

1 The genetic history of admixture across inner Eurasia

2
3 Choongwon Jeong^{1,2,†,*}, Oleg Balanovsky^{3,4,†}, Elena Lukianova³, Nurzhibek Kahbatkyzy^{5,6}, Pavel
4 Flegontov^{7,8}, Valery Zaporozhchenko^{3,4}, Alexander Immel¹, Chuan-Chao Wang^{1,9}, Olzhas Ixan⁵, Elmira
5 Khussainova⁵, Bakhytzhan Bekmanov^{5,6}, Victor Zaibert¹⁰, Maria Lavryashina¹¹, Elvira Pocheshkhova¹²,
6 Yuldash Yusupov¹³, Anastasiya Agdzhoyan^{3,4}, Sergey Koshelev¹⁴, Andrei Bukin¹⁵, Pagbajabyn Nymadawa¹⁶,
7 Shahlo Turdikulova¹⁷, Dilbar Dalimova¹⁷, Mikhail Churnosov¹⁸, Roza Skhalyakho⁴, Denis Daragan⁴, Yuri
8 Bogunov^{3,4}, Anna Bogunova⁴, Alexandr Shtrunov⁴, Nadezhda Dubova¹⁹, Maxat Zhabagin^{20,21}, Levon
9 Yepiskoposyan²², Vladimir Churakov²³, Nikolay Pislegin²³, Larissa Damba²⁴, Ludmila Saroyants²⁵,
10 Khadizhat Dibirova^{3,4}, Lubov Atramentova²⁶, Olga Utevska²⁶, Eldar Idrisov²⁷, Evgeniya Kamenshchikova⁴,
11 Irina Evseeva²⁸, Mait Metspalu²⁹, Alan K. Outram³⁰, Martine Robbeets², Leyla Djansugurova^{5,6}, Elena
12 Balanovska⁴, Stephan Schiffels¹, Wolfgang Haak¹, David Reich^{31,32}, Johannes Krause^{1,*}

13
14 ¹ Department of Archaeogenetics, Max Planck Institute for the Science of Human History, Jena, Germany
15 ² Eurasia3angle Research Group, Max Planck Institute for the Science of Human History, Jena, Germany
16 ³ Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russia
17 ⁴ Federal State Budgetary Institution «Research Centre for Medical Genetics», Moscow, Russia
18 ⁵ Department of Population Genetics, Institute of General Genetics and Cytology, SC MES RK, Almaty, Kazakhstan
19 ⁶ Department of Molecular Biology and Genetics, Kazakh National University by al-Farabi, Almaty, Kazakhstan
20 ⁷ Department of Biology and Ecology, Faculty of Science, University of Ostrava, Ostrava, Czech Republic
21 ⁸ Faculty of Science, University of South Bohemia and Biology Centre, Czech Academy of Sciences, České Budějovice,
22 Czech Republic
23 ⁹ Department of Anthropology and Ethnology, Xiamen University, Xiamen 361005, China
24 ¹⁰ Institute of Archeology and Steppe Civilization, Kazakh National University by al-Farabi, Almaty, Kazakhstan
25 ¹¹ Kemerovo State Medical University, Krasnaya 3, Kemerovo, Russia
26 ¹² Kuban State Medical University, Krasnodar, Russia
27 ¹³ Institute of Strategic Research of the Republic of Bashkortostan, Ufa, Russia
28 ¹⁴ Faculty of Geography, Lomonosov Moscow State University, Moscow, Russia
29 ¹⁵ Transbaikalian State University, Chita, Russia
30 ¹⁶ Mongolian Academy of Medical Sciences, Ulaanbaatar, Mongolia
31 ¹⁷ Center for Advanced Technologies under the Ministry of Innovational Development, Tashkent, Uzbekistan
32 ¹⁸ Belgorod State University, Belgorod, Russia
33 ¹⁹ The Institute of Ethnology and Anthropology of the Russian Academy of Sciences, Moscow, Russia
34 ²⁰ National Laboratory Astana, Nazarbayev University, Astana, Kazakhstan
35 ²¹ National Center for Biotechnology, Astana, Kazakhstan
36 ²² Laboratory of Ethnogenomics, Institute of Molecular Biology of National Academy of Sciences, Yerevan, Armenia
37 ²³ Udmurt Institute of History, Language and Literature of Udmurt Federal Research Center of the Ural Branch of the Russian
38 Academy of Sciences, Russia
39 ²⁴ Research Institute of Medical and Social Problems and Control of the Healthcare Department of Tuva Republic, Kyzyl, Russia
40 ²⁵ Leprosy Research Institute, Astrakhan, Russia
41 ²⁶ V. N. Karazin Kharkiv National University, Kharkiv, Ukraine
42 ²⁷ Astrakhan branch of the Russian Academy of National Economy and Public Administration under the President of the Russian
43 Federation, Astrakhan, Russia
44 ²⁸ Northern State Medical University, Arkhangelsk, Russia
45 ²⁹ Estonian Biocentre, Institute of Genomics, University of Tartu, Tartu 51010, Estonia
46 ³⁰ Department of Archaeology, University of Exeter, Exeter EX4 4QE, UK
47 ³¹ Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA
48 ³² Howard Hughes Medical Institute, Harvard Medical School, Boston, Massachusetts 02115, USA
49 [†] These authors contributed equally to this work
50 ^{*} Correspondence: jeong@shh.mpg.de (C.J.), krause@shh.mpg.de (J.K.)

51 **Abstract**

52 The indigenous populations of inner Eurasia, a huge geographic region covering the central
53 Eurasian steppe and the northern Eurasian taiga and tundra, harbor tremendous diversity in their genes,
54 cultures and languages. In this study, we report novel genome-wide data for 763 individuals from Armenia,
55 Georgia, Kazakhstan, Moldova, Mongolia, Russia, Tajikistan, Ukraine, and Uzbekistan. We furthermore
56 report additional damage-reduced genome-wide data of two previously published individuals from the
57 Eneolithic Botai culture in Kazakhstan (~5,400 BP). We find that present-day inner Eurasian populations
58 are structured into three distinct admixture clines stretching between various western and eastern Eurasian
59 ancestries, mirroring geography. The Botai and more recent ancient genomes from Siberia show a
60 decrease in contribution from so-called “ancient North Eurasian” ancestry over time, detectable only in the
61 northern-most “forest-tundra” cline. The intermediate “steppe-forest” cline descends from the Late Bronze
62 Age steppe ancestries, while the “southern steppe” cline further to the South shows a strong West/South
63 Asian influence. Ancient genomes suggest a northward spread of the southern steppe cline in Central Asia
64 during the first millennium BC. Finally, the genetic structure of Caucasus populations highlights a role of
65 the Caucasus Mountains as a barrier to gene flow and suggests a post-Neolithic gene flow into North
66 Caucasus populations from the steppe.

67

68

69 Present-day human population structure is often marked by a correlation between geographic and genetic
70 distances^{1,2}, reflecting continuous gene flow among neighboring groups, a process known as “isolation by
71 distance”. However, there are also striking failures of this model, whereby geographically proximate
72 populations can be quite distantly related. Such barriers to gene flow often correspond to major geographic
73 features, such as the Himalayas³ or the Caucasus Mountains⁴. Many cases also suggest the presence of
74 social barriers to gene flow. For example, early Neolithic farming populations in Central Europe show a
75 remarkable genetic homogeneity suggesting minimal genetic exchange with local hunter-gatherer
76 populations through the initial expansion; mixing of these two gene pools became evident only after
77 thousands of years in the middle Neolithic⁵. Present-day Lebanese populations provide another example
78 by showing a population stratification reflecting their religious community⁶. There are also examples of
79 geographically very distant populations that are closely related: for example, people buried in association
80 with artifacts of the Yamnaya horizon in the Pontic-Caspian steppe and the contemporaneous Afanasievo
81 culture 3,000 km east in the Altai-Sayan Mountains^{7,8}.

82 The vast region of the Eurasian inland (“inner Eurasia” herein) is split into distinct ecoregions,
83 such as the Eurasian steppe in central Eurasia, boreal forests (taiga) in northern Eurasia, and the Arctic
84 tundra at the periphery of the Arctic Ocean (Fig. 1). These ecoregions stretch in an east-west direction
85 within relatively narrow north-south bands. Various cultural features show a distribution that broadly
86 mirrors the eco-geographic distinction in inner Eurasia. For example, indigenous peoples of the Eurasian
87 steppe traditionally practice nomadic pastoralism^{9,10}, while northern Eurasian peoples in the taiga mainly
88 rely on reindeer herding and hunting¹¹. The subsistence strategies in each of these ecoregions are often
89 considered to be adaptations to the local environments¹².

90 At present there is limited information about how environmental and cultural influences are
91 mirrored in the genetic structure of inner Eurasians. Recent genome-wide studies of inner Eurasians
92 mostly focused on detecting and dating genetic admixture in individual populations¹³⁻¹⁶. So far only three
93 studies have reported recent genetic sharing between geographically distant populations based on the
94 analysis of “identity-by-descent” segments^{13,17,18}. One study reports a long-distance extra genetic sharing

95 between Turkic populations based on a detailed comparison between Turkic-speaking groups and their
96 non-Turkic neighbors¹³. The other two studies extend this approach to some Uralic and Yeniseian-
97 speaking populations^{17,18}. However, a comprehensive spatial genetic analysis of inner Eurasian
98 populations is still lacking.

99 Ancient DNA studies have already shown that human populations of this region have dramatically
100 transformed over time. For example, the Upper Paleolithic genomes from the Mal'ta and Afontova Gora
101 sites in southern Siberia revealed a genetic profile, often called "Ancient North Eurasians (ANE)", which
102 is deeply related to Paleolithic/Mesolithic hunter-gatherers in Europe and also substantially contributed to
103 the gene pools of present-day Native Americans, Siberians, Europeans and South Asians^{19,20}. Studies of
104 Bronze Age steppe populations found the appearance of additional Western Eurasian-related ancestries
105 across the steppe from the Pontic-Caspian to the Altai-Sayan regions, here we collectively refer to as
106 "Western Steppe Herders (WSH)": the earlier populations associated with the Yamnaya and Afanasievo
107 cultures (often called "steppe Early and Middle Bronze Age"; "steppe_EMBA") and the later ones
108 associated with many cultures such as Potapovka, Sintashta, Srubnaya and Andronovo to name a few
109 (often called "steppe Middle and Late Bronze Age"; "steppe_MLBA")⁸. The steppe_MLBA gene pool
110 was largely descended from the preceding steppe_EMBA gene pool, with a substantial contribution from
111 Late Neolithic Europeans.²¹ Also, recent archaeogenetic studies trace multiple large-scale trans-Eurasian
112 migrations over the last several millennia using ancient inner Eurasian genomes^{22,23}, including individuals
113 from the Eneolithic Botai culture in northern Kazakhstan in the 4th millennium BC²⁴. These studies now
114 provide a rich context to interpret present-day population structure of inner Eurasians and to characterize
115 ancient admixtures in fine resolution.

116 In this study, we analyzed newly produced genome-wide data for 763 individuals belonging to 60
117 self-reported ethnic groups to provide a dense portrait of the genetic structure of inner Eurasians. We also
118 produced damage-reduced genome-wide data of two ancient Botai individuals, whose genome-wide data
119 were recently published²³, to explore the genetic structure of pre-Bronze Age populations in inner Eurasia
120 (Table 1). We aimed at characterizing the genetic composition of inner Eurasians in fine resolution by

121 applying both allele frequency- and haplotype-based methods. Based on the fine-scale genetic profile, we
122 further explored if and where the barriers and conduits of gene flow exist in inner Eurasia.

123

124

125 **Results**

126

127 **Present-day Inner Eurasians form distinct east-west genetic clines mirroring geography.** We

128 generated genome-wide genotype data of 763 participants who represent a majority of large ethnic groups
129 in Armenia, Georgia, Kazakhstan, Moldova, Mongolia, Russia, Tajikistan, Ukraine, and Uzbekistan (Fig.

130 1 and Table S1). We merged new data with published data of present-day^{20,25,26} and ancient

131 individuals^{3,8,19-23,27-42} (Table S2). The final data set covers 581,230 autosomal single nucleotide

132 polymorphisms (SNPs) in the Affymetrix Axiom® Genome-wide Human Origins 1 (“HumanOrigins”)

133 array platform⁴³.

