



This is a self-archived – parallel published version of an original article. This version may differ from the original in pagination and typographic details. When using please cite the original.

This is a post-peer-review, pre-copyedit version of an article published in

Journal	European Journal of Human Genetics
DOI	The final authenticated version is available online at https://doi.org/10.1038/s41431-021-00897-8
CITATION	Ilumäe, AM., Post, H., Flores, R. <i>et al.</i> Phylogenetic history of patrilineages rare in northern and eastern Europe from large-scale re-sequencing of human Y-chromosomes. <i>Eur J Hum Genet</i> (2021). https://doi.org/10.1038/s41431-021-00897-8

1 **Phylogenetic history of patrilineages rare in northern and eastern Europe from large-scale**
2 **re-sequencing of human Y-chromosomes**

3 *Anne-Mai Ilumäe,^{1,2,*} Helen Post,^{1,3,*} Rodrigo Flores,¹ Monika Karmin,^{1,4} Hovhannes*
4 *Sahakyan,^{1,5} Mayukh Mondal,¹ Francesco Montinaro,^{1,6} Lauri Saag,¹ Concetta Bormans,⁷ Luisa*
5 *Fernanda Sanchez,⁷ Adam Ameer,^{8,9} Ulf Gyllensten,⁸ Mart Kals,¹⁰ Reedik Mägi,¹⁰ Luca*
6 *Pagani,^{1,11} Doron M Behar,^{1,7} Siiri Rootsi,¹ Richard Villems^{1,3}*

7 ¹ Estonian Biocentre, Institute of Genomics, University of Tartu, Tartu 51010, Estonia

8 ² Department of Biology, University of Turku, Turku 20014, Finland

9 ³ Department of Evolutionary Biology, Institute of Molecular and Cellular Biology, University of Tartu,
10 Tartu 51010, Estonia

11 ⁴ Computational Biology Research Group, School of Fundamental Sciences, Massey University,
12 Palmerston North 4474, New Zealand

13 ⁵ Laboratory of Evolutionary Genomics, Institute of Molecular Biology of National Academy of Sciences,
14 Yerevan 0014, Armenia

15 ⁶ Department of Biology-Genetics, University of Bari, Bari 70125, Italy

16 ⁷ Genomic Research Center, Gene by Gene, Houston, Texas, 77008, USA

17 ⁸ Science for Life Laboratory, Department of Immunology, Genetics and Pathology, Uppsala University,
18 Uppsala 75108, Sweden

19 ⁹ Department of Epidemiology and Preventive Medicine, Monash University, Melbourne VIC 3004,
20 Australia

21 ¹⁰ Estonian Genome Centre, Institute of Genomics, University of Tartu, Tartu 51010, Estonia

22 ¹¹ Department of Biology, University of Padova, Padova 35131, Italy

23 Corresponding author: Anne-Mai Ilumäe

24 Institute of Genomics, University of Tartu

25 23 Riia Street, Tartu 51010, Estonia

26 annemai.ilumae@ut.ee

27 *These authors contributed equally to this work

28 **Abstract**

29 The most frequent Y-chromosomal (chrY) haplogroups in northern and eastern Europe
30 (NEE) are well-known and thoroughly characterized. Yet a considerable number of men
31 in every population carry rare paternal lineages with estimated frequencies around 5%.
32 So far, limited sample-sizes and insufficient resolution of genotyping have obstructed a
33 truly comprehensive look into the variety of rare paternal lineages segregating within
34 populations and potential signals of population history that such lineages might convey.
35 Here we harness the power of massive re-sequencing of human Y chromosomes to
36 identify previously unknown population-specific clusters among rare paternal lineages
37 in NEE. We construct dated phylogenies for haplogroups E2-M215, J2-M172, G-M201
38 and Q-M242 on the basis of 421 (of them 282 novel) high-coverage chrY sequences
39 collected from large-scale databases focusing on populations of NEE. Within these
40 otherwise rare haplogroups we disclose lineages that began to radiate ~1-3 thousand
41 years ago in Estonia and Sweden and reveal male phylogenetic patterns testifying of
42 comparatively recent local demographic expansions. Conversely, haplogroup Q lineages
43 bear evidence of ancient Siberian influence lingering in the modern paternal gene pool
44 of northern Europe. We assess the possible direction of influx of ancestral carriers for
45 some of these male lineages. In addition, we demonstrate the congruency of paternal
46 haplogroup composition of our dataset with two independent population-based cohorts
47 from Estonia and Sweden.

48 **Keywords:** Human population genetics, Y chromosome variety, North-East Europe, Y
49 rare haplogroups

50 **Introduction**

51 Genetic studies investigating uniparental and fine-scale autosomal variation in Estonia
52 [1] and in its neighboring populations in NEE [2–6] observed that the regional genetic
53 structure correlates closely with geography. In addition, recent ancient DNA studies
54 have begun to uncover the settlement history of NEE, which is distinct from that of
55 central and southern parts of the continent [7–9].

56 The four most common chrY haplogroups (hgs) with incidence above 5% (R1a-M198,
57 N3-TAT, I-M170, R1b-M343) constitute over 90% of the chrY pool in NEE [3, 10–12].
58 Several studies have analysed these hgs in a wide phylogeographic context.

59 Besides the four most common hgs, several paternal lineages belonging on the basic
60 level to hgs E2, J2, G and Q with frequency up to 5%, complement the pool of Y-
61 chromosomes in NEE [3, 4, 13–15]. In Europe, hgs E2a, J2 and G are common in the
62 southern Mediterranean populations and form 20-30 % of their chrY lineages. In NEE,
63 the frequency of hg E2a'd is ~2-3%, hgs J2 and G respectively reach ~1-2% and ~1% of
64 the total pool of chrY lineages [3, 5, 6, 15, 16]. Hg Q has a frequency of 1-3% in most
65 European populations with the highest incidence in Sweden [3, 4]. Hg Q is otherwise
66 widely spread in Siberian populations and is among the major founding male lineages in
67 the peopling of the Americas [17, 18]. These rare hgs that make up less than 10% of
68 NEE male lineages, are mostly left unexplained and are often regarded as recent
69 scattered entries into populations. The small sample sizes and low phylogenetic
70 resolution has not allowed separation of rare lineages beyond the major hg labels. The
71 sequencing of complete Y-chromosomes provides a way to resolve the inner structure
72 of lineages on the phylogenetic tree regardless of their prevalence in populations [13,
73 14, 19, 20]. Sequencing a considerable number of well dispersed samples from NEE

74 reveals the distribution of rare lineages on the entire phylogenetic tree and provides
75 sufficiently granular data to estimate their split times. This builds the necessary
76 geographic and chronological context for surveying patterns of uncovered lineage
77 clusters stemming from a single node and hallmarking local expansions. The
78 coalescence ages of ancestral internal nodes and phylogenetically well-defined clusters
79 nested within disclose the geography and timeframe of local expansions as well as
80 possible gene flow involving ancestral carriers of rare male lineages in Estonia, Sweden
81 and their neighbouring populations.

