# The Genomic Formation of Human Populations in East Asia

Chuan-Chao Wang[1,2,3,4,*], Hui-Yuan Yeh[5,*], Alexander N Popov[6,*], Hu-Qin Zhang[7,*], Hirofumi Matsumura[8], Kendra Sirak[2,9], Olivia Cheronet[10], Alexey Kovalev[11], Nadin Rohland[2], Alexander M. Kim[2,12], Rebecca Bernardos[2], Dashtseveg Tumen[13], Jing Zhao[7], Yi-Chang Liu[14], Jiun-Yu Liu[15], Matthew Mah[2,16,17], Swapan Mallick[2, 9,16,17], Ke Wang[3], Zhao Zhang[2], Nicole Adamski[2,17], Nasreen Broomandkhoshbacht[2,17], Kimberly Callan[2,17], Brendan J. Culleton[18], Laurie Eccles[19], Ann Marie Lawson[2,17], Megan Michel[2,17], Jonas Oppenheimer[2,17], Kristin Stewardson[2,17], Shaoqing Wen[20], Shi Yan[21], Fatma Zalzala[2,17], Richard Chuang[14], Ching-Jung Huang[14], Chung-Ching Shiung[14], Yuri G. Nikitin[22], Andrei V. Tabarev[23], Alexey A. Tishkin[24], Song Lin[7], Zhou-Yong Sun[25], Xiao-Ming Wu[7], Tie-Lin Yang[7], Xi Hu[7], Liang Chen[26], Hua Du[27], Jamsranjav Bayarsaikhan[28], Enkhbayar Mijiddorj[29], Diimaajav Erdenebaatar[29], Tumur-Ochir Iderkhangai[29], Erdene Myagmar[13], Hideaki Kanzawa-Kiriyama[30], Msato Nishino[31], Ken-ichi Shinoda[30], Olga A. Shubina[32], Jianxin Guo[1], Qiongying Deng[33], Longli Kang[34], Dawei Li[35], Dongna Li[36], Rong Lin[36], Wangwei Cai[37], Rukesh Shrestha[4], Ling-Xiang Wang[4], Lanhai Wei[1], Guangmao Xie[38,39], Hongbing Yao[40], Manfei Zhang[4], Guanglin He[1], Xiaomin Yang[1], Rong Hu[1], Martine Robbeets[3], Stephan Schiffels[3], Douglas J. Kennett[41], Li Jin[4], Hui Li[4], Johannes Krause[3], Ron Pinhasi[10], David Reich[2,9,16,17]

1.  Department of Anthropology and Ethnology, Institute of Anthropology, School of Sociology and Anthropology and State Key Laboratory of Cellular Stress Biology, School of Life Sciences, Xiamen University, Xiamen 361005, China
2.  Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA
3.  Max Planck Institute for the Science of Human History, 07745 Jena, Germany
4.  MOE Key Laboratory of Contemporary Anthropology, Department of Anthropology and Human Genetics, School of Life Sciences, Fudan University, Shanghai 200438, China
5.  School of Humanities, Nanyang Technological University, Nanyang 639798, Singapore
6.  Scientific Museum, Far Eastern Federal University, 690950 Vladivostok, Russia
7.  Key Laboratory of Biomedical Information Engineering of Ministry of Education, School of Life Science and Technology, Xi'an Jiaotong University, Xi'an 710049, China
8.  School of Health Science, Sapporo Medical University, S1 W17, Chuo-ku, Sapporo, 060-8556, Japan
9.  Department of Human Evolutionary Biology, Harvard Unviersity, Cambridge, MA 02138, USA
10. Department of Evolutionary Anthropology, University of Vienna, 1090 Vienna, Austria
11. Institute of Archaeology, Russian Academy of Sciences, Moscow, Russia
12. Department of Anthropology, Harvard University, Cambridge, Massachusetts 02138, USA
13. Department of Anthropology and Archaeology, National University of Mongolia, Ulaanbaatar 46, Mongolia
14. Institute of Archaeology, National Cheng Kung University, Tainan 701, Taiwan