134 In a Principal Component Analysis (PCA) of Eurasian individuals, we find that PC1 separates

135 eastern and western Eurasian populations, PC2 splits eastern Eurasians along a north-south cline, and PC3

136 captures variation in western Eurasians with Caucasus and northeastern European populations at opposite

137 ends (Fig. 2a and Supplementary Figs. 1-2). Inner Eurasians are scattered across PC1 in between,

138 mirroring their geographic locations. Strikingly, they seem to be structured into three distinct west-east

139 genetic clines running between different western and eastern Eurasian groups, instead of being evenly

140 spaced in PC space. The uppermost cline, composed of individuals from northern Eurasia, mostly

141 speaking Uralic or Yeniseian languages, connects northeast Europeans and the Uralic (Samoyedic)

142 speaking Nganasans from northern Siberia. The other two lower clines are occupied by individuals from

143 the Eurasian steppe, mostly speaking Turkic and Mongolic languages. Both clines run into

144 Turkic/Mongolic-speaking populations in southern Siberia and Mongolia, and further into Tungusic-

145 speaking populations in Manchuria and the Russian Far East in the East; however, they diverge in the west,

146 one heading to the Caucasus and the other heading to populations of the Volga-Ural area (Fig. 2 and

147 [Supplementary Fig. 2](#)). Four groups, Daur, Mongola, Tu and Dungans, are located alongside other East
148 Asian populations and displaced from the three inner Eurasian clines.

149 A model-based clustering analysis using ADMIXTURE shows a similar pattern ([Fig. 2b](#) and
150 [Supplementary Fig. 3](#)). Overall, the proportions of ancestry components associated with eastern or western
151 Eurasians are well correlated with longitude in inner Eurasians ([Fig. 3](#)). Notable outliers include known
152 historical migrants such as Kalmyks, Nogais and Dungans. The Uralic- and Yeniseian-speaking
153 populations, as well as Russians from multiple locations, derive most of their eastern Eurasian ancestry
154 from a component most enriched in Nganasans, while Turkic/Mongolic-speakers have this component
155 together with another component most enriched in populations from the Russian Far East, such as Ulchi
156 and Nivkh ([Supplementary Fig. 3](#)). Turkic/Mongolic-speakers comprising the bottom-most cline **have** a
157 distinct western Eurasian ancestry profile: they have a high proportion of a component most enriched in
158 Mesolithic Caucasus hunter-gatherers (“CHG”)³⁰ and Neolithic Iranians (“Iran_N”)²⁰ and frequently
159 harbor another component enriched in present-day South Asians ([Supplementary Fig. 4](#)). Based on the
160 PCA and ADMIXTURE results, we heuristically assign inner Eurasians into three clines: the “forest-
161 tundra” cline includes Russians and all Uralic- and Yeniseian-speakers, the “steppe-forest” cline includes
162 Turkic- and Mongolic-speaking populations from the Volga and the Altai-Sayan regions and southern
163 Siberia, and the “southern steppe” cline includes the rest of populations. We separate four groups (Daur,
164 Mongola, Tu and Dungans) as “others” ([Supplementary Table 2](#)).

165 The genetic barriers splitting the inner Eurasians are also found in the EEMS (“estimated effective
166 migration surface”) analysis⁴⁴ ([Supplementary Fig. 5](#)). Inferred barriers to gene flow are often co-localized
167 with geographic features or genetic gaps. We observe a barrier overlapping with the Urals, one separating
168 Beringian populations from the rest, one separating southern Siberians from central and northern Siberians,
169 and one separating Caucasus populations from those further to the north. The southern Siberian barrier
170 matches with our distinction between the steppe-forest and forest-tundra populations, with the exception
171 of two northern-most Turkic speaking populations, Yakuts and Dolgans. The Caucasus barrier also
172 matches with our distinction between the southern steppe and steppe-forest populations. A local EEMS

173 analysis on the Caucasus shows fine-scale barriers and conduits of gene flow, matching with the fine-scale
174 structure within Caucasus populations ([Supplementary Note 1](#)).

175
176 **High-resolution tests of admixture distinguish the genetic profile of source populations in the inner**
177 **Eurasian clines.** We performed both allele frequency-based three-population (f_3) tests and a haplotype-
178 sharing-based GLOBETROTTER analysis to characterize the admixed gene pools of inner Eurasian
179 groups. For these group-based analyses, we manually removed 87 outliers based on PCA results
180 ([Supplementary Table 1](#)). We also split a few inner Eurasian groups showing genetic heterogeneity into
181 subgroups based on PCA results and their sampling locations ([Supplementary Table 1](#)). This was done to
182 minimize false positive admixture signals. Including two Aleut populations as positive control targets, we
183 chose a total of 73 groups as the targets of admixture tests and another 260 groups (167 present-day and 93
184 ancient groups) as the “sources” to represent world-wide genetic diversity ([Supplementary Table 2](#)).

185 Testing all possible pairs of 167 present-day “source” groups as references, we detect highly
186 significant f_3 statistics for 66 of 73 targets (< -3 SE; standard error; [Supplementary Table 3](#)). Negative f_3
187 values mean that allele frequencies of the target group are on average intermediate between those of the
188 references, providing unambiguous evidence that the target population is a mixture of groups related,
189 perhaps deeply, to the source populations.⁴³ Extending the references to include 93 ancient groups, the
190 **remaining** seven groups also have small f_3 statistics around zero (-5.1 SE to $+2.7$ SE). Reference pairs with
191 the most negative f_3 statistics for the most part involve one eastern and one western Eurasian groups
192 supporting the qualitative impression of east-west admixture from PCA and ADMIXTURE analysis. To
193 highlight the difference between the distinct inner Eurasian clines, we looked into f_3 results with
194 representative reference pairs comprising two ancient western (Srubnaya to represent MLBA_steppe
195 ancestry²¹ and Chalcolithic Iranians (“Iran_ChL”) to represent West/South Asian-related ancestry²⁰;
196 [Supplementary Table 1](#)) and three eastern Eurasian groups (Mixe, Nganasan and Ulchi). In the southern
197 steppe cline populations, reference pairs with Chalcolithic Iranians tend to produce more negative f_3
198 statistics than those with Srubnaya while the opposite pattern is uniformly observed for the steppe-forest

199 and forest-tundra populations (Fig. 4a). Reference pairs with Nganasans mostly result in more negative f_3
200 statistic than those with Ulchi in the forest-tundra populations, but the opposite pattern is dominant in the
201 southern steppe populations. The steppe-forest cline populations show an intermediate pattern: seven
202 northern groups (Chuvash, Bashkir_north, Tatar_Zabolotniye, Todzin, Tofalar, Dolgan and Yakut) have
203 more negative f_3 with Nganasans while the others have more negative f_3 with Ulchi. Most of these seven
204 groups are also upward-shifted in PCA toward the forest-tundra cline, suggesting a cross-talk between two
205 clines.

206 To perform a higher resolution characterization of the admixture landscape, we performed a
207 haplotype-based GLOBETROTTER analysis. We took a “regional” approach, meaning that all 73 target
208 groups were modeled as a patchwork of haplotypes from the 167 reference groups but not those from any
209 target. The goal of this approach was to minimize false negative results due to sharing of admixture
210 history between targets. All 73 targets show a robust signal of admixture: i.e. a correlation of ancestry
211 status shows a distinct pattern of decay over genetic distance in all bootstrap replicates (bootstrap $p < 0.01$
212 for all 73 targets; Supplementary Table 4). When the relative contribution of references, categorized to 12
213 groups (Supplementary Table 2), into the two main sources of the admixture signal (“date 1 PC 1”) is
214 considered, we observe a pattern comparable to PCA, ADMIXTURE and f_3 results (Fig. 4b). The
215 European references provide a major contribution for the western Eurasian-related source in the forest-
216 tundra and steppe-forest populations while the Caucasus/Iranian references do so in the southern steppe
217 populations. Similarly, Siberian references make the highest contribution to the eastern Eurasian-related
218 source in the forest-tundra populations, followed by the steppe-forest and southern steppe ones. Admixture
219 date estimates from GLOBETROTTER range 7-55 generations (200-1600 BP; years before present; using
220 29 years per generation⁴⁵; Supplementary Fig. 6 and Supplementary Note 2). These match with previous
221 reports using similar methodologies¹³, but much younger observed admixtures in the Late Bronze and Iron
222 Ages^{8,39}.

223

224 **Admixture modeling of inner Eurasians shows multiple different temporal layers for present-day**
225 **admixture clines.** Using F -statistic-based approaches, we show that the Eneolithic Botai gene pool was
226 closely related to the ANE ancestry and substantially contributed to the later Okunevo individuals
227 (Supplementary Note 3). To test if this ancient layer left a genetic legacy in later populations of inner
228 Eurasia, we systematically explored diverse qpAdm-based admixture models to inner Eurasian
229 populations.

230 Two-way mixture of Ulchi/Nganasan and Srubnaya approximates the steppe-forest populations
231 surprisingly well ($\chi^2 p \geq 0.05$ and ≥ 0.01 for 12/24 and 18/24 populations, respectively; Supplementary
232 Table 5). A more complex three-way model of Ulchi+Srubnaya+AG3 fits all steppe-forest populations (χ^2
233 $p \geq 0.05$ for 24/24 populations; Fig. 5 and Supplementary Table 5). Similarly, Nganasan+Srubnaya+AG3
234 provides a good fit to most populations, but with negative contribution from AG3 ($\chi^2 p \geq 0.05$ for 19/24
235 populations). We interpret this as reflecting a minor heterogeneity in the eastern Eurasian source, with
236 average affinity to the ANE ancestry is intermediate between Ulchi and Nganasan. Based on this
237 admixture modeling, we suggest that the steppe-forest cline does not keep a detectable level of
238 contribution from the older clines, the sources of which have higher ANE ancestry in both western and
239 eastern Eurasian parts.

240 In contrast, the southern steppe populations do not match with the Ulchi+Srubnaya model ($\chi^2 p \leq$
241 1.34×10^{-7} ; Supplementary Table 6). Adding Chalcolithic Iranians as the third ancestry significantly
242 improves model fit with substantial contribution from them ($\chi^2 p \leq 5.10 \times 10^{-5}$ with 7.0-64.6% contribution;
243 Fig. 5 and Supplementary Table 6), although the three-way model still does not adequately explain data.
244 Ancient individuals from the Tian Shan region²², dated to 2,200-1,100 BP, show a similar pattern
245 (Supplementary Table 7). However, older individuals from Central Kazakhstan dated to 2,500 BP
246 (“Saka_Kazakhstan_2500BP”)²² are adequately modeled as Nganasan+Srubnaya or Ulchi+Srubnaya+AG3
247 ($\chi^2 p = 0.057$ and 0.824 , respectively; Supplementary Table 7).

248 For the forest-tundra populations, the Nganasan+Srubnaya model is adequate only for the two
249 Volga region populations, Udmurts and Besermians (Fig. 5 and Supplementary Table 8). For the other

250 populations west of the Urals, six from the northeastern corner of Europe are modeled with additional
251 Mesolithic western European hunter-gatherers (“WHG”) contribution (8.2-11.4%; [Supplementary Table 8](#)),
252 while the rest need both WHG and early Neolithic European farmers (EEF; represented by “LBK_EN”;
253 [Supplementary Table 2](#))^{5,21}. Nganasan-related ancestry substantially contributes to their gene pools and
254 cannot be removed from the model without a significant decrease in model fit (4.1% to 29.0% contribution;
255 $\chi^2 p \leq 1.68 \times 10^{-5}$; [Supplementary Table 8](#)). For the four populations east of the Urals (Enets, Selkups, Kets
256 and Mansi), for which the above models are not adequate, Nganasan+Srubnaya+AG3 provide a good fit
257 ($\chi^2 p \geq 0.018$; [Fig. 5](#) and [Supplementary Table 8](#)). Substituting Nganasan to early Bronze Age populations
258 from the Baikal Lake region (“Baikal_EBA”; [Supplementary Table 2](#))²³, the two-way model of
259 Baikal_EBA+Srubnaya provides a reasonable fit ($\chi^2 p \geq 0.016$; [Supplementary Table 8](#)) and three-way
260 model of Baikal_EBA+Srubnaya+AG3 are adequate but with negative AG3 contribution for Enets and
261 Mansi ($\chi^2 p \geq 0.460$; [Supplementary Table 8](#)). Bronze/Iron Age populations from southern Siberia also
262 show a similar ancestry composition with high ANE affinity ([Supplementary Table 9](#)). The additional
263 ANE contribution beyond the Nganasan+Srubnaya model suggests a legacy from ANE-ancestry-rich
264 clines prior to Late Bronze Age.