82 Here we aim to analyze the previously understudied rare chrY lineages with a focus on
83 Estonia and Sweden together with their NEE neighbors and Germany to account for the
84 historic influence of the Baltic Germans. Additional populations are included to widen
85 the geographic context. We combined full sequences of Y-chromosomes from
86 populations inhabiting Estonia, Sweden, Finland, Latvia, Lithuania, Poland, Germany,
87 Ukraine and the Russian Federation to build updated phylogenetic trees for haplogroups
88 rare in NEE. In order to mitigate sampling bias that might influence any conclusions
89 drawn from such a rare substratum present among the populations, we tested the
90 representativeness of our two largest cohorts sampled from the Estonian and Swedish
91 populations by comparing their frequency compositions with sample sets independently
92 obtained from the same two populations.

93 **Materials and methods**

94 **Samples**

95 We screened the occurrence of rare hgs in a sample of 1,160 chrY sequences from male
96 donors (selected randomly by county of birth) from the population-based Estonian

97 Biobank [21]. The Estonian chrY sequences are part of the whole genome sequencing
98 (WGS) data set autosomally first described in Mitt et al. (2017) [22] for constructing a
99 population-specific imputation panel. Only chrY sequences of the haplogroups rare in
100 NEE (N=64) are included in the current study. Next, in scientific collaboration with the
101 commercial genetic testing company Gene by Gene (Houston, Texas, USA), we
102 screened the collection of customers who had provided informed consent for their data
103 to be used in scientific inquiry. This resulted in a total of 2018 male donors with self-
104 reported ancestry from Sweden, Finland, Latvia, Lithuania, Poland, Germany, Ukraine
105 and the Russian Federation. If the database contained more than 500 samples from a
106 respective country, individuals with identical self-reported paternal and maternal origin
107 were preferably selected. In case of smaller available sample sets, all samples with self-
108 reported paternal origin from the respective country were selected. From the resulting
109 set of 2,018 samples, we detected 222 Y-chromosomes belonging to the rare NEE
110 haplogroups and these samples were included in the current study. We collected
111 additional 139 chrY sequences from published sources resulting in the final set of 421
112 chrY sequences (Supplemental Table S1) used for reconstruction of phylogenetic trees
113 for rare hgs E2 (129 samples), J2 (136 samples), Q (83 samples) and G (71 samples)
114 (Figures 1-2 and Supplemental Figures S1-S7).

115 To test for possible sampling bias in the two largest sequencing cohorts, we screened
116 the haplogroup frequencies of two independent datasets – a total of 505 chrY sequences
117 available from the SweGen project (samples specifically selected to be representative of
118 the historic Swedish population [23]) and a randomly selected non-overlapping set of
119 genotyped 7,949 Estonian male donors from the Estonian Biobank.

120 **Data availability**

121 The Estonian WGS data are available on demand through the Estonian Biobank:
122 <https://www.geenivaramu.ee/en/biobank.ee/data-access>. In accordance to the consent
123 form signed by the customers of Gene by Gene commercial genetic testing company,
124 the sequencing data included in this study is used for the sole purpose of scientific
125 inquiry and is reported here on an aggregate level in the form of phylogenetic trees. For
126 both the Estonian Biobank and the Gene by Gene samples, summary-level data
127 including variable positions and their frequency in the cohort population have been
128 deposited to dbSNP with links to BioProject accession number PRJNA718714 in the
129 NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>). The Swedish
130 data from the SweGen Project is available upon request from the original authors of the
131 project [23].

132

133

134 **Sequencing, mapping and genotyping**

135 ChrY sequences from the Estonian Biobank and the SweGen project were generated
136 with Illumina Inc. (Illumina, San Diego, CA, USA) using HiSeq instruments (PCR-free
137 protocol) and targeted 30x genome-wide coverage. The personal genetic testing
138 company dataset was generated using the proprietary BigY Illumina-based targeted
139 chrY capture sequencing service ([https://learn.familytreedna.com/wp-](https://learn.familytreedna.com/wp-content/uploads/2014/08/BIG_Y_WhitePager.pdf)
140 [content/uploads/2014/08/BIG_Y_WhitePager.pdf](https://learn.familytreedna.com/wp-content/uploads/2014/08/BIG_Y_WhitePager.pdf)).

141 We used the same processing pipeline for all Illumina data. Fastq files were mapped
142 with BWA-MEM (v0.7.12) [24] on the human reference hs37d5
143 ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assem](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence)
144 [bly_sequence](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence)). Read duplicates were removed with Picard (v2.12.0)
145 (<http://broadinstitute.github.io/picard/>) and remaining unique reads realigned around
146 known indels, followed by base quality score recalibration (BQSR) using GATK (v3.8)
147 [25]. Variant calling was performed with GATK tool HaplotypeCaller in haploid mode.
148 All-sites VCF files were filtered with bcftools (v1.9) [26]. The Illumina data and
149 previously filtered data from Complete Genomics (Supplemental Table S1) were
150 merged with CombineVariants from GATK (v3.8) [25]. We extracted the effective
151 overlap between the two datasets by masking out all positions with 5% or higher
152 proportion of missing genotypes in either Illumina or Complete Genomics datasets. We
153 additionally excluded regions with poor mappability as described previously [13]
154 resulting in a total of 9.7 Mb of analysed sequence. Within this sequence, the resulting
155 numbers of variant positions used for phylogenetic reconstruction in each haplogroup
156 are given in Supplemental Tables S4-S7.