15. Department of Anthropology, University of Washington, 314 Denny Hall, Seattle, USA
16. Broad Institute of Harvard and MIT, Cambridge, MA, 02142, USA
17. Howard Hughes Medical Institute, Harvard Medical School, Boston, MA 02115, USA
18. Institutes of Energy and the Environment, The Pennsylvania State University, University Park, PA 16802, USA.
19. Department of Anthropology, Pennsylvania State University, University Park, PA 16802, USA
20. Institute of Archaeological Science, Fudan University, Shanghai 200433, China
21. School of Ethnology and Sociology, Minzu University of China, Beijing 100081, China
22. Museum of Archaeology and Ethnology of Institute of History, Far Eastern Branch of Russian Academic of Sciences, Vladivostok 690001, Russia
23. Institute of Archaeology and Ethnography, Siberian Branch of Russian Academy of Sciences, Novosibirsk 630090, Russia
24. Department of Archeology, Ethnography and Museology, Altai State University, Barnaul, Altaisky Kray 656049, Russia
25. Shaanxi Provincial Institute of Archaeology, Xi'an 710054, China
26. College of Cultural Heritage, Northwest University, Xi'an 710069, China
27. Xi'an AMS Center, Institute of Earth Environment, Chinese Academy of Sciences, Xi'an 710061, China
28. Research Center at the National Museum of Mongolia, Ulaanbaatar, Region of Sukhbaatar 14201, Mongolia
29. Department of Archaeology, Ulaanbaatar State University, Ulaanbaatar, Region of Bayanzurkh 13343, Mongolia
30. Department of Anthropology, National Museum of Nature and Science, Tsukuba City, Ibaraki Prefecture 305-0005, Japan
31. Archeological Center of Chiba City, Chiba 260-0814 Japan
32. Department of Archeology, Sakhalin Regional Museum, Yuzhno-Sakhalinsk, Russia
33. Department of Human Anatomy and Center for Genomics and Personalized Medicine, Guangxi Medical University, Nanning 530021, China
34. Key Laboratory for Molecular Genetic Mechanisms and Intervention Research on High Altitude Disease of Tibet Autonomous Region, Key Laboratory of High Altitude Environment and Gene Related to Disease of Tibet, Ministry of Education, School of Medicine, Xizang Minzu University, Xianyang 712082, Shaanxi, China
35. Guangxi Museum of Nationalities, Nanning 530028, Guangxi, China
36. Department of Biology, Hainan Medical University, Haikou 571199, Hainan, China
37. Department of Biochemistry and Molecular Biology, Hainan Medical University, Haikou 571199, Hainan, China
38. College of History, Culture and Tourism, Guangxi Normal University, Guilin 541001, China
39. Guangxi Institute of Cultural Relics Protection and Archaeology, Nanning 530003, Guangxi, China
40. Belt and Road Research Center for Forensic Molecular Anthropology, Key Laboratory of Evidence Science of Gansu Province, Gansu Institute of Political Science and Law, Lanzhou 730070, China
41. Department of Anthropology, University of California, Santa Barbara, CA 93106, USA

* Contributed equally.

Correspondence: wang@xmu.edu.cn (C-C.W), krause@shh.mpg.de (J.K.), ron.pinhasi@univie.ac.at (R.P.), and reich@genetics.med.harvard.edu (D.R.).

**The deep population history of East Asia remains poorly understood due to a lack of ancient DNA data and sparse sampling of present-day people. We report genome-wide data from 191 individuals from Mongolia, northern China, Taiwan, the Amur River Basin and Japan dating to 6000 BCE - 1000 CE, many from contexts never previously analyzed with ancient DNA. We also report 383 present-day individuals from 46 groups mostly from the Tibetan Plateau and southern China. We document how 6000-3600 BCE people of Mongolia and the Amur River Basin were from populations that expanded over Northeast Asia, likely dispersing the ancestors of