265

266

267 **Discussion**

268 In this study, we analyzed new genome-wide data of indigenous peoples from inner Eurasia,
269 providing a dense representation for human genetic diversity in this vast region. Our finding of inner
270 Eurasian populations being structured into three largely distinct clines shows a striking correlation
271 between genes, geography and language ([Figs. 1-2](#)). Ecoregion-wide, the three clines match boreal forests
272 and tundra, forest-steppe zone and steppe/shrub-land further to the south, respectively. Language-wide,
273 they match the distribution of the Uralic, northern and southern Turkic-speaking languages. We
274 acknowledge that the distinction of three clines is far from complete and that there are cases of
275 intermediate patterns. For example, Turkic- and Uralic-speakers from the Volga region are genetically

276 quite similar, but the Uralic speakers still have extra affinity with the Uralic speakers further to the east
277 (e.g. Nganasans; [Supplementary Fig. 4b](#)). Likewise, a number of Turkic-speaking populations (e.g.
278 Dolgans, Todzins, Tofalars and Tatar_Zabolotniye), living at the periphery or even inside of the taiga belt,
279 do show a genetic influence from the forest-tundra cline ([Fig. 4](#)).

280 It may be viewed that our sampling scheme is not uniform geographically, although gathering the
281 vast majority of ethnic groups and quite dense geographically. Indeed, the gaps between distinct genetic
282 clines (with only a few groups located in between) tend to correspond to the gaps in sampling locations
283 ([Fig. 1-2](#)). Although this non-uniformity of sampling largely results from the non-uniformity in the density
284 of (language-defined) ethnic groups, it is important to organize a future study for further sampling on
285 sparsely populated regions between the clines (e.g. central Kazakhstan or East Siberia).

286 The steppe cline populations derive their eastern Eurasian ancestry from a gene pool similar to
287 contemporary Tungusic speakers from the Amur river basin ([Figs. 2 and 4](#)), thus suggesting a genetic
288 connection among the speakers of languages belonging to the Altaic macrofamily (Turkic, Mongolic and
289 Tungusic families). Based on our results as well as early Neolithic genomes from the Russian Far East³⁸,
290 we speculate that such a gene pool may represent the genetic profile of prehistoric hunter-gatherers in the
291 Amur river basin. On the other hand, a distinct Nganasan-related eastern Eurasian ancestry in the forest-
292 tundra cline suggests a substantial separation between these two eastern ancestries. Nganasans have high
293 genetic affinity with prehistoric individuals with the “ANE” ancestry in North Eurasia, such as the Upper
294 Paleolithic Siberians or the Mesolithic EHG, which is exceeded only by Native Americans and by
295 Beringians among eastern Eurasians ([Supplementary Fig. 7](#)). Also, Northeast Asians are closer to
296 Nganasans than they are to either Beringians, Native Americans or ancient Baikal populations, and the
297 ANE affinity in East Asians is correlated well with their affinity with Nganasans ([Supplementary Fig. 8](#)).
298 We hypothesize that Nganasans may be relatively isolated descendants of a prehistoric Siberian meta-
299 population with high ANE affinity, which formed present-day Northeast Asians by mixing with
300 populations related to the Neolithic Northeast Asians³⁸.

301 Forest-tundra populations to the east of the Urals, such as Selkups and Kets, show excess ANE
302 affinity, suggesting a legacy from the ANE-ancestry-rich pre-Bronze Age gene pools ([Supplementary](#)
303 [Table 8](#)). In contrast, admixture modeling finds that no contemporary steppe-forest cline population is
304 required to have additional ANE ancestry beyond what a mixture model of Bronze Age steppe plus
305 present-day Eastern Eurasians can explain ([Supplementary Table 5](#)). This suggests that both western and
306 eastern Eurasian ancestries of the steppe-forest populations are largely inherited from later gene flows
307 since Late Bronze Age: Srubnaya-like WSH ancestry for the western Eurasian part and present-day
308 Tungusic speaker-related ancestry for the eastern Eurasian part. Additional ancient genomes from Siberia
309 will be critical to reconstruct changes in the ANE-related ancestries in Siberia over time and to understand
310 the formation of Nganasan gene pool.

311 The southern steppe populations differentiate from the steppe-forest ones to the north by having a
312 strong genetic affinity broadly to West/ South Asian ancestries ([Supplementary Fig. 4](#) and [Supplementary](#)
313 [Table 6](#)). Ancient Tian Shan populations dating back up to 2,200 BP show the same property
314 ([Supplementary Table 7](#)), while Sintashta culture-related WSH ancestry was widely reported in this region
315 during the Late Bronze Age⁴⁶. Together with the lack of West/South Asian affinity in the Saka culture
316 individuals in Kazakhstan around 2,500 BP ([Supplementary Table 7](#)), we suggest a northward influx of
317 West/South Asian-related ancestry into the Tian Shan region during the first half of the first millennium
318 BC and into Kazakhstan further to the north slightly later.

319 It will be extremely important to expand the set of available ancient genomes across inner Eurasia.
320 Inner Eurasia has functioned as a conduit for human migration and cultural transfer since the first
321 appearance of modern humans in this region. As a result, we observe deep sharing of genes between
322 western and eastern Eurasian populations in multiple layers: the Pleistocene ANE ancestry in Mesolithic
323 EHG and contemporary Native Americans, Bronze Age steppe ancestry from Europe to Mongolia, and
324 Nganasan-related ancestry extending from western Siberia into Eastern Europe. More recent historical
325 migrations, such as the westward expansions of Turkic and Mongolic groups, further complicate genomic
326 signatures of admixture and have overwritten those from older events. Ancient genomes of Iron Age

327 steppe individuals, already showing signatures of west-east admixture in the 5th to 2nd century BC³⁹,
328 provide further direct evidence for the hidden old layers of admixture, which is often difficult to appreciate
329 from present-day populations as shown in our finding of a discrepancy between the estimates of admixture
330 dates from contemporary individuals and those from ancient genomes.

331

332

333 **Methods**

334

335 **Study participants and genotyping.** We collected samples from 763 participants from nine countries
336 (Armenia, Georgia, Kazakhstan, Moldova, Mongolia, Russia, Tajikistan, Ukraine, and Uzbekistan). The
337 sampling strategy included sampling a majority of large ethnic groups in the studied countries. Within
338 groups, we sampled subgroups if they were known to speak different dialects; for ethnic groups with large
339 area, we sampled within several districts across the area. We sampled individuals whose grandparents
340 were all self-identified members of the given ethnic groups and were born within the studied district(s).
341 Most of the ethnic Russian samples were collected from indigenous Russian areas (present-day Central
342 Russia) and had been stored for years in the Estonian Biocenter; samples from Mongolia, Tajikistan,
343 Uzbekistan, and Ukraine were collected partially in the framework of the Genographic project. Most DNA
344 samples were extracted from venous blood via the phenol-chloroform method. For this study we identified
345 112 subgroups (belonging to 60 ethnic group labels) which were not previously genotyped on the
346 Affymetrix Axiom® Genome-wide Human Origins 1 (“HumanOrigins”) array platform⁴³ and selected on
347 average 7 individuals per subgroup (Fig. 1 and Supplementary Table 1). Genome-wide genotyping
348 experiments were performed on the HumanOrigins array platform. We removed 18 individuals from
349 further analysis either due to high genotype missing rate (> 0.05 ; $n=2$) or due to being outliers in principal
350 component analysis (PCA) relative to other individuals from the same group ($n=16$). The remaining 745
351 individuals assigned to 60 group labels were merged to published HumanOrigins data sets of world-wide
352 contemporary populations²⁰ and of four Siberian ethnic groups (Enets, Kets, Nganasans and Selkups)²⁵.

353 Diploid genotype data of six contemporary individuals (two Saami, two Sherpa and two Tibetans) were
354 obtained from the Simons Genome Diversity Panel data set²⁶. We also added ancient individuals from
355 published studies^{3,8,19-23,27-42}, by randomly sampling a single allele for 581,230 autosomal single nucleotide
356 polymorphisms (SNPs) in the HumanOrigins array ([Supplementary Table 2](#)).

357
358 **Sequencing of the ancient Botai genomes.** We extracted genomic DNA from four skeletal remains
359 belonging to two individuals and built sequencing libraries either with no uracil-DNA glycosylase (UDG)
360 treatment or with partial treatment following published protocols^{47,48} ([Table 1](#)). Radiocarbon dating of
361 BKZ001 was conducted by the CEZ Archaeometry gGmbH (Mannheim, Germany) for one of two bone
362 samples used for DNA extraction. All libraries were barcoded with two library-specific 8-mer indices⁴⁹.
363 The samples were manipulated in dedicated clean room facilities at the University of Tübingen or at the
364 Max Planck Institute for the Science of Human History (MPI-SHH). Indexed libraries were enriched for
365 about 1.24 million informative nuclear SNPs using the in-solution capture method (“1240K capture”)^{5,21}.

366 Libraries were sequenced on the Illumina HiSeq 4000 platform with either single-end 75 bp (SE75)
367 or paired-end 50 bp (PE50) cycles following manufacturer’s protocols. Output reads were demultiplexed
368 by allowing up to 1 mismatch in each of two 8-mer indices. FASTQ files were processed using EAGER
369 v1.92⁵⁰. Specifically, Illumina adapter sequences were trimmed using AdapterRemoval v2.2.0⁵¹, aligned
370 reads (30 base pairs or longer) onto the human reference genome (hg19) using BWA aln/samse v0.7.12⁵²
371 with relaxed edit distance parameter (“-n 0.01”). Seeding was disabled for reads from non-UDG libraries
372 by adding an additional parameter (“-l 9999”). PCR duplicates were then removed using DeDup v0.12.2⁵⁰
373 and reads with Phred-scaled mapping quality score < 30 were filtered out using Samtools v1.3⁵³. We did
374 several measurements to check data authenticity. First, patterns of chemical damages typical to ancient
375 DNA were tabulated using mapDamage v2.0.6⁵⁴. Second, mitochondrial contamination for all libraries
376 was estimated by Schmutzi⁵⁵. Third, nuclear contamination for libraries derived from males was estimated
377 by the contamination module in ANGSD v0.910⁵⁶. Prior to genotyping, the first and last 3 bases of each
378 read were masked for libraries with partial UDG treatment using the trimBam module in bamUtil

379 v1.0.13⁵⁷. To obtain haploid genotypes, we randomly chose one high-quality base (Phred-scaled base
380 quality score ≥ 30) for each of the 1.24 million target sites using pileupCaller
381 (<https://github.com/stschiff/sequenceTools>). We used masked reads from libraries with partial UDG
382 treatment for transition (Ts) SNPs and used unmasked reads from all libraries for transversions (Tv).
383 Mitochondrial consensus sequences were obtained by the log2fasta program in Schmutzi with the quality
384 cutoff 10 and subsequently assigned to haplogroups using HaploGrep2⁵⁸. Y haplogroup R1b was assigned
385 using the yHaplo program⁵⁹. To estimate the phylogenetic position of the Botai Y haplogroup more
386 precisely, Y chromosomal SNPs were called with Samtools mpileup using bases with quality score ≥ 30 : a
387 total of 2,481 SNPs out of $\sim 30,000$ markers included in the 1240K capture panel were called with mean
388 read depth of 1.2. Twenty-two SNP positions relevant to the up-to-date haplogroup R1b tree
389 (www.isogg.org; www.yfull.com) confirmed that the sample was positive for the markers of R1b-P297
390 branch but negative for its R1b-M269 sub-branch.

391 The frequency distribution map of this Y chromosomal clade was created by the GeneGeo
392 software^{60,61} using the average weighed interpolation procedure with the weight function of degree 3 and
393 radius 1,200 km. The initial frequencies were calculated as proportion of samples positive for “root” R1b
394 marker M343 but negative for M269; these proportions were calculated for the 577 populations from the
395 in-home *Y-base* database, which was compiled mainly from the published datasets.

396
397 **Analysis of population structure.** We performed principal component analysis (PCA) of various groups
398 using smartpca v13050 in the EIGENSOFT v6.0.1 package⁶². We used the “*lsqproject: YES*” option to
399 project individuals not used for calculating PCs (this procedure avoids bias due to missing genotypes). We
400 performed unsupervised model-based genetic clustering as implemented in ADMIXTURE v1.3.0⁶³. For
401 that purpose, we used 118,387 SNPs with minor allele frequency (maf) 1% or higher in 3,507 individuals
402 after pruning out linked SNPs ($r^2 > 0.2$) using the “--indep-pairwise 200 25 0.2” command in PLINK
403 v1.90⁶⁴. For each value of K ranging from 2 to 20, we ran 5 replicates with different random seeds and
404 took one with the highest log likelihood value.