157 **Haplogroup assignment**

158 We assigned chrY haplogroups using yHaplo [27] for the Illumina capture and WGS
159 data. We used SNAPPY [28] for chrY haplogroup assignment of the genome-wide array
160 genotyping data.

161 **Comparisons with an independent Estonian cohort**

162 To validate the representativeness of sequenced Estonian chrY samples (N=1,160), we
163 compared the hg frequencies of this cohort against a ~7 times larger cohort of 7,949
164 Estonian male samples genotyped with the Illumina Infinium Global Screening Array
165 v2 (Illumina, San Diego, CA, USA) containing 6,638 Y-specific single nucleotide
166 variants (SNVs). To do this, we first assessed the accuracy of haplogroup assignments
167 obtained from this particular set of SNVs. We sub-sampled the 6,638 array-specific Y-
168 SNVs from the Estonian WGS data and used SNAPPY software to determine the
169 haplogroups from the extracted set of SNVs. We compared the results against those
170 from the software yHaplo [27]. The latter utilises the full set of SNPs in the WGS
171 samples. The results are identical on the highest level of the major branches and only
172 differ slightly at the finest resolution due to the lower number of array-genotyped SNPs
173 available to SNAPPY for detecting the haplogroups. However, this shows that hg
174 assignments based on the 6,638 array-specific Y-positions are accurate enough to be
175 compared to hg assignments based on full sequencing. The comparison of the hg
176 frequencies of the WGS-based and array-based Estonian datasets was performed using a
177 Wilcoxon signed rank test with continuity correction. We only used array-based data for
178 comparing haplogroup frequencies between two independent cohorts. For the
179 phylogeny reconstruction and phylogeographic analysis full sequencing data were used.

180 **Phylogeny reconstruction of rare paternal haplogroups**

181 We reconstructed phylogenies and estimated the coalescent times with the software
182 package BEAST v.1.7.5 [29]. We used a Bayesian skyline coalescent tree prior, the
183 general time reversible (GTR) substitution model with gamma-distributed rates and a
184 relaxed lognormal clock. The run was performed with the piecewise-constant coalescent
185 model. The mutation rate used was 0.74×10^{-9} (95% CI: $0.63 - 0.95 \times 10^{-9}$) per base

186 pair per year [13]. The results were visualized and checked for effective sample size
187 above 200 in Tracer v.1.4. Coalescence time estimates were computed with normally
188 distributed age priors with 10% standard deviation from previously published
189 phylogeny [13] and are in Supplemental Table S3. Lineages from hg R and I were used
190 as outgroups for hgs Q and G, respectively. Each run had thirty million chains logged
191 every 3000 steps and 10% discarded as burn-in. Two parallel with different random
192 number seeds were combined with LogCombiner.

193 The manually annotated phylogenetic trees, mutation lists and coalescence age
194 estimates are available in Supplemental Figures S1-S4 and Supplemental Tables S4-
195 S11. This study's updated nomenclature follows the criteria set in Karmin et al. (2015)
196 [13].

197 **Bayesian phylogeographic analysis**

198 To illustrate the potential direction of influx of the primarily Estonian subclades in hgs
199 E2a1-CTS1273 and J2b2-L283, we performed Bayesian phylogeographic analyses in
200 continuous space. For this we used available geographic coordinates for 59 sequences
201 belonging to hg E2a1-CTS1273 and for 41 sequences belonging to hg J2b2-L283. This
202 method has been originally developed and successfully used to reveal the ancestral
203 location and spatial dynamics of viruses in continuous space [30, 31]. We conducted the
204 analysis according to the publication exploring the history of Y-chromosomal hg J1 [32]
205 in BEAST v1.10.4 [33] using BEAGLE library v3.1.2 [34] for accelerated likelihood
206 evaluation. This statistically robust and absolutely data-driven method uses molecular
207 sequence data and geographic coordinates of the samples to infer phylogeography in a
208 continuous landscape while simultaneously reconstructing the evolutionary history in

209 time. It draws the confidence area of ancestral locations where the root and internal
210 nodes originated together with the directions and the speed of the diffusion (Figure 3).
211 The uncertainties of the maximum clade credibility tree node locations were visualised
212 with Spread3 v0.9.7.1rc software [35]. This inference approach accounts for the
213 coalescent, phylogenetic, molecular clock, location, and other uncertainties within a
214 single framework. Additional details are provided in Supplemental Note 1.

215 **Ethics approval**

216 All donors have provided informed consent and all experiments were performed in
217 accordance with the relevant guidelines and regulations of collaborating institutions.
218 Access to genetic data in Estonian Biobank was approved by the Research Ethics
219 Committee of the University of Tartu (permission number 1.1.-12/659 granted by the
220 Research Ethics Committee of the University of Tartu, Estonia). The chrY sequences
221 included from customers of the commercial personal genetic testing service were only
222 from individuals who had provided informed consent for the use of their data in
223 scientific research and for publication in aggregated form. The list of IDs along with
224 additional sample information is presented in Supplemental Table S1.

225

226 **Results**

227 *Phylogeny of rare lineages in NEE*

228 The studied 1160 high coverage sequences of Y-chromosomes from Estonia disclose 64
229 samples carrying male lineages rare in NEE (frequency of each under 3%), amounting
230 to ~6% of the total paternal lineage pool in Estonia. The most frequent minor lineage in

231 Estonia belongs to hg E2 (2.5%), followed by hgs J2 (1.9%) and hg G (0.9%), whereas
232 hg Q is the rarest (0.3%) (Supplemental Table S2). Our second largest sample set
233 consists of a total of 746 males from Sweden and discloses 78 samples with rare NEE
234 chrY lineages. The most common minor haplogroup in the Swedish cohort is hg Q
235 (4.6%); followed by hgs G (3%), E2 (1.7%) and J2 (1.2%) (Supplemental Table S2).

236 To verify the robustness of our frequency estimates, we compared hg frequencies of our
237 Swedish sample set and the SweGen cohort (N=505) [23]. The Wilcoxon signed rank
238 test showed no statistically significant differences between the two, either considering
239 all hgs (p-value=0.4689) or minor hgs with major hgs collapsed (p-value 0.6602).
240 Similarly, a comparison of hg frequencies between the Estonian sample set and an
241 independent non-overlapping set of 7,949 genotyped male samples from the Estonian
242 Biobank yielded no statistically significant differences in their hg composition, either
243 considering all haplogroups (Wilcoxon signed rank test p-value=0.4896) or rare hgs
244 with major hgs collapsed (Wilcoxon signed rank test p-value=0.9219).

