrum.txt: Counts of derived alleles classified into the 96 mutation types for the extant samples of the SGDP, per chromosome. Each line has the following attributes:

1. ind: individual identifier
2. reg: region that the individual belongs to as defined by [3].
3. sex: individual sex defined by [3]. M = male, F = female.
4. chrom: chromosome which the counts belong to.
5. fiv: contiguous 5' base pair of the focal SNP
6. anc: ancestral allele of the mutation
7. thr: contiguous 3' base pair of the focal SNP
8. der: ancestral allele of the mutation
9. counts: number of mutation types found

**Bibliography**

1. Green, R. E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**, 710–722 (2010).

2. Bergström, A. *et al.* Insights into human genetic variation and population history from 929 diverse genomes. doi:10.1101/674986.

3. Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature* **538**, 201–206 (2016).

4. Jónsson, H. *et al.* Parental influence on human germline de novo mutations in 1,548 trios from Iceland. *Nature* **549**, 519–522 (2017).

5. Harris, K. & Pritchard, J. K. Rapid evolution of the human mutation spectrum. *Elife* **6**, (2017).

6. Villanea, F. A. & Schraiber, J. G. Multiple episodes of interbreeding between Neanderthal and modern humans. *Nat Ecol Evol* **3**, 39–44 (2019).

7. Wall, J. D. *et al.* Higher levels of neanderthal ancestry in East Asians than in Europeans. *Genetics* **194**, 199–209 (2013).

8. Lazaridis, I. *et al.* Genomic insights into the origin of farming in the ancient Near East. *Nature* **536**, 419–424 (2016).

9. Skov, L. *et al.* The nature of Neanderthal introgression revealed by 27,566 Icelandic genomes. *Nature* (2020) doi:10.1038/s41586-020-2225-9.

10. Carlson, J., DeWitt, W. S. & Harris, K. Inferring evolutionary dynamics of mutation rates through the lens of mutation spectrum variation. *Curr. Opin. Genet. Dev.* **62**, 50–57 (2020).

11. Skov, L. *et al.* Detecting archaic introgression using an unadmixed outgroup. *PLoS Genet.* **14**, e1007641 (2018).

12. Prüfer, K. *et al.* A high-coverage Neandertal genome from Vindija Cave in Croatia. *Science* **358**, 655–658 (2017).

13. Browning, S. R., Browning, B. L., Zhou, Y., Tucci, S. & Akey, J. M. Analysis of Human Sequence Data Reveals Two Pulses of Archaic Denisovan Admixture. *Cell* **173**, 53–

1104      61.e9 (2018).

1105   14. Harris, K. & Nielsen, R. The Genetic Cost of Neanderthal Introgression. *Genetics* **203**,

1106      881–891 (2016).

1107   15. Petr, M., Pääbo, S., Kelso, J. & Vernot, B. Limits of long-term selection against

1108      Neandertal introgression. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 1639–1644 (2019).

1109   16. Fu, Q. *et al.* Genome sequence of a 45,000-year-old modern human from western

1110      Siberia. *Nature* vol. 514 445–449 (2014).

1111   17. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for

1112      present-day Europeans. *Nature* **513**, 409–413 (2014).

1113   18. Moorjani, P. & Others. Molecular clock helps estimate age of ancient genomes. *Proc.*

1114      *Natl. Acad. Sci. U. S. A.* **113**, 5459–5460 (2016).

1115   19. Schiffels, S. & Durbin, R. Inferring human population size and separation history from

1116      multiple genome sequences. *Nat. Genet.* **46**, 919–925 (2014).

1117   20. Fu, Q. *et al.* DNA analysis of an early modern human from Tianyuan Cave, China. *Proc.*

1118      *Natl. Acad. Sci. U. S. A.* **110**, 2223–2227 (2013).

1119   21. Seguin-Orlando, A. *et al.* Paleogenomics. Genomic structure in Europeans dating back

1120      at least 36,200 years. *Science* **346**, 1113–1118 (2014).

1121   22. Harris, K. Evidence for recent, population-specific evolution of the human mutation rate.

1122      *Proc. Natl. Acad. Sci. U. S. A.* **112**, 3439–3444 (2015).

1123   23. Mathieson, I. & Reich, D. Differences in the rare variant spectrum among human

1124      populations. *PLoS Genet.* **13**, e1006581 (2017).

1125   24. Moorjani, P., Amorim, C. E. G., Arndt, P. F. & Przeworski, M. Variation in the molecular

1126      clock of primates. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 10607–10612 (2016).

1127   25. Halldorsson, B. V. *et al.* Characterizing mutagenic effects of recombination through a

1128      sequence-level genetic map. *Science* **363**, (2019).

1129   26. DeWitt, W. S., Harris, K. D. & Harris, K. Joint nonparametric coalescent inference of

1130      mutation spectrum history and demography. *bioRxiv* (2020).

1131   27. Fenner, J. N. Cross-cultural estimation of the human generation interval for use in

1132   genetics-based population divergence studies. *American Journal of Physical*

1133   *Anthropology* vol. 128 415–423 (2005).

1134  28. Amster, G. & Sella, G. Life History Effects on Neutral Diversity Levels of Autosomes and

1135   Sex Chromosomes. *Genetics* (2020) doi:10.1534/genetics.120.303119.

1136  29. Goldmann, J. M. *et al.* Parent-of-origin-specific signatures of de novo mutations. *Nat.*

1137   *Genet.* **48**, 935–939 (2016).

1138  30. Consortium, T. 1000 G. P. & The 1000 Genomes Project Consortium. A global

1139   reference for human genetic variation. *Nature* vol. 526 68–74 (2015).

1140  31. Skoglund, P. & Jakobsson, M. Archaic human ancestry in East Asia. *Proc. Natl. Acad.*

1141   *Sci. U. S. A.* **108**, 18301–18306 (2011).

1142  32. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic

1143   features. *Bioinformatics* **26**, 841–842 (2010).

1144  33. Prüfer, K. *et al.* The complete genome sequence of a Neanderthal from the Altai

1145   Mountains. *Nature* **505**, 43–49 (2014).

1146  34. Meyer, M. *et al.* A high-coverage genome sequence from an archaic Denisovan

1147   individual. *Science* **338**, 222–226 (2012).

1148  35. Keinan, A., Mullikin, J. C., Patterson, N. & Reich, D. Accelerated genetic drift on

1149   chromosome X during the human dispersal out of Africa. *Nat. Genet.* **41**, 66–70 (2009).

1150  36. Amster, G., Murphy, D. A., Milligan, W. M. & Sella, G. Changes in life history and

1151   population size can explain relative neutral diversity levels on X and autosomes in

1152   extant human populations. doi:10.1101/763524.

1153  37. Hammer, M. F. *et al.* The ratio of human X chromosome to autosome diversity is

1154   positively correlated with genetic distance from genes. *Nat. Genet.* **42**, 830–831 (2010).

1155  38. Skov, L., Danish Pan Genome Consortium & Schierup, M. H. Analysis of 62 hybrid

1156   assembled human Y chromosomes exposes rapid structural changes and high rates of

1157   gene conversion. *PLoS Genet.* **13**, e1006834 (2017).