405
406 **F-statistics analysis.** We computed various f_3 and f_4 statistics using the qp3Pop (v400) and qpDstat (v711)
407 programs in the ADMIXTOOLS package⁴³. We computed f_4 -statistics with the “*f4mode: YES*” option. For
408 these analyses, we studied a total of 301 groups, including 73 inner Eurasian target groups and 167
409 contemporary and 93 ancient reference groups (Supplementary Table 2). We included two groups from the
410 Aleutian Islands (“Aleut” and “Aleut_Tlingit”; Supplementary Table 2) as positive control targets with
411 known recent admixture. Aleut_Tlingits are Aleut individuals whose mitochondrial haplogroup lineages
412 are related to Tlingits³¹. For each target, we calculated outgroup f_3 statistic of the form $f_3(\text{Target}, X; \text{Mbuti})$
413 against all targets and references to quantify overall allele sharing and performed admixture f_3 test of the
414 form $f_3(\text{Ref}_1, \text{Ref}_2; \text{Target})$ for all pairs of references to explore the admixture signal in targets. We
415 estimated standard error (SE) using a block jackknife with 5 centiMorgan (cM) block⁶².

416 We performed f_4 statistic-based admixture modeling using the qpAdm (v632) program²⁰ in the
417 ADMIXTOOLS package. We used a basic set of 7 outgroups, unless specified otherwise, to provide high
418 enough resolution to distinguish various western and eastern Eurasian ancestries: Mbuti (n=10; central
419 African), Natufian (n=6; early Holocene Levantine)²⁰, Onge (n=11; from the Andaman Islands), Iran_N
420 (n=5; Neolithic Iranian)²⁰, Villabruna (n=1; Paleolithic European)²⁸, Ami (n=10; Taiwanese aborigine) and
421 Mixe (n=10; Central American). Prior to qpAdm modeling, we checked if the reference groups are well
422 distinguished by their relationship with the outgroups using the qpWave (v400) program⁶⁵.

423 We used the qpGraph (v6065) program in the ADMIXTOOLS package for graph-based admixture
424 modeling. Starting with a graph of (Mbuti, Ami, WHG), we iteratively added AG3 (n=1; Paleolithic
425 Siberian)²⁸, EHG (n=4; Mesolithic hunter-gatherers from Karelia or Samara)^{5,23,28}, and Botai onto the
426 graph by testing all possible topologies allowing up to one additional gene flow. After obtaining the best
427 two-way admixture model for Botai, we tested additional three-way admixture models.

428
429 **GLOBETROTTER analysis.** We performed a GLOBETROTTER analysis of admixture for 73 inner
430 Eurasian target populations to obtain haplotype sharing based evidence of admixture, independent of the

431 allele frequency based f -statistics, as well as estimates of admixture dates and a fine-scale profile of their
432 admixture sources¹⁴. We followed the “regional” approach described in Hellenthal et al.¹⁴, in which target
433 haplotypes can only be copied from the haplotypes of 167 contemporary reference groups, but not from
434 those of the other target groups. This approach is recommended when multiple target groups share a
435 similar admixture history¹⁴, which is likely to be the case for our inner Eurasian populations.

436 We jointly phased the contemporary genome data without a pre-phased set of reference haplotypes,
437 using SHAPEIT2 v2.837 in its default setting⁶⁶. We used a genetic map for the 1000 Genomes Project
438 phase 3 data, downloaded from: https://mathgen.stats.ox.ac.uk/impute/1000GP_Phase3.html. We used
439 haplotypes from a total of 2,615 individuals belonging to 240 groups (73 recipients and 167 donors;
440 [Supplementary Table 2](#)) for the GLOBETROTTER analysis. To reduce computational burden and to
441 provide more balanced set of donor populations, we randomly sampled 20 individuals if a group contained
442 more than 20 individuals. Using these haplotypes, we performed GLOBETROTTER analysis following
443 the recommended workflow¹⁴. We first ran 10 rounds of the expectation-maximization (EM) algorithm for
444 chromosomes 4, 10, 15 and 22 in ChromoPainter v2 with “-in” and “-iM” switches to estimate chunk size
445 and switch error rate parameters⁶⁷. Both recipient and donor haplotypes were modeled as a patchwork of
446 donor haplotypes. The “chunk length” output was obtained by running ChromoPainter v2 across all
447 chromosomes with the estimated parameters averaged over both recipient and donor individuals (“-n
448 238.05 -M 0.000617341”). We also generated 10 painting samples for each recipient group by running
449 ChromoPainter with the parameters averaged over all recipient individuals (“-n 248.455 -M
450 0.000535236”). Using the chunklength output and painting samples, we ran GLOBETROTTER with the
451 “prop.ind: 1” and “null.ind: 1” options. We estimated significance of estimated admixture date by running
452 100 bootstrap replicates using the “prop.ind: 0” and “bootstrap.date.ind: 1” options; we considered date
453 estimates between 1 and 400 generations as evidence of admixture¹⁴. For populations that gave evidence
454 of admixture by this procedure, we repeated GLOBETROTTER analysis with the “null.ind: 0” option¹⁴.
455 We also compared admixture dates from GLOBETROTTER analysis with those based on weighted
456 admixture linkage disequilibrium (LD) decay, as implemented in ALDER v1.3⁶⁸. As the reference pair, we

457 used (French, Eskimo_Naukan), (French, Nganasan), (Georgian, Ulchi), (French, Ulchi) and (Georgian,
458 Ulchi) for the target group categories 1 to 5, respectively, based on their genetic profile ([Supplementary
459 Table 2](#)). We used a minimum inter-marker distance of 1.0 cM to account for LD in the references.

460
461 **EEMS analysis.** To visualize the heterogeneity in the rate of gene flow across inner Eurasia, we
462 performed the EEMS (“estimated effective migration surface”) analysis⁴⁴. We included a total of 1,214
463 individuals from 98 groups in the analysis ([Supplementary Table 2](#)). In this dataset, we kept 101,370 SNPs
464 with $\text{maf} \geq 0.01$ after LD pruning ($r^2 \leq 0.2$). We computed the mean squared genetic difference matrix
465 between all pairs of individuals using the “bed2diffs_v1” program in the EEMS package. To reduce
466 distortion in northern latitudes due to map projection, we used geographic coordinates in the Albers equal
467 area conic projection (“+proj=aea +lat_1=50 +lat_2=70 +lat_0=56 +lon_0=100 +x_0=0 +y_0=0
468 +ellps=WGS84 +datum=WGS84 +units=m +no_defs”). We converted geographic coordinates of each
469 sample and the boundary using the “spTransform” function in the R package rgdal v1.2-5. We ran five
470 initial MCMC runs of 2 million burn-ins and 4 million iterations with different random seeds and took a
471 run with the highest likelihood. Starting from the best initial run, we set up another five MCMC runs of 2
472 million burn-ins and 4 million iterations as our final analysis. We used the following proposal variance
473 parameters to keep the acceptance rate around 30-40%, as recommended by the developers⁴⁴:
474 $\text{qSeedsProposalS2} = 5000$, $\text{mSeedsProposalS2} = 1000$, $\text{qEffctProposalS2} = 0.0001$, $\text{mrateMuProposalS2} =$
475 0.00005 . We set up a total of 532 demes automatically with the “nDemes = 600” parameter. We visualized
476 the merged output from all five runs using the “eems.plots” function in the R package rEEMSplots⁴⁴.

477 We performed the EEMS analysis for Caucasus populations in a similar manner, including a total
478 of 237 individuals from 21 groups ([Supplementary Table 2](#)). In this dataset, we kept 95,442 SNPs with
479 $\text{maf} \geq 0.01$ after LD pruning ($r^2 \leq 0.2$). We applied the Mercator projection of geographic coordinates to
480 the map of Eurasia (“+proj=merc +datum=WGS84”). We ran five initial MCMC runs of 2 million burn-
481 ins and 4 million iterations with different random seeds and took a run with the highest likelihood. Starting
482 from the best initial run, we set up another five MCMC runs of 1 million burn-in and 4 million iterations

483 as our final analysis. We used the default following proposal variance parameters: $qSeedsProposalS2 = 0.1$,
484 $mSeedsProposalS2 = 0.01$, $qEffctProposalS2 = 0.001$, $mrRateMuProposalS2 = 0.01$. A total of 171 demes
485 were automatically set up with the “nDemes = 200” parameter.

486

487 **Life Science Reporting Summary.** Further information on experimental design is available in the Life
488 Sciences Reporting Summary.

489

490 **Ethics Statement.** The study protocol was approved by the Ethics Committee of the Research Centre for
491 Medical Genetics, Moscow, Russia. All 763 participants who contributed their genetic materials provided
492 a signed written informed consent.

493

494 **Data Availability.** Genome-wide sequence data of two Botai individuals (BAM format) are available at
495 the European Nucleotide Archive under the accession number PRJEB31152 (ERP113669). Eigenstrat-
496 format array genotype data of 763 present-day individuals and 1240K pulldown genotype data of two
497 ancient Botai individuals are available at the Edmond data repository of the Max Planck Society
498 (<https://edmond.mpdl.mpg.de/imeji/collection/Aoh9c69DscnxSNjm?q=>).

499

500

501 **References**

- 502 1 Li, J. Z. *et al.* Worldwide human relationships inferred from genome-wide patterns of variation.
 503 *Science* **319**, 1100-1104 (2008).
- 504 2 Wang, C., Zöllner, S. & Rosenberg, N. A. A quantitative comparison of the similarity between
 505 genes and geography in worldwide human populations. *PLoS Genet.* **8**, e1002886 (2012).
- 506 3 Jeong, C. *et al.* Long-term genetic stability and a high altitude East Asian origin for the peoples of
 507 the high valleys of the Himalayan arc. *Proc. Natl. Acad. Sci. USA* **113**, 7485-7490 (2016).
- 508 4 Yunusbayev, B. *et al.* The Caucasus as an asymmetric semipermeable barrier to ancient human
 509 migrations. *Mol. Biol. Evol.* **29**, 359-365 (2012).
- 510 5 Haak, W. *et al.* Massive migration from the steppe was a source for Indo-European languages in
 511 Europe. *Nature* **522**, 207-211 (2015).
- 512 6 Haber, M. *et al.* Genome-wide diversity in the Levant reveals recent structuring by culture. *PLoS*
 513 *Genet.* **9**, e1003316 (2013).
- 514 7 Martiniano, R. *et al.* The population genomics of archaeological transition in west Iberia:
 515 Investigation of ancient substructure using imputation and haplotype-based methods. *PLoS Genet.*
 516 **13**, e1006852 (2017).
- 517 8 Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature* **522**, 167-172 (2015).
- 518 9 Barfield, T. J. *The nomadic alternative.* (Prentice Hall, Englewood Cliffs, NJ, 1993).
- 519 10 Frachetti, M. D. *Pastoralist landscapes and social interaction in Bronze Age Eurasia.* (Univ of
 520 California Press, Berkeley, CA, 2009).
- 521 11 Burch, E. S. The caribou/wild reindeer as a human resource. *Am. Antiquity* **37**, 339-368 (1972).
- 522 12 Sherratt, A. The secondary exploitation of animals in the Old World. *World Archaeol.* **15**, 90-104
 523 (1983).
- 524 13 Yunusbayev, B. *et al.* The genetic legacy of the expansion of Turkic-speaking nomads across
 525 Eurasia. *PLoS Genet.* **11**, e1005068 (2015).
- 526 14 Hellenthal, G. *et al.* A genetic atlas of human admixture history. *Science* **343**, 747-751 (2014).
- 527 15 Flegontov, P. *et al.* Genomic study of the Ket: a Paleo-Eskimo-related ethnic group with
 528 significant ancient North Eurasian ancestry. *Sci. Rep.* **6**, 20768 (2016).
- 529 16 Pugach, I. *et al.* The complex admixture history and recent southern origins of Siberian
 530 populations. *Mol. Biol. Evol.* **33**, 1777-1795 (2016).
- 531 17 Triska, P. *et al.* Between Lake Baikal and the Baltic Sea: genomic history of the gateway to
 532 Europe. *BMC Genet.* **18**, 110 (2017).
- 533 18 Tambets, K. *et al.* Genes reveal traces of common recent demographic history for most of the
 534 Uralic-speaking populations. *Genome Biology* **19**, 139 (2018).
- 535 19 Raghavan, M. *et al.* Upper Palaeolithic Siberian genome reveals dual ancestry of Native
 536 Americans. *Nature* **505**, 87-91 (2014).
- 537 20 Lazaridis, I. *et al.* Genomic insights into the origin of farming in the ancient Near East. *Nature*
 538 **536**, 419-424 (2016).
- 539 21 Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient Eurasians. *Nature* **528**, 499-
 540 503 (2015).
- 541 22 Damgaard, P. d. B. *et al.* 137 ancient human genomes from across the Eurasian steppes. *Nature*
 542 **557**, 369-374 (2018).
- 543 23 Damgaard, P. d. B. *et al.* The first horse herders and the impact of early Bronze Age steppe
 544 expansions into Asia. *Science*, 10.1126/science.aar7711 (2018).
- 545 24 Levine, M. & Kislenko, A. New Eneolithic and early Bronze Age radiocarbon dates for north
 546 Kazakhstan and south Siberia. *Camb. Archaeol.* **7**, 297-300 (1997).
- 547 25 Flegontov, P. *et al.* Paleo-Eskimo genetic legacy across North America. *bioRxiv*, 203018 (2017).
- 548 26 Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142 diverse
 549 populations. *Nature* **538**, 201-206 (2016).