245 Hg E originated in Africa with its sublineage E1 distributed solely on the African
246 continent, whereas the neighbour-lineage hg E2 displays a notably wider distribution.
247 Subclade E2-V13 is common (~10-20%) among south-eastern European populations [4,
248 6, 14, 16], falling to 10% in Anatolia and the Middle East [36] and declining towards
249 northern Europe to 1-2% in Scandinavia [4].

250 Here we reconstruct the phylogeny of hg E2a'b'c'd-M35. Its subclade E2a-M78 is
251 largely confined to Europe with a coalescence time of ~14 kya (95% CI: 10,432-18,566)
252 (Figure 1a, Supplemental Table S8). Within this subclade, L618 marker unites almost
253 all European samples that split ~13 (95% CI: 9,682-17,360) kya from the neighbouring

254 clade E2a2-V22. The latter consists primarily of samples from the Middle East with
255 deeper diversification times (Supplemental Figure S5). The absolute majority (25/29) of
256 hg E samples from Estonia belong to subclade E2a1-V13 (Supplemental Table S2). The
257 bulk of Estonian samples form clearly distinguishable clusters: lineage E2a1-S7461
258 contains an Estonian founding cluster that splits from the neighbour lineages with
259 Swedish and Middle Eastern origin ~ 4 kya (95% CI: 3,146 – 5,752, Supplemental Table
260 S8) and a radiation time of ~ 2 kya (95% CI: 1,398 – 2,999; Supplemental Table S8).
261 Similar pattern can be seen in the hg E2a1-B409 that has lineages from Germany and
262 Sweden and an exclusively Estonian cluster defined by marker Z37869 with a radiation
263 time of ~ 2 kya (95% CI: 1,150 – 2,428; Supplemental Table S8).

264 Hg J is one of the most common haplogroups in Western Asia and in regions
265 surrounding the Mediterranean Sea and thus was initially connected to the dispersion of
266 male farmers from the Fertile Crescent. Phylogenetic studies of hg J have shown
267 surviving ancient sublineages with radiation signs in the Bronze Age [37, 38].
268 Additionally, hg J2a and an unresolved hg J have been discovered in ancient DNA from
269 hunter-gatherer samples excavated in the Caucasus [39] and Karelia [40]. In southern
270 Europe, the most common hg J subclade is J2-M172, which, however, becomes rare
271 throughout the northern latitudes [4, 16].

272 Here we reconstruct the phylogenetic tree of hg J2-M172 (Figure 1b and Supplemental
273 Figure S6) with 134 individuals. A substantial part of NEE individuals belong to
274 sublineages within hg J2b2-L283 (Figure 1b) which splits from its neighboring clade at
275 ~ 16 kya (95% CI: 11,860 – 20,018; Supplemental Table S9). Hg J2b2-L283 itself split
276 ~ 7 kya (95% CI: 5,000 – 8,912) into two major sublineages J2b2-Z2505 and J2b2-
277 YP29. The latter is an exclusively Estonian cluster encompassing over half of all hg J

278 samples from Estonia (12 of 22) with an expansion time of ~2 kya (95% CI: 1,446 –
279 3,027) (Figure 1b, Supplemental Figure S6, Supplemental Table S9).

280 The other major hg J subbranch – J2a-M410 – contains samples from broad Eurasian
281 background which are distributed in subclades mostly coalescing during the early post-
282 Last Glacial Maximum – a much deeper time estimate than in the neighbouring hg J2b-
283 M12 phylogeny (Supplemental Figure S6). Lack of information on detailed geographic
284 or ethnic origin hinders any further conclusions regarding the single-origin clusters from
285 the Russian Federation (Supplemental Figure S6). Based on published research,
286 lineages of hg J2a-M67 are among the most common (~20%) paternal haplogroups of
287 the North Caucasus region [41], whereas in ethnic Russians this haplogroup amounts to
288 less than 2% [5, 6].

289 Hg Q is frequent in Siberian populations and is carried by over 85% of male Native
290 Americans [16–18, 42]. In Europe, the occurrence of hg Q is uneven and the general
291 frequency is low (~0.42%) [42], but hg Q is somewhat more frequent in the populations
292 of Sweden and Norway [3, 4]. It is the most numerous minor haplogroup in both of our
293 Swedish sample sets with frequencies of 2.6% and 4.6% (Supplemental Table S2). In
294 the datasets of Karlsson et al. (2006) [4] and Lappalainen et al. (2008) [3] the frequency
295 of hg Q fluctuates between 1% and 5% in different regions of Sweden. On the updated
296 phylogenetic tree, Swedish samples fall into two main clusters that separated from each
297 other around the peak of the Last Glacial Maximum ~20 kya (Figure 2). About a third
298 of the Swedish hg Q samples are defined by marker L804. Hg Q1a-L804 coalesces ~16
299 kya (95% CI: 12,456 – 19,874; Supplemental Table S10) with haplogroup Q1a-M3,
300 which today describes the overwhelming majority of Native American Y-chromosomes

301 [42]. The rapid diversification among Swedes in the L804-defined clade began ~3 kya
302 (95% CI: 1,961 – 3,917; Supplemental Table S10).

303 Haplogroup G-M201 is common in the Caucasus and the Middle East. Hg G is one of
304 the most prevalent male lineages in Sardinia and Corsica, but displays low frequencies
305 elsewhere in Europe [4, 14, 15]. Hg G splits into two basal lineages – hgs G1 and G2, of
306 which the former occurs infrequently in Western and Central Asia and is almost absent
307 in Europe [15]. Almost all hg G samples from NEE belong to hg G2-P287 that ~22 kya
308 (95% CI: 17,620 – 26,973) split into two main subclades – G2a-P15 and G2b-M377
309 (Supplemental Figure S7, Supplemental Table S11). The bulk of sampled European
310 individuals belong to subclade G2a2-P303 (Supplemental Table S2). Downstream, in hg
311 G2a2-Z727, the absolute majority of Swedish hg G samples forms localised clusters
312 with a variety of coalescence times (Supplemental Figure S7, Supplemental Table S11).

313

314 **Discussion**

315 In case of Estonia, our sequenced samples were collected across the country avoiding
316 large settlements with recorded extensive migration history. Considering a census size
317 of roughly 1 million, rare lineages amount to a total of 30,000 men evenly sampled
318 across the country and thus cannot be exclusively ascribed to any random influx of
319 recent migrants.

320 From the screened sample of 506 Finnish males we did not detect any rare NEE
321 lineages as almost all Finnish samples belong to hgs common among neighbouring

322 populations – a probable reflection of either differing migration history or of
323 demographic bottleneck(s) that have affected the Finnish population [43, 44].

324 Hg E sublineages have been associated with Neolithic demic diffusion into Europe [16],
325 but current ancient DNA data has shown this haplogroup to be uncommon among the
326 first agriculturalists in Europe [40]. In the resolved phylogenetic tree, the primarily
327 Middle Eastern neighbouring clade with deeply diverged lineages supports a possible
328 Levantine source of the European hg E2a1-V13. However, the split time predates the
329 Neolithic transition in Europe and matches better with the age of the Villabruna hunter-
330 gatherer cluster that displays earliest autosomal affinities to the Middle Eastern
331 populations detected in ancient samples from Europe [45]. The coalescence age of the
332 primarily European clades of hgs E3a1-V13 and J2b2-Z2505 underpins mid-Holocene
333 as the starting point of chrY variation growth in Western Europe (Figure 1) and
334 indicates a possible influx of male lineages from the Levant or the Caucasus.

335 The coalescence ages of Estonia-specific clades J2b2-YP29, E2a1-Z37869 and E2a1-
336 Y28220 broadly correspond to the Late Bronze Age and Iron Age period in Northern
337 Europe (Supplemental Figure S8). Additional sampling might certainly affect the
338 coalescence age of these clusters. However, the geographical spread across all Estonian
339 counties and current age estimates suggest that these expansions are not the result of
340 any migratory events from the recent recorded (last ~800 years) history of this region.
341 To infer the potential directions of influx of the clades J2b2-YP29, E2a1-Z37869, and
342 E2a1-Y28220, we conducted continuous Bayesian phylogeographic analysis of parent
343 hgs J2b2-L283 and E2a1-CTS1273. The estimated diffusion rate of hg J2b2-L283
344 equals 0.27 (95% HPDs: 0.1992 – 0.3478) and for hg E2a1-CTS1273 0.231 (95%

345 HPDs: 0.175 – 0.295) kilometers/year. The 80% HPD of the putative geographic centre
346 of diffusion for the hg J2b2 covers the area focused in present-day Poland, with a partial
347 covering of central and southeastern Europe, spreading further north and south (Figure
348 3a). The area for hg E2a1 ancestral location similarly covers central and eastern Europe
349 with a focus on Poland (Figure 3b), but the focal point appears to be more condensed.

350 From a conservative standpoint, all three subclades most probably arrived to present-
351 day Estonia from the direction of central Europe. However, based on currently available
352 data, it is not possible to say whether the evident local expansions initially began in
353 Estonia or were the carriers already sufficiently diversified on arrival.

354 Within hg Q, clusters defined by L804 and Y4838 capture almost all of Swedish hg Q
355 diversity, marking these lineages as an inherent, albeit scarce, part of the pool of male
356 lineages in Sweden. The scarcity of internal nodes on the branches leading to the two
357 now predominantly Swedish clusters hinders any discussion regarding a potential
358 direction of influx or ancestral center of diffusion. Due to the glacial coverage, the split
359 between lineages Q1a-L804 and the Native American Q1a-M3 could not have happened
360 in Scandinavia. Ancient DNA research confirms the presence of hg Q in the remains of
361 hunter-gatherers (~8 kya) from Latvia and Lower Volga Region in Russia [46]. Today,
362 European Q1 lineages are restricted to NEE with occasional findings in other
363 populations (single L804 derived English chrY sample in Grugni et al. (2019) [47]).
364 Precursors to current European hg Q1 sublineages could have been widely present in
365 North Eurasia during the Last Glacial Maximum and followed a primarily northern
366 (Siberian) route of dispersal into Europe. The presumptively common ancient gene pool
367 is reflected in the autosomal European affinity of 24,000-old Mal'ta sample from the
368 vicinity of Lake Baikal [48]. Alternatively, the prevalence of hg Q in Sweden could

369 testify of a more recent Siberian influence deduced both from modern and ancient DNA
370 analysis in northeastern Europe [2, 8]. Studies have demonstrated minor eastern
371 affinities in the autosomes and in the maternal lineages of the modern Saami, but small
372 sample sizes have not revealed Saami male lineages belonging to hg Q [2]. Further
373 sampling across Northern Eurasia might provide additional insights about these peculiar
374 North Eurasian hg Q lineages. A total of two out of the three Estonian hg Q samples
375 form a subset of the Swedish Y4838-defined cluster. It is most parsimonious to assume
376 that the paternal ancestors of the two Y4838-derived individuals arrived in Estonia
377 around 1-2 kya from Scandinavia.

378 Hg G2a has become firmly associated with the early Neolithic farmers of Europe [40,
379 46, 49]. Most of European hg G2a inner lineages started to diverge around 5-7 kya
380 (Supplemental Figure S7, Supplemental Table S11) – within the timeframe of the
381 European agricultural transition. In Sweden, it is the second most frequent minor chrY
382 haplogroup. The majority of Swedish carriers demonstrate a strong expansion signal
383 approximated to ~1 kya (nodes 52 and 60 in Supplemental Figure S7), whereas Estonian
384 samples are not part of the Swedish hg G2a diversity.

385 In conclusion, we demonstrate that in NEE, rare paternal lineages are not just single
386 lineages scattered across different subclades in the phylogeny. We identified several
387 population-specific clusters among less common haplogroups, which testify of radiation
388 events that have occurred in various timeframes and can be used to tentatively suggest
389 possible influx directions.

390 This study demonstrates the power of large-scale re-sequencing of Y-chromosomes to
391 explore and compare the male demographic history of single populations. Current
392 survey of rare lineages paves the way for future research involving large datasets of re-

393 sequenced genomes with a focus on those maternal and paternal lineages that have left a
394 major demographic impact on modern populations in NEE and elsewhere.