550 27 Fu, Q. *et al.* Genome sequence of a 45,000-year-old modern human from western Siberia. *Nature*
551 **514**, 445-449 (2014).

552 28 Fu, Q. *et al.* The genetic history of Ice Age Europe. *Nature* **534**, 200-205 (2016).

553 29 Haber, M. *et al.* Continuity and admixture in the last five millennia of Levantine history from
554 ancient Canaanite and present-day Lebanese genome sequences. *Am. J. Hum. Genet.* **101**, 274-282
555 (2017).

556 30 Jones, E. R. *et al.* Upper Palaeolithic genomes reveal deep roots of modern Eurasians. *Nat.*
557 *Commun.* **6**, 8912 (2015).

558 31 Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for present-day
559 Europeans. *Nature* **513**, 409-413 (2014).

560 32 Lazaridis, I. *et al.* Genetic origins of the Minoans and Mycenaeans. *Nature* **548**, 214-218 (2017).

561 33 Raghavan, M. *et al.* The genetic prehistory of the New World Arctic. *Science* **345**, 1255832
562 (2014).

563 34 Rasmussen, M. *et al.* The genome of a Late Pleistocene human from a Clovis burial site in
564 western Montana. *Nature* **506**, 225-229 (2014).

565 35 Rasmussen, M. *et al.* Ancient human genome sequence of an extinct Palaeo-Eskimo. *Nature* **463**,
566 757-762 (2010).

567 36 Rasmussen, M. *et al.* The ancestry and affiliations of Kennewick Man. *Nature* **523**, 455-458
568 (2015).

569 37 Saag, L. *et al.* Extensive farming in Estonia started through a sex-biased migration from the
570 Steppe. *Curr. Biol.* **27**, 2185-2193. e2186 (2017).

571 38 Siska, V. *et al.* Genome-wide data from two early Neolithic East Asian individuals dating to 7700
572 years ago. *Sci. Adv.* **3**, e1601877 (2017).

573 39 Unterländer, M. *et al.* Ancestry and demography and descendants of Iron Age nomads of the
574 Eurasian Steppe. *Nat. Commun.* **8**, 14615 (2017).

575 40 Yang, M. A. *et al.* 40,000-year-old individual from Asia provides insight into early population
576 structure in Eurasia. *Curr. Biol.* **27**, 3202-3208. e3209 (2017).

577 41 Kılınç, G. M. *et al.* The demographic development of the first farmers in Anatolia. *Curr. Biol.* **26**,
578 2659-2666 (2016).

579 42 McColl, H. *et al.* The prehistoric peopling of Southeast Asia. *Science* **361**, 88-92 (2018).

580 43 Patterson, N. *et al.* Ancient admixture in human history. *Genetics* **192**, 1065-1093 (2012).

581 44 Petkova, D., Novembre, J. & Stephens, M. Visualizing spatial population structure with estimated
582 effective migration surfaces. *Nat. Genet.* **48**, 94-100 (2016).

583 45 Fenner, J. N. Cross-cultural estimation of the human generation interval for use in genetics-
584 based population divergence studies. *Am. J. Phys. Anthropol.* **128**, 415-423 (2005).

585 46 Narasimhan, V. M. *et al.* The genomic formation of South and Central Asia. *bioRxiv*, 292581
586 (2018).

587 47 Dabney, J. *et al.* Complete mitochondrial genome sequence of a Middle Pleistocene cave bear
588 reconstructed from ultrashort DNA fragments. *Proc. Natl. Acad. Sci. USA* **110**, 15758-15763
589 (2013).

590 48 Rohland, N., Harney, E., Mallick, S., Nordenfelt, S. & Reich, D. Partial uracil-DNA-glycosylase
591 treatment for screening of ancient DNA. *Phil. Trans. R. Soc. B* **370**, 20130624 (2015).

592 49 Kircher, M. in *Ancient DNA: methods and protocols* (eds Beth Shapiro & Michael Hofreiter) 197-
593 228 (Humana Press, 2012).

594 50 Peltzer, A. *et al.* EAGER: efficient ancient genome reconstruction. *Genome Biol.* **17**, 60 (2016).

595 51 Schubert, M., Lindgreen, S. & Orlando, L. AdapterRemoval v2: rapid adapter trimming,
596 identification, and read merging. *BMC Res. Notes* **9**, 88 (2016).

597 52 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
598 *Bioinformatics* **25**, 1754-1760 (2009).

599 53 Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078-2079
600 (2009).

601 54 Jónsson, H., Ginolhac, A., Schubert, M., Johnson, P. L. & Orlando, L. mapDamage2.0: fast
602 approximate Bayesian estimates of ancient DNA damage parameters. *Bioinformatics* **29**, 1682-
603 1684 (2013).

604 55 Renaud, G., Slon, V., Duggan, A. T. & Kelso, J. Schmutzi: estimation of contamination and
605 endogenous mitochondrial consensus calling for ancient DNA. *Genome Biol.* **16**, 224 (2015).

606 56 Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. ANGSD: analysis of next generation
607 sequencing data. *BMC Bioinformatics* **15**, 356 (2014).

608 57 Jun, G., Wing, M. K., Abecasis, G. R. & Kang, H. M. An efficient and scalable analysis
609 framework for variant extraction and refinement from population-scale DNA sequence data.
610 *Genome Res.* **25**, 918-925 (2015).

611 58 Weissensteiner, H. *et al.* HaploGrep 2: mitochondrial haplogroup classification in the era of high-
612 throughput sequencing. *Nucleic Acids Res.* **44**, W58-W63 (2016).

613 59 Poznik, G. D. Identifying Y-chromosome haplogroups in arbitrarily large samples of sequenced or
614 genotyped men. *bioRxiv*, 088716 (2016).

615 60 Balanovsky, O. *et al.* Parallel evolution of genes and languages in the Caucasus region. *Mol. Biol.*
616 *Evol.* **28**, 2905-2920 (2011).

617 61 Koshel, S. in *Sovremennaya geograficheskaya kartografiya (Modern Geographic Cartography)*
618 (eds I. Lourie & V. Kravtsova) 158-166 (Data+, 2012).

619 62 Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS Genet.* **2**,
620 e190 (2006).

621 63 Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated
622 individuals. *Genome Res.* **19**, 1655-1664 (2009).

623 64 Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.
624 *Gigascience* **4**, 7 (2015).

625 65 Reich, D. *et al.* Reconstructing native American population history. *Nature* **488**, 370-374 (2012).

626 66 Delaneau, O., Zagury, J.-F. & Marchini, J. Improved whole-chromosome phasing for disease and
627 population genetic studies. *Nat. Methods* **10**, 5-6 (2013).

628 67 Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of population structure using
629 dense haplotype data. *PLoS Genet.* **8**, e1002453 (2012).

630 68 Loh, P.-R. *et al.* Inferring admixture histories of human populations using linkage disequilibrium.
631 *Genetics* **193**, 1233-1254 (2013).

632 69 Sedghifar, A., Brandvain, Y., Ralph, P. & Coop, G. The spatial mixing of genomes in secondary
633 contact zones. *Genetics* **201**, 243-261 (2015).

634 70 Levine, M. Botai and the origins of horse domestication. *J. Anthropol. Archaeol.* **18**, 29-78 (1999).

635 71 Bronk Ramsey, C. Bayesian analysis of radiocarbon dates. *Radiocarbon* **51**, 337-360 (2009).

636 72 Reimer, P. J. *et al.* IntCal13 and Marine13 radiocarbon age calibration curves 0–50,000 years cal
637 BP. *Radiocarbon* **55**, 1869-1887 (2016).

638

639

640 **Acknowledgements**

641 We thank Iain Mathieson and Iosif Lazaridis for their helpful comments. The research leading to these
642 results has received funding from the Max Planck Society, the Max Planck Society Donation Award and
643 the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation
644 programme (grant agreement No 646612 granted to M.R.). Analysis of the Caucasus dataset was
645 supported by RFBR grant 16-06-00364 and analysis of the Far East dataset was supported by Russian
646 Scientific Fund project 17-14-01345. D.R. was supported by the U.S. National Science Foundation
647 HOMINID grant BCS-1032255, the U.S. National Institutes of Health grant GM100233, by an Allen
648 Discovery Center grant, and is an investigator of the Howard Hughes Medical Institute. P.F. was
649 supported by IRP projects of the University of Ostrava and by the Czech Ministry of Education, Youth
650 and Sports (project OPVVV 16_019/0000759). C.C.W. was funded by Nanqiang Outstanding Young
651 Talents Program of Xiamen University and the Fundamental Research Funds for the Central Universities.
652 M.Z. has been funded by research grants from the MES RK No. AP05134955 and No. 0114RK00492.

653

654 **Author Contributions**

655 C.J., O.B., E.B., S.S., W.H., D.R., J.K. conceived and coordinated the study. O.B., M.L., E.P., Y.Y., A.A.,
656 K.S., A.Bu., P.N., S.T., D.Dal., M.C., R.S., D.Dar., Y.B., A.Bo., A.S., N.D., M.Z., L.Y., V.C., N.P., L.Da.,
657 L.S., K.D., L.A., O.U., E.I., E.Ka., I.E., M.M., E.B. contributed the present-day samples. N.K., O.I., E.Kh.,
658 B.B., V.Zai., L.Dj. A.K.O contributed the ancient Botai samples. N.K., A.I. performed ancient DNA
659 laboratory works. C.J., O.B., E.L., V.Zap., C.C.W. conducted population genetic analyses. C.J., O.B., S.S.,
660 W.H., J.K. wrote the paper with input from P.F., M.R., L.Dj., D.R. and co-authors.

661

662 **Competing Interests**

663 The authors declare no competing interests.

664

665 **Figure Legends**

666

667 **Fig. 1. Geographic locations of the Eneolithic Botai site (red triangle), 65 groups including newly**
668 **sampled individuals (filled diamonds) and nearby groups with published data (filled squares).** Mean
669 latitude and longitude values across all individuals under each group label were used. Two zoom-in plots
670 for the Caucasus (blue) and the Altai-Sayan (magenta) regions are presented in the lower left corner. A list
671 of new groups, their three-letter codes, and the number of new individuals (in parenthesis) are provided at
672 the bottom. Present-day populations are color-coded based on the language family for Figs. 1-3, following
673 key codes listed in Fig. 2. Corresponding information for the previously published groups is provided in
674 Supplementary Table 2. The map is overlaid with ecoregional information, divided into 14 biomes,
675 downloaded from <https://ecoregions2017.appspot.com/> (credited to Ecoregions 2017 © Resolve). The
676 main inner Eurasian map is on the Albers equal area projection and was produced using the spTransform
677 function in the R package rgdal v1.2-5.

678

679 **Fig. 2. The genetic structure of inner Eurasian populations.** (a) The first two PCs of 2,077 Eurasian
680 individuals separate western and eastern Eurasians (PC1) and Northeast and Southeast Asians (PC2). Most
681 inner Eurasians are located between western and eastern Eurasians on PC1. Ancient individuals (color-
682 filled shapes) are projected onto PCs calculated based on contemporary individuals. Present-day
683 individuals are marked by grey dots, with their per-group mean coordinates marked by three-letter codes
684 listed in Supplementary Table 2. Individuals are colored by their language family. (b) ADMIXTURE
685 results for a chosen set of ancient and present-day groups (K = 14). The top row shows ancient inner
686 Eurasians and representative present-day eastern Eurasians. The following three rows show forest-tundra,
687 steppe-forest and southern steppe cline populations. Most inner Eurasians are modeled as a mixture of
688 components primarily found in eastern or western Eurasians. Results for the full set of individuals are
689 provided in Supplementary Fig. 3.