395

396 **Conflict of Interest**

397 D.M.B and C.B. declare stock ownership at Gene by Gene, Ltd. L.S. is an employee of
398 Gene by Gene.

399 **Funding**

400 This work was supported by institutional research funding IUT24-1 of the Estonian
401 Ministry of Education and Research, Estonian Research Council grants PRG243,
402 PRG1071 and project No. 2014-2020.4.01.16-0024 (MOBTT53) granted by the
403 European Regional Development Fund. A-M.I. is supported by Finnish Academy
404 (DIGIHUM project URKO, decision number 329257). High coverage genome data for
405 five 1000 Genomes samples were generated at the New York Genome Center with
406 funds provided by NHGRI Grant 3UM1HG008901-03S1.

407 **References**

- 408 1 Pankratov V, Montinaro F, Kushniarevich A, Hudjashov G, Jay F, Saag L *et al.*
409 Differences in local population history at the finest level: the case of the Estonian
410 population. *Eur J Hum Genet* 2020; **28**.
- 411 2 Tambets K, Yunusbayev B, Hudjashov G, Ilumäe AM, Rootsi S, Honkola T *et al.*
412 Genes reveal traces of common recent demographic history for most of the
413 Uralic-speaking populations. *Genome Biol* 2018; **19**: 1–20.

- 414 3 Lappalainen T, Laitinen V, Salmela E, Andersen P, Huoponen K, Savontaus ML
415 *et al.* Migration waves to the baltic sea region. *Ann Hum Genet* 2008; **72**: 337–
416 348.
- 417 4 Karlsson AO, Wallerström T, Götherström A, Holmlund G. Y-chromosome
418 diversity in Sweden - A long-time perspective. *Eur J Hum Genet* 2006; **14**: 963–
419 970.
- 420 5 Balanovsky O, Rootsi S, Pshenichnov A, Kivisild T, Churnosov M, Evseeva I *et*
421 *al.* Two Sources of the Russian Patrilineal Heritage in Their Eurasian Context.
422 *Am J Hum Genet* 2008; **82**: 236–250.
- 423 6 Kushniarevich A, Utevska O, Chuhryaeva M, Agdzhoyan A, Dibirova K,
424 Uktveryte I *et al.* Genetic heritage of the balto-slavic speaking populations: A
425 synthesis of autosomal, mitochondrial and Y-chromosomal data. *PLoS One* 2015;
426 **10**. doi:10.1371/journal.pone.0135820.
- 427 7 Jones ER, Zarina G, Moiseyev V, Lightfoot E, Nigst PR, Manica A *et al.* The
428 Neolithic Transition in the Baltic Was Not Driven by Admixture with Early
429 European Farmers. *Curr Biol* 2017; **27**: 576–582.
- 430 8 Lamnidis TC, Majander K, Jeong C, Salmela E, Wessman A, Moiseyev V *et al.*
431 Ancient Fennoscandian genomes reveal origin and spread of Siberian ancestry in
432 Europe. *Nat Commun* 2018; **9**. doi:10.1038/s41467-018-07483-5.
- 433 9 Saag L, Laneman M, Varul L, Malve M, Valk H, Razzak MA *et al.* The Arrival
434 of Siberian Ancestry Connecting the Eastern Baltic to Uralic Speakers further
435 East. *Curr Biol* 2019; **29**: 1701-1711.e16.

- 436 10 Myres NM, Rootsi S, Lin AA, Järve M, King RJ, Kutuev I *et al.* A major Y-
437 chromosome haplogroup R1b Holocene era founder effect in Central and
438 Western Europe. *Eur J Hum Genet* 2011; **19**: 95–101.
- 439 11 Underhill PA, Poznik GD, Rootsi S, Järve M, Lin AA, Wang J *et al.* The
440 phylogenetic and geographic structure of Y-chromosome haplogroup R1a. *Eur J*
441 *Hum Genet* 2015; **23**: 124–131.
- 442 12 Ilumäe AM, Reidla M, Chukhryaeva M, Järve M, Post H, Karmin M *et al.*
443 Human Y chromosome haplogroup N: a non-trivial time-resolved
444 phylogeography that cuts across language families. *Am J Hum Genet* 2016; **99**:
445 163–173.
- 446 13 Karmin M, Saag L, Vicente M, Wilson Sayres MA, Järve M, Gerst Talas U *et al.*
447 A recent bottleneck of Y chromosome diversity coincides with a global change in
448 culture. *Genome Res* 2015; **25**: 459–466.
- 449 14 Batini C, Hallast P, Zadik D, Delser PM, Benazzo A, Ghirotto S *et al.* Large-
450 scale recent expansion of European patrilineages shown by population
451 resequencing. *Nat Commun* 2015; **6**. doi:10.1038/ncomms8152.
- 452 15 Rootsi S, Myres NM, Lin AA, Järve M, King RJ, Kutuev I *et al.* Distinguishing
453 the co-ancestries of haplogroup G Y-chromosomes in the populations of Europe
454 and the Caucasus. *Eur J Hum Genet* 2012; **20**: 1275–1282.
- 455 16 Cruciani F, La Fratta R, Trombetta B, Santolamazza P, Sellitto D, Colomb EB *et*
456 *al.* Tracing past human male movements in northern/eastern Africa and western
457 Eurasia: New clues from Y-chromosomal haplogroups E-M78 and J-M12. *Mol*
458 *Biol Evol* 2007; **24**: 1300–1311.