690

691 **Fig. 3. Correlation of longitude and ancestry proportion across inner Eurasian populations.** Across
692 inner Eurasian populations, mean longitudinal coordinates (x-axis) and mean eastern Eurasian ancestry
693 proportions (y-axis) are strongly correlated. Eastern Eurasian ancestry proportions are estimated from
694 ADMIXTURE results with K=14 by summing up six components maximized in Surui, Chipewyan,
695 Itelmen, Nganasan, Atayal and early Neolithic Russian Far East individuals (“Devil’s Gate”), respectively
696 (Supplementary Fig. 3). The yellow curve shows a probit regression fit following the model in Sedghifar
697 et al.⁶⁹. Three groups (Kalmyks, Dungans, Nogai2) are marked with grey square due to their substantial
698 deviation from the curve as well as their historically known migration history.

699

700 **Fig. 4. Characterization of the western and eastern Eurasian source ancestries in inner Eurasian**
701 **populations.** (a) Admixture f_3 values are compared for different eastern Eurasian references (Mixe,
702 Nganasan, Ulchi; left) or western Eurasian ones (Srubnaya, Iran_ChL; right). For each target group, darker
703 shades mark more negative f_3 values. (b) Weights of donor populations in two sources characterizing the
704 main admixture signal (“date 1 PC 1”) in the GLOBETROTTER analysis. We merged 167 donor
705 populations into 12 groups, as listed on the top right side. Target populations are split into five groups:
706 Aleuts, the forest-tundra cline populations, the steppe-forest cline populations, the southern steppe cline
707 populations and the rest of four populations (“others”), from the top to bottom.

708

709 **Fig. 5. qpAdm-based admixture models for the forest-tundra and steppe-forest cline populations.**
710 For the forest-tundra population to the west of the Urals, Nganasan+Srubnaya+WHG+LBK_EN or its
711 submodel provides a good fit, while additional ANE-related contribution (AG3) is required for those to the
712 east of the Urals (Enets, Selkups, Kets, and Mansi). For the steppe-forest populations, Srubnaya+Ulchi,
713 Srubnaya+Ulchi+AG3, or Srubnaya+Nganasan provides a good fit. 5 cM jackknifing standard errors are
714 marked by the horizontal bar. Models with p -value between 0.01 and 0.05 are marked by grey color and

715 those with p -value < 0.01 are marked by grey color and italic font. Details of the model information are
716 presented in [Supplementary Tables 5 and 8](#).

717 **Table 1. Sequencing statistics and radiocarbon dates of two Eneolithic Botai individuals analyzed in this study.** For Botai individuals we
718 produced additional data, we provide corresponding individual ID from a previous publication²³ (“Published ID”), radiocarbon date, the number of
719 total reads sequenced, mean autosomal coverage for the 1240K target sites, the number of SNPs covered at least once for the 1240K and
720 HumanOrigins panels, uniparental haplogroup and contamination estimates.
721

ID	Published ID	Genetic Sex	Uncal. ¹⁴ C Date	Cal. ¹⁴ C Date (2-sigma) ^b	# of reads sequenced	Mean autosomal coverage	# of SNPs covered ^c	MT / Y haplogroup	MT.cont ^d	X.cont ^e
TU45	BOT14	M	4620 ± 80 ^a	3632-3100 cal. BCE	84,170,835	0.827x	169,053 (77,363)	K1b2 / R1b1a1	0.02 (0.01-0.03)	0.0122 (0.0050)
BKZ001	BOT2016	F	4660 ± 25	3517-3367 cal. BCE	69,678,735	2.420x	825,332 (432,078)	Z1 / NA	0.01 (0.00-0.02)	NA

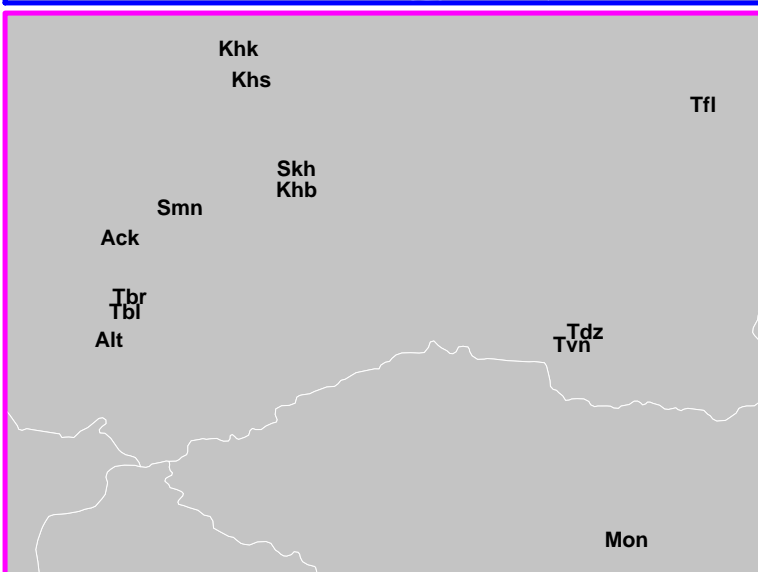
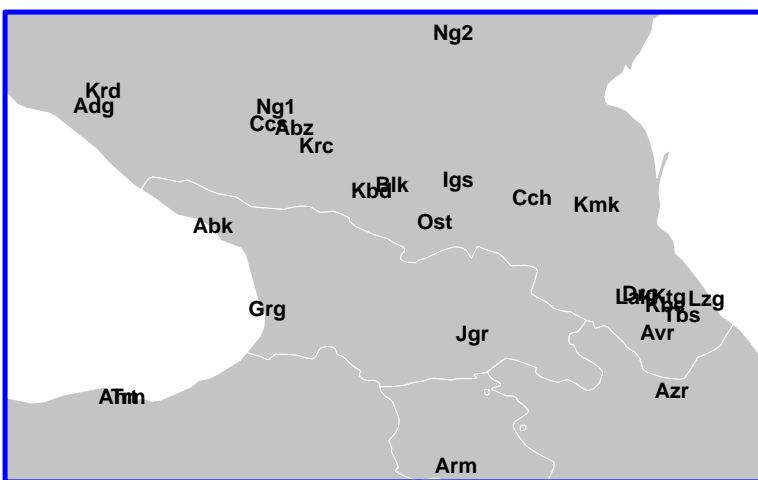
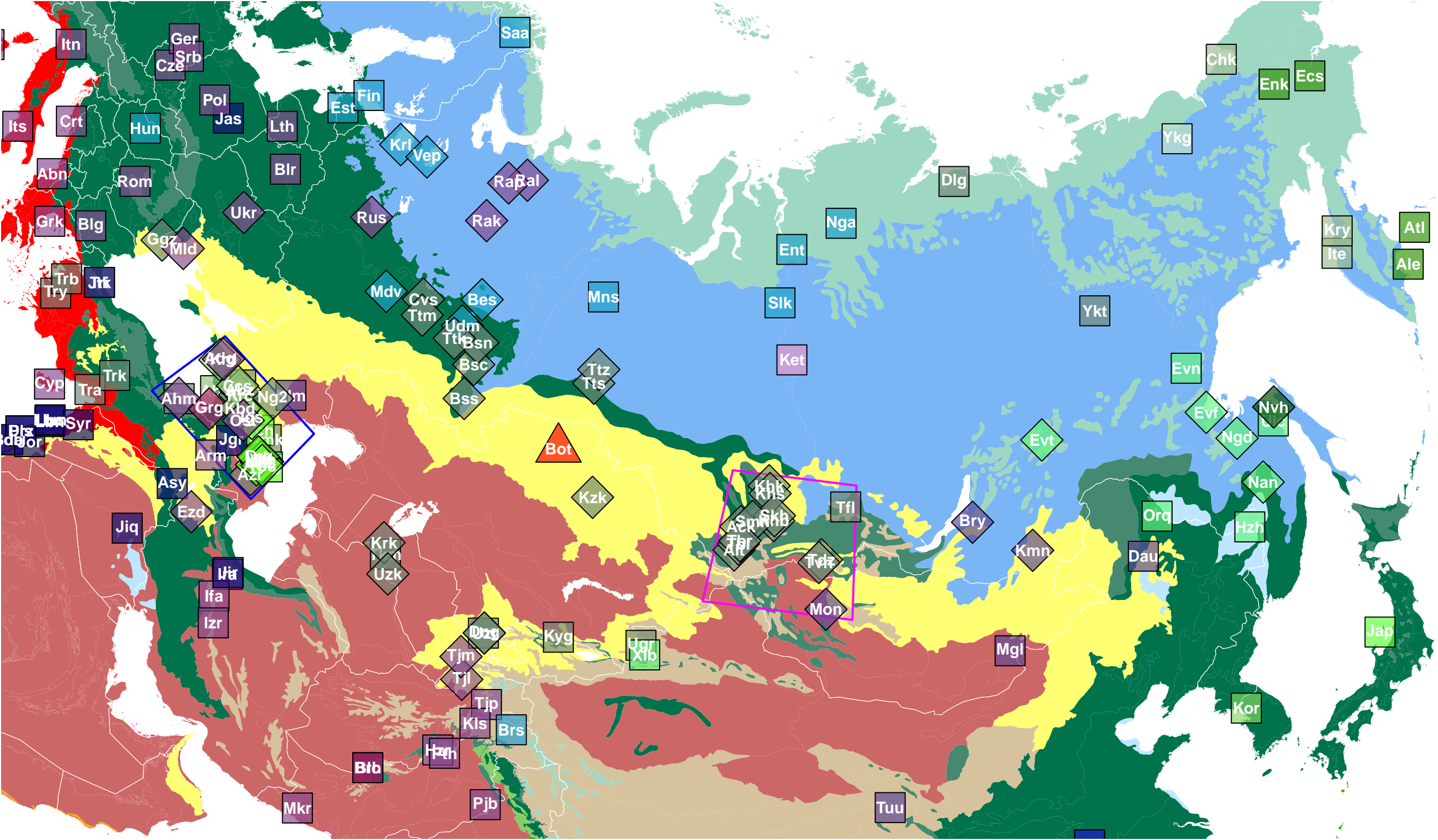
722 ^a The uncalibrated date of TU45 was published in Levine (1999) under the ID OxA-4316⁷⁰.

723 ^b The calibrated ¹⁴C dates are calculated based on uncalibrated dates, by the OxCal v4.3.2 program⁷¹ using the INTCAL13 atmospheric curve⁷².

724 ^c The number of SNPs in the 1240K panel (out of 1,233,013) or autosomal SNPs in the HumanOrigins array (out of 581,230; within the parenthesis) covered at
725 least by one read. Only transversion SNPs are considered for the non-UDG libraries (both of the TU45 libraries, one of two BKZ001 libraries).

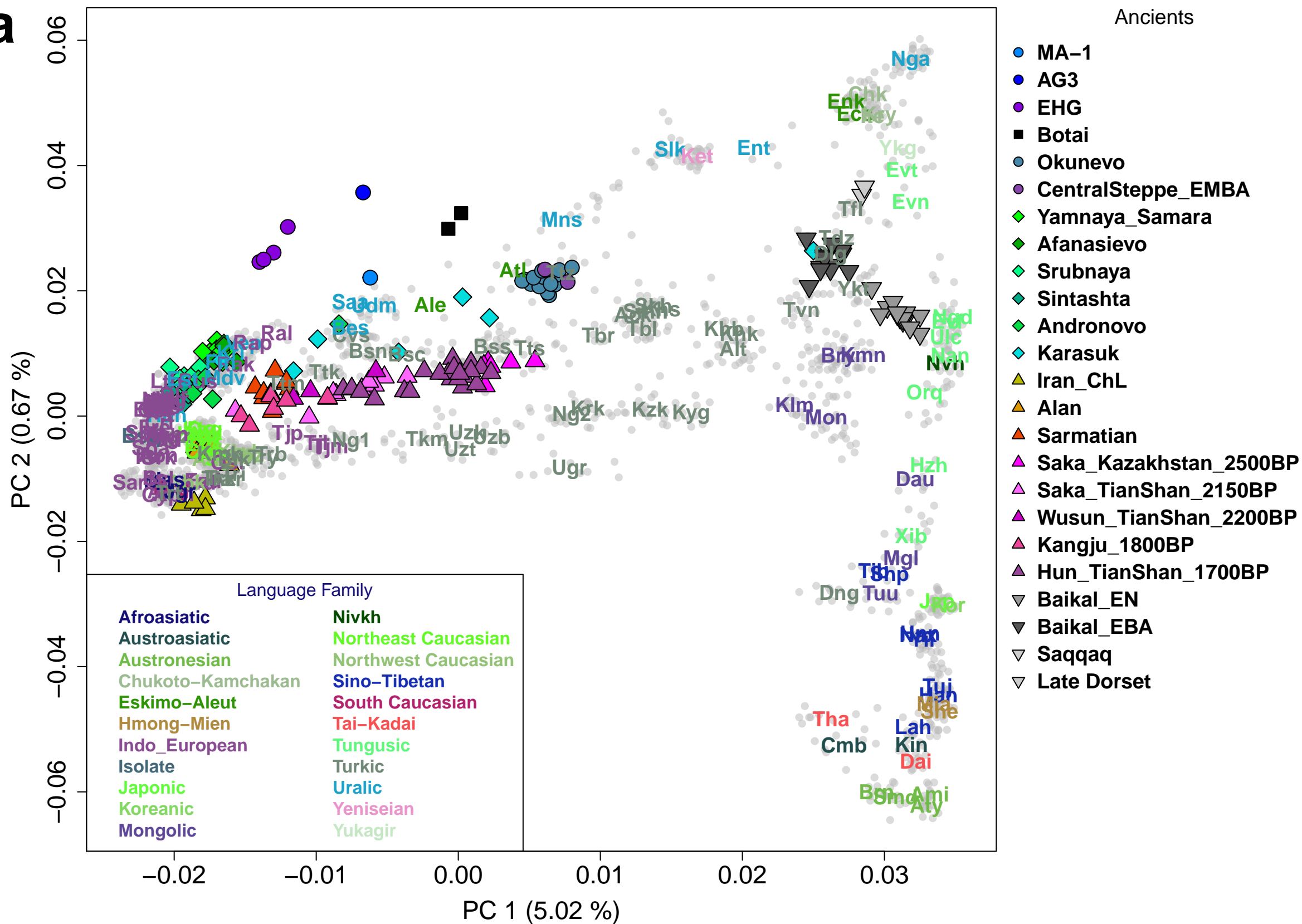
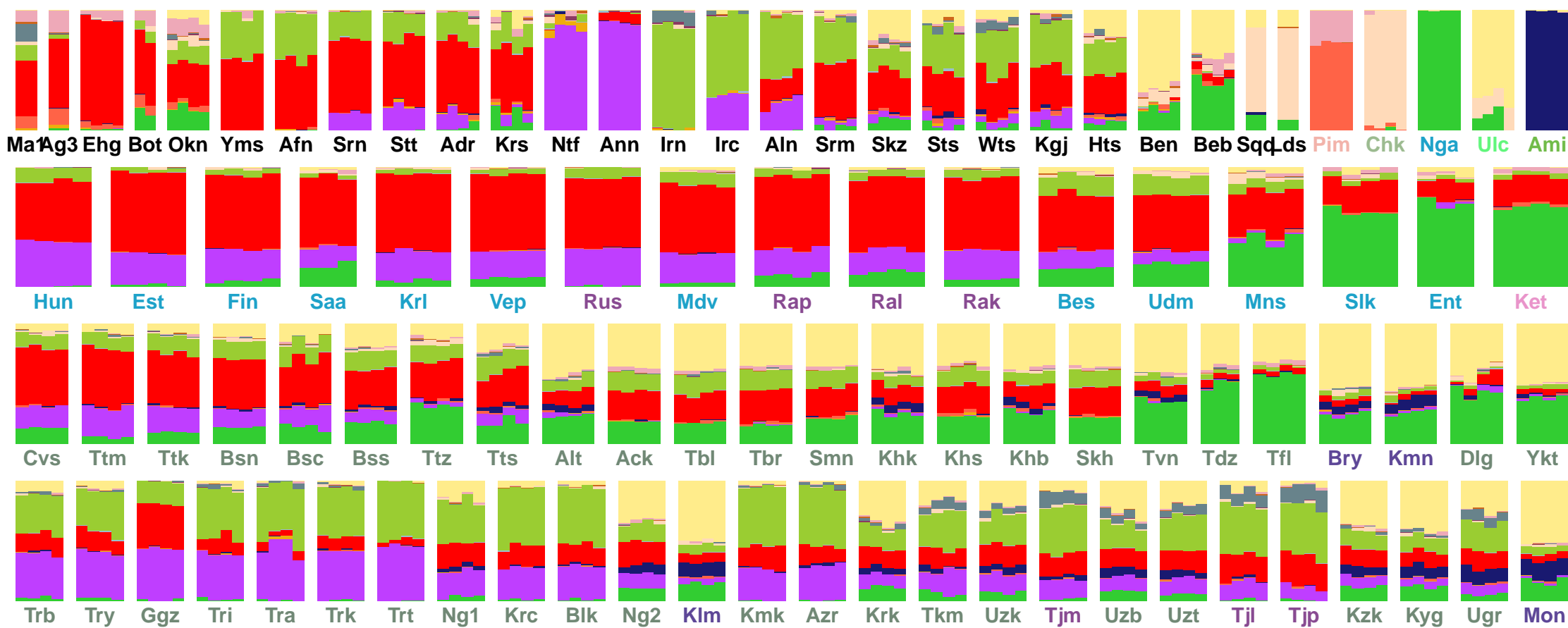
726 ^d The contamination rate of mitochondrial reads estimated by the Schmutzi program (95% confidence interval in parentheses)

727 ^e The nuclear contamination rate for the male (TU45) estimated based on X chromosome data by ANGSD software (standard error in parentheses)
728