- 459 17 Karafet TM, Osipova LP, Gubina MA, Posukh OL, Zegura SL, Hammer MF.
460 High levels of Y-chromosome differentiation among native Siberian populations
461 and the genetic signature of a boreal Hunter-Gatherer way of life. *Hum Biol*
462 2002; **74**: 761–789.
- 463 18 Dulik MC, Zhadanov SI, Osipova LP, Askapuli A, Gau L, Gokcumen O *et al.*
464 Mitochondrial DNA and Y chromosome variation provides evidence for a recent
465 common ancestry between Native Americans and indigenous Altaians. *Am J*
466 *Hum Genet* 2012; **90**: 229–246.
- 467 19 Hallast P, Batini C, Zadik D, Delser PM, Wetton JH, Arroyo-Pardo E *et al.* The
468 Y-chromosome tree bursts into leaf: 13,000 high-confidence SNPs covering the
469 majority of known clades. *Mol Biol Evol* 2014; **32**: 661–673.
- 470 20 Poznik GD, Xue Y, Mendez FL, Willems TF, Massaia A, Wilson Sayres MA *et*
471 *al.* Punctuated bursts in human male demography inferred from 1,244 worldwide
472 Y-chromosome sequences. *Nat Genet* 2016; **48**: 593–599.
- 473 21 Leitsalu L, Haller T, Esko T, Tammesoo ML, Alavere H, Snieder H *et al.* Cohort
474 profile: Estonian biobank of the Estonian genome center, university of Tartu. *Int*
475 *J Epidemiol* 2015; **44**: 1137–1147.
- 476 22 Mitt M, Kals M, Pärn K, Gabriel SB, Lander ES, Palotie A *et al.* Improved
477 imputation accuracy of rare and low-frequency variants using population-specific
478 high-coverage WGS-based imputation reference panel. *Eur J Hum Genet* 2017;
479 **25**: 869–876.
- 480 23 Ameer A, Dahlberg J, Olason P, Vezzi F, Karlsson R, Martin M *et al.* SweGen:
481 A whole-genome data resource of genetic variability in a cross-section of the

- 482 Swedish population. *Eur J Hum Genet* 2017; **25**: 1253–1260.
- 483 24 Li H. Aligning sequence reads, clone sequences and assembly contigs with
484 BWA-MEM. 2013. <https://arxiv.org/abs/1303.3997>.
- 485 25 Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Auwera GA
486 Van der *et al.* Scaling accurate genetic variant discovery to tens of thousands of
487 samples. *bioRxiv* 2017. doi:10.1101/201178.
- 488 26 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N *et al.* The sequence
489 alignment/map format and SAMtools. *Bioinformatics* 2009; **25**: 2078–2079.
- 490 27 Poznik GD. Identifying Y-chromosome haplogroups in arbitrarily large samples
491 of sequenced or genotyped men. *bioRxiv* 2016. doi:10.1101/088716.
- 492 28 Severson AL, Shortt JA, Mendez FL, Wojcik GL, Bustamante CD, Gignoux CR.
493 SNAPPY: Single Nucleotide Assignment of Phylogenetic Parameters on the Y
494 chromosome. *bioRxiv* 2018. doi:10.1101/454736.
- 495 29 Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with
496 BEAUti and the BEAST 1.7. *Mol Biol Evol* 2012; **29**: 1969–1973.
- 497 30 Lemey P, Rambaut A, Welch JJ, Suchard MA. Phylogeography Takes a Relaxed
498 Random Walk in Continuous Space and Time. *Mol Biol Evol* 2010; **27**: 1877–
499 1885.
- 500 31 Pybus OG, Suchard MA, Lemey P, Bernardin FJ, Rambaut A, Crawford FW *et*
501 *al.* Unifying the spatial epidemiology and molecular evolution of emerging
502 epidemics. *Proc Natl Acad Sci* 2012; **109**: 15066–15071.
- 503 32 Sahakyan H, Margaryan A, Saag L, Karmin M, Bahmanimehr A, Parik J *et al.*

- 504 Origin and diffusion of human Y chromosome haplogroup J1-M267. *Sci Rep*
505 2021. doi:10.1038/s41598-021-85883-2.
- 506 33 Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A.
507 Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10.
508 *Virus Evol* 2018; **4**: 1–5.
- 509 34 Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO *et al.*
510 BEAGLE: An Application Programming Interface and High-Performance
511 Computing Library for Statistical Phylogenetics. *Syst Biol* 2012; **61**: 170–173.
- 512 35 Bielejec F, Baele G, Vrancken B, Suchard MA, Rambaut A, Lemey P. Spread3:
513 interactive visualization of spatiotemporal history and trait evolutionary
514 processes. *Mol Biol Evol* 2016; **33**: 2167–2169.
- 515 36 Trombetta B, D’Atanasio E, Massaia A, Myres NM, Scozzari R, Cruciani F *et al.*
516 Regional differences in the accumulation of SNPs on the male-specific portion of
517 the human y chromosome replicate autosomal patterns: Implications for genetic
518 dating. *PLoS One* 2015; **10**: 1–18.
- 519 37 Finocchio A, Trombetta B, Messina F, D’Atanasio E, Akar N, Loutradis A *et al.*
520 A finely resolved phylogeny of y chromosome Hg J illuminates the processes of
521 Phoenician and Greek colonizations in the Mediterranean. *Sci Rep* 2018; **8**: 3–11.
- 522 38 Zalloua PA, Platt DE, El Sibai M, Khalife J, Makhoul N, Haber M *et al.*
523 Identifying Genetic Traces of Historical Expansions: Phoenician Footprints in the
524 Mediterranean. *Am J Hum Genet* 2008; **83**: 633–642.
- 525 39 Jones ER, Gonzalez-Fortes G, Connell S, Siska V, Eriksson A, Martiniano R *et*

526 *al.* Upper Palaeolithic genomes reveal deep roots of modern Eurasians. *Nat*
527 *Commun* 2015; **6**. doi:10.1038/ncomms9912.

528 40 Mathieson I, Lazaridis I, Rohland N, Mallick S, Patterson N, Roodenberg SA *et*
529 *al.* Genome-wide patterns of selection in 230 ancient Eurasians. *Nature* 2015;
530 **528**: 499–503.

531 41 Yunusbayev B, Metspalu M, Ja M, Kutuev I, Rootsi S, Metspalu E *et al.* The
532 Caucasus as an Asymmetric Semipermeable Barrier to Ancient Human
533 Migrations Research article. *Mol Biol Evol* 2012; **29**: 359–365.

534 42 Zegura SL, Karafet TM, Zhivotovsky LA, Hammer MF. High-Resolution SNPs
535 and Microsatellite Haplotypes Point to a Single, Recent Entry of Native
536 American Y Chromosomes into the Americas. *Mol Biol Evol* 2004; **21**: 164–175.

537 43 Kittles RA, Bergen AW, Urbanek M, Virkkunen M, Linnoila M, Goldman D *et*
538 *al.* Autosomal, mitochondrial, and Y chromosome DNA variation in Finland:
539 Evidence for a male-specific bottleneck. *Am J Phys Anthropol* 1999; **108**: 381–
540 399.