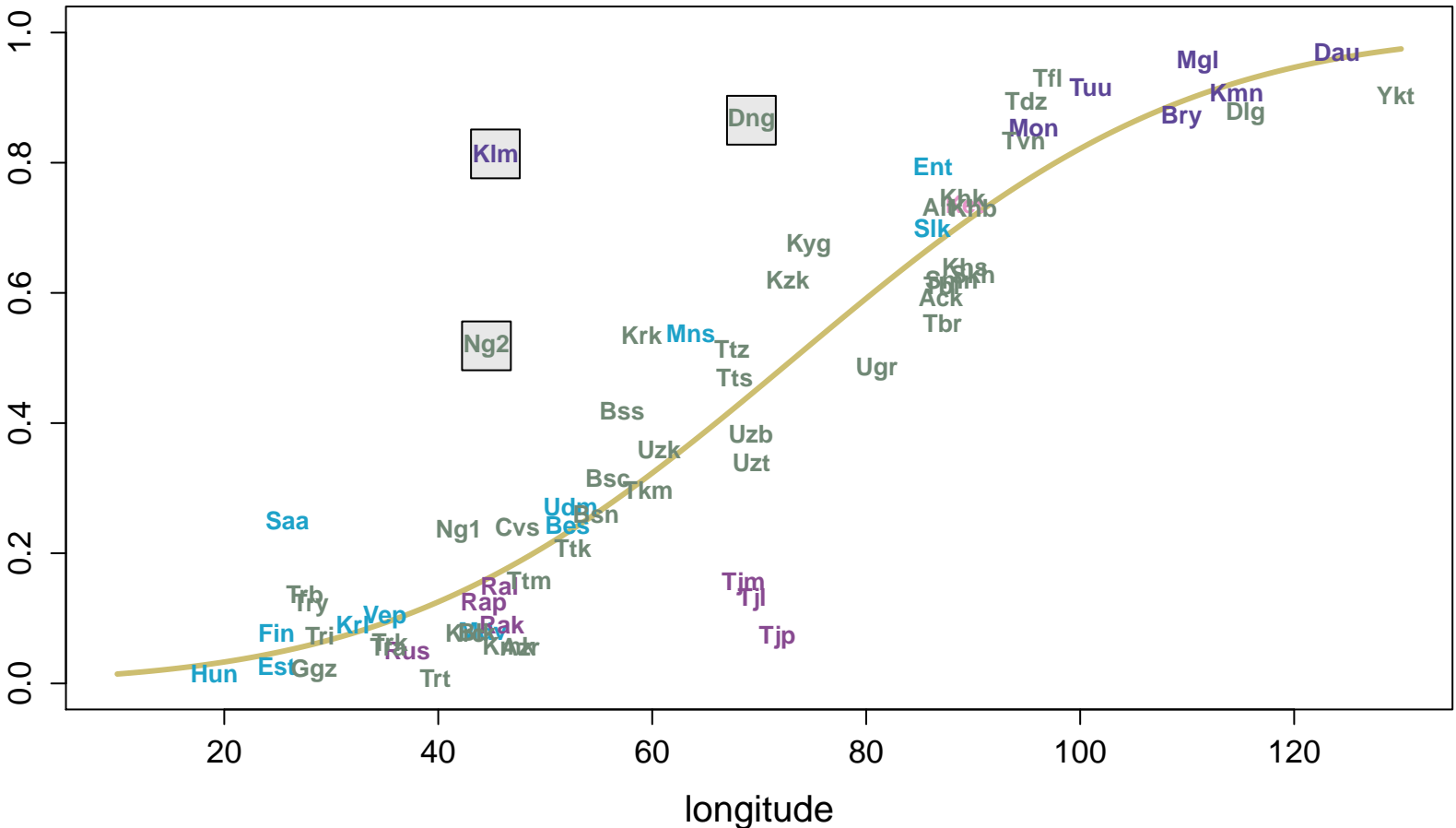


New groups

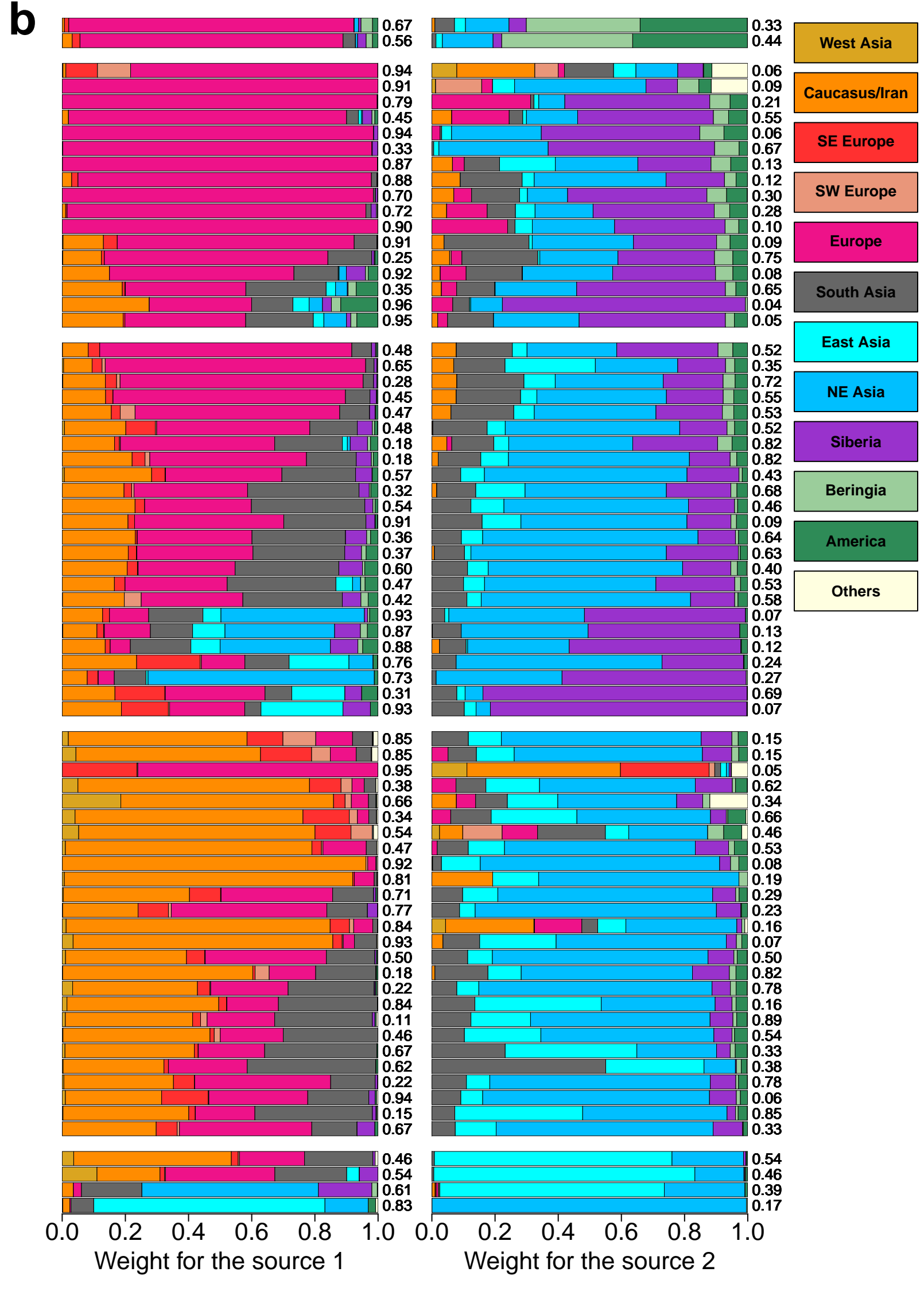
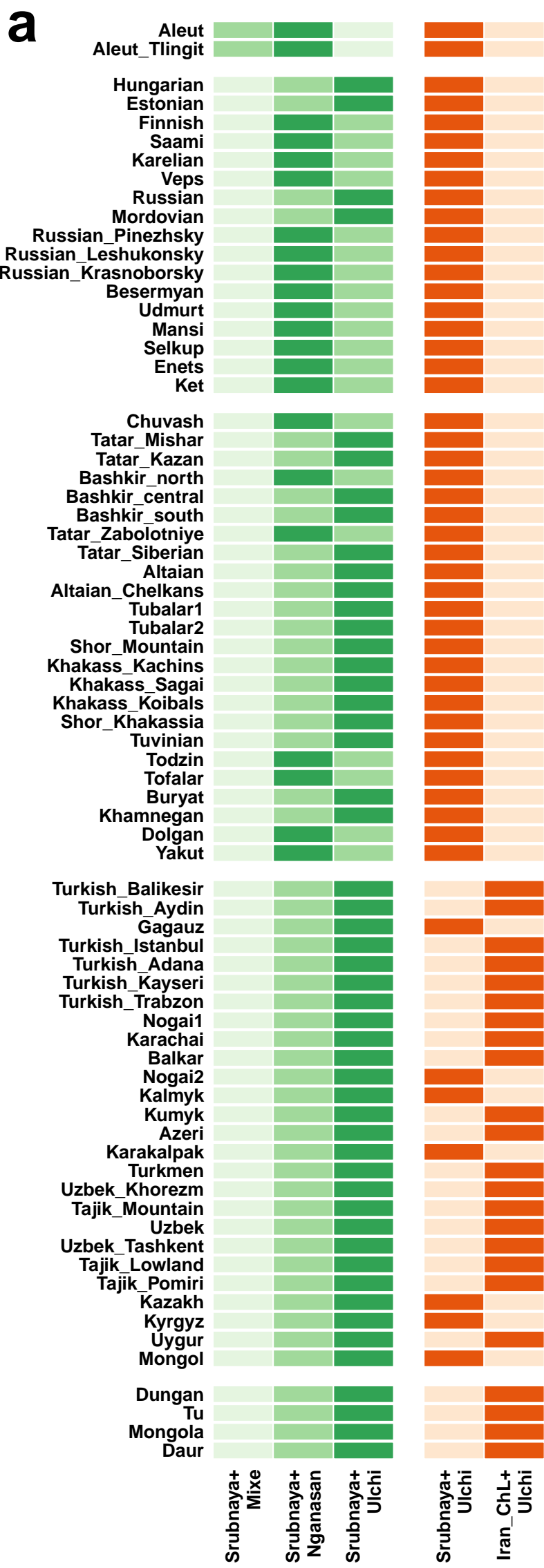
- Abz** Abazin (8)
- Ack** Altaian_Chelkans (6)
- Adg** Adygei (14)
- Ahm** Armenian_Hemsheni (7)
- Alt** Altaian (17)
- Avr** Avar (9)
- Azr** Azeri (17)
- Bes** Besermyan (6)
- Bry** Buryat (36)
- Bsc** Bashkir_central (16)
- Bsn** Bashkir_north (18)
- Bss** Bashkir_south (19)
- Ccs** Circassian (9)
- Cvs** Chuvash (4)
- Dng** Dungan (13)
- Drg** Darginian (8)
- Evf** Evenk_FarEast (2)
- Evt** Evenk_Transbaikal (8)
- Ezd** Ezid (8)
- Ggz** Gagauz (7)
- Grg** Georgian (12)
- Igs** Ingushian (10)
- Kbc** Kubachinian (6)
- Kbd** Kabardinian (8)
- Khb** Khakass_Koibals (5)
- Khk** Khakass_Kachins (7)
- Khs** Khakass_Sagai (9)
- Kmn** Khamnegan (8)
- Krc** Karachai (11)
- Krd** Kurd (8)
- Krk** Karakalpak (14)
- Krl** Karelian (15)
- Ktg** Kaitag (8)
- Kzk** Kazakh (18)
- Lak** Lak (10)
- Mdv** Mordovian (22)
- Mld** Moldavian (10)
- Mon** Mongol (34)
- Nan** Nanai (10)
- Ng2** Nogai2 (13)
- Ngd** Negidal (3)
- Nvh** Nivh (10)
- Ost** Ossetian (6)
- Rak** Russian_Krasnoborsky (6)
- Ral** Russian_Leshukonsky (5)
- Rap** Russian_Pinezhsky (5)
- Rus** Russian (49)
- Skh** Shor_Khakassia (5)
- Smn** Shor_Mountain (6)
- Tbl** Tubalar1 (3)
- Tbr** Tubalar2 (2)
- Tbs** Tabasaran (10)
- Tdz** Todzin (3)
- Tjl** Tajik_Lowland (11)
- Tjm** Tajik_Mountain (12)
- Ttk** Tatar_Kazan (13)
- Ttm** Tatar_Mishar (10)
- Tts** Tatar_Siberian (18)
- Ttz** Tatar_Zabolotniye (5)
- Tvn** Tuvinian (10)
- Udm** Udmurt (10)
- Ukr** Ukrainian (12)
- Uzk** Uzbek_Khorezm (6)
- Uzt** Uzbek_Tashkent (9)
- Vep** Veps (9)
- Bot** Botai (2)

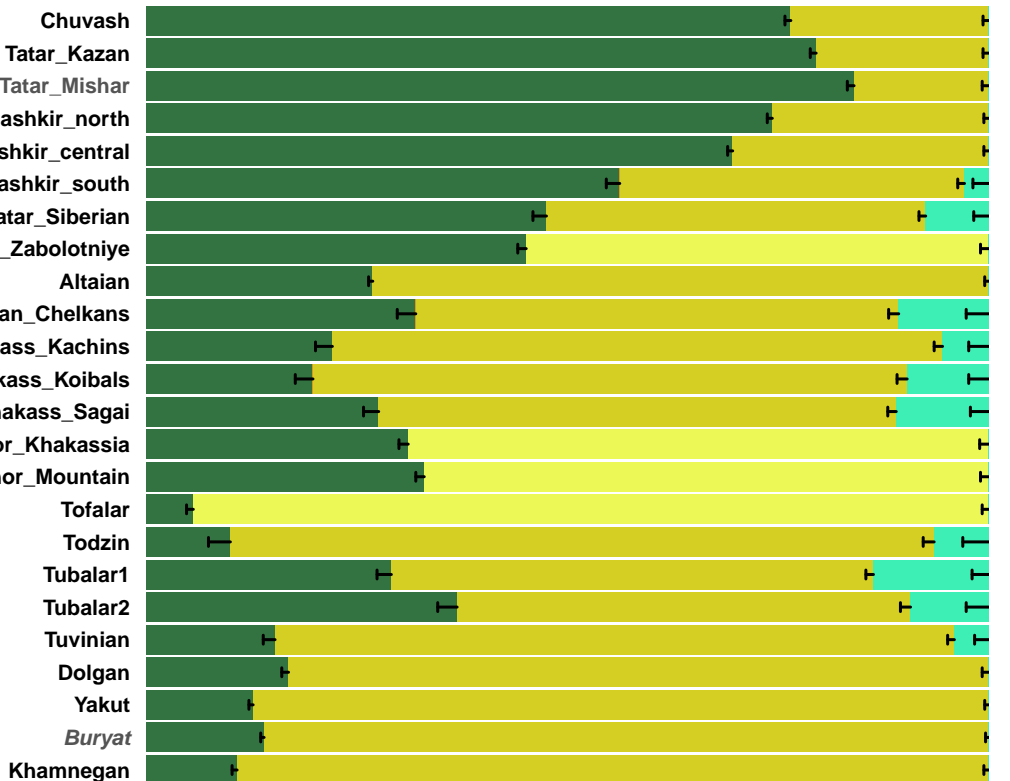
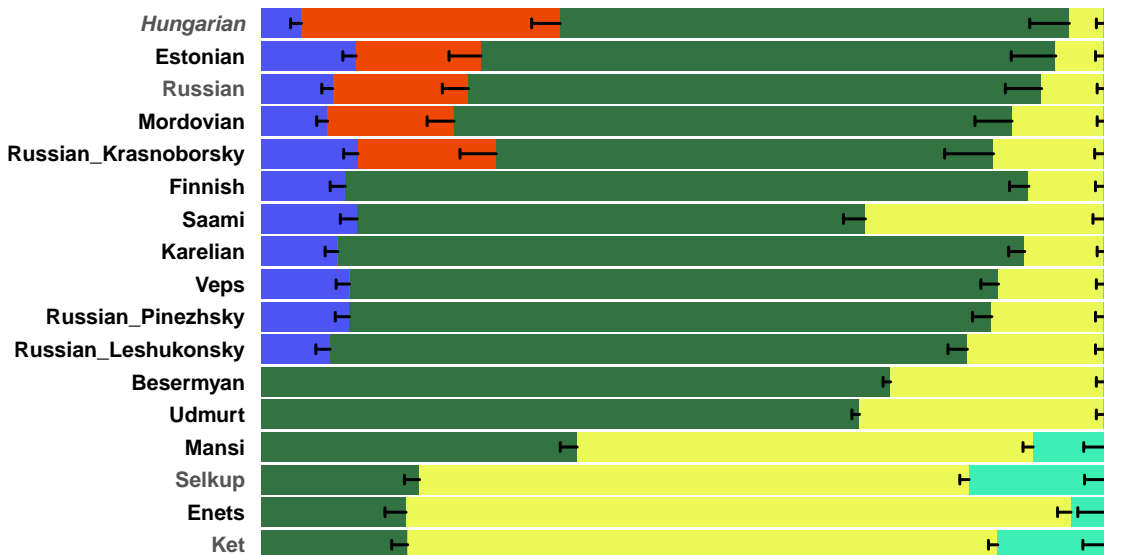
a**b**

Eastern Eurasian ancestry



- Bes Besermyan
- Ent Enets
- Est Estonian
- Fin Finnish
- Hun Hungarian
- Kri Karelian
- Mns Mansi
- Mdv Mordovian
- Saa Saami
- Sik Selkup
- Udm Udmurt
- Vep Veps
- Ket Ket
- Alt Altaian
- Ack Altaian_Chelkans
- Bsc Bashkir_central
- Bsn Bashkir_north
- Bss Bashkir_south
- Cvs Chuvash
- Dlg Dolgan
- Dng Dungan
- Khk Khakass_Kachins
- Khb Khakass_Koibals
- Khs Khakass_Sagai
- Shk Shor_Khakassia
- Smn Shor_Mountain
- Ttk Tatar_Kazan
- Ttm Tatar_Mishar
- Tts Tatar_Siberian
- Ttz Tatar_Zabolotniye
- Tdz Todzin
- Tbl Tubalar1
- Tbr Tubalar2
- Tvn Tuvinian
- Ykt Yakut
- Azr Azeri
- Blk Balkar
- Ggz Gagauz
- Krc Karachai
- Krk Karakalpak
- Kzk Kazakh
- Kmk Kumyk
- Kyg Kyrgyz
- Ng1 Nogai1
- Ng2 Nogai2
- Tra Turkish_Adana
- Try Turkish_Aydin
- Trb Turkish_Balikesir
- Tri Turkish_Istanbul
- Trk Turkish_Kayseri
- Trt Turkish_Trabzon
- Tkm Turkmen
- Ugr Uygur
- Uzb Uzbek
- Uzk Uzbek_Khorezm
- Uzt Uzbek_Tashkent
- Bry Buryat
- Dau Daur
- Klm Kalmyk
- Kmn Khamnegan
- Mon Mongol
- Mgl Mongola
- Tuu Tu
- Rus Russian
- Rak Russian_Krasnoborsky
- Ral Russian_Leshukonsky
- Rap Russian_Pinezhsky
- Tjl Tajik_Lowland
- Tjm Tajik_Mountain
- Tjp Tajik_Pomiri



WHG**LBK_EN****Srubnaya****Ulchi****Nganasan****AG3**

0.0 0.2 0.4 0.6 0.8 1.0

Ancestry proportion