541 44 Martin AR, Karczewski KJ, Kerminen S, Kurki MI, Sarin AP, Artomov M *et al.*
542 Haplotype Sharing Provides Insights into Fine-Scale Population History and
543 Disease in Finland. *Am J Hum Genet* 2018; **102**: 760–775.

544 45 Fu Q, Posth C, Hajdinjak M, Petr M, Mallick S, Fernandes D *et al.* The genetic
545 history of Ice Age Europe. *Nature* 2016; **534**: 200–205.

546 46 Mathieson I, Alpaslan-Roodenberg S, Posth C, Szécsényi-Nagy A, Rohland N,
547 Mallick S *et al.* The genomic history of southeastern Europe. *Nature* 2018; **555**:

548 197–203.

549 47 Grugni V, Raveane A, Ongaro L, Battaglia V, Trombetta B, Colombo G *et al.*

550 Analysis of the human Y-chromosome haplogroup Q characterizes ancient

551 population movements in Eurasia and the Americas. *BMC Biol* 2019; **17**: 1–14.

552 48 Raghavan M, Skoglund P, Graf KE, Metspalu M, Albrechtsen A, Moltke I *et al.*

553 Upper palaeolithic Siberian genome reveals dual ancestry of native Americans.

554 *Nature* 2014; **505**. doi:10.1038/nature12736.

555 49 Marchi N, Winkelbach L, Schulz I, Brami M, Hofmanová Z. The mixed genetic

556 origin of the first farmers of Europe. *bioRxiv* 2020.

557 doi:10.1101/2020.11.23.394502.

558

559

560

561

562

563

564

565

566

567

568 **Titles and legends to figures**

569 **Figure 1.** Schematic phylogenetic trees of hg E2a and J2b

570 The calibrated trees were constructed using BEAST v.1.7.5 software package. Internal
571 nodes, sub-clade names and population names (numbers show the number of samples)
572 are indicated. Internal nodes with posterior probabilities <0.73 are not shown. Samples
573 from Estonia and Sweden are marked in blue and orange, respectively.

574 A) A schematic phylogenetic tree of hg E2a is based on 132 high coverage chrY
575 sequences. Neighbor-clade E2b and its sublineages are marked in grey. Detailed tree
576 can be found in Supplemental materials (Supplemental Figure S5). Age estimates can be
577 found in Supplemental Table S8. All the subclade (node) defining mutations and marker
578 names are presented in Supplemental Table S4.

579 B) A schematic phylogenetic tree of hg J2b is based on 136 high coverage chrY
580 sequences. Neighbor-clade J2a and its sublineages are marked in grey. Detailed tree can
581 be found in Supplemental materials (Supplemental Figure S6). Age estimates can be
582 found in Supplemental Table S9. All the sub-clade (node) defining mutations and
583 marker names are presented in Supplemental Table S5.

584 **Figure 2.** Detailed phylogenetic tree of hg Q

585 A detailed phylogenetic tree of hg Q-M242 is based on 84 high coverage chrY
586 sequences. Two hg R1a sequences were used for an outgroup. The detailed calibrated
587 tree was constructed using BEAST v.1.7.5 software package. Internal nodes, sub-clade
588 names and population names are indicated. Internal nodes with posterior probabilities
589 <0.73 are not shown. Age estimates can be found in Supplemental Table S10. All the
590 subclade (node) defining mutations and marker names are presented in Supplemental

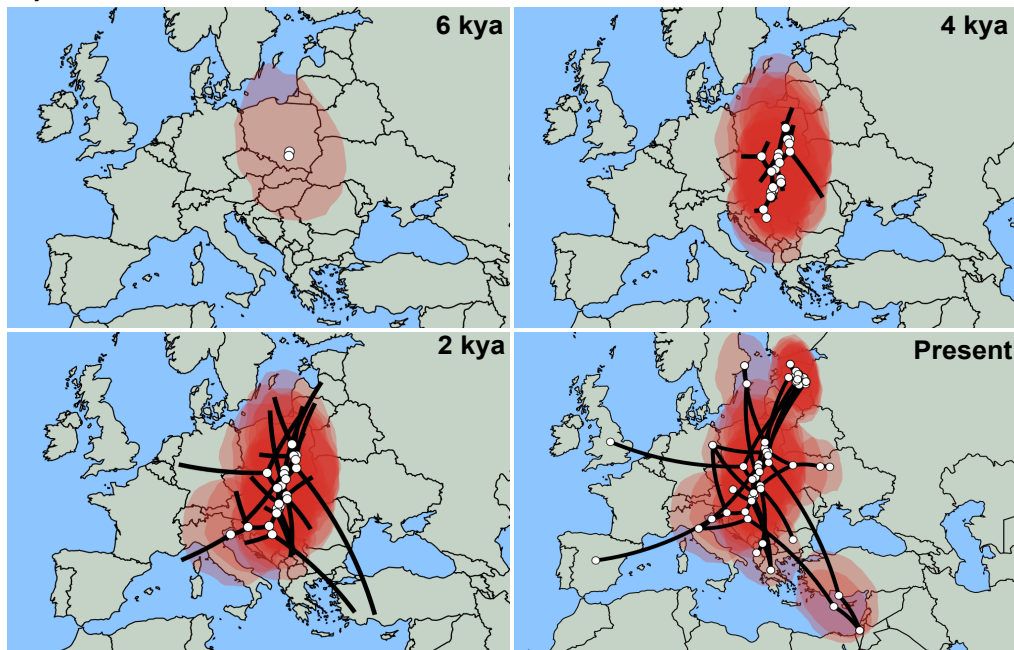
591 Table S6. Samples from Estonia and Sweden are marked in blue and orange,
592 respectively.

593

594 **Figure 3.** Phylogeographic spread maps of hgs J2b2-L283 and E2a1-CTS1273 in
595 Europe

596 Maps indicate the phylogeographic spread of A) J2b2-L283 around 6 kya, 4 kya, 2 kya
597 and in the present, and B) E2a1-CTS1273 around 5-6 kya, 4 kya, 2 kya and in the
598 present. Shaded in pink are the 80% HPD areas of the node locations inferred by
599 Bayesian continuous phylogeographic analysis with Beast v1.10.4 software. White
600 circles indicate the median locations of the nodes, while black lines indicate the
601 branches of the maximum clade credibility tree.

A) J2b2-L283



B) E2a1-CTS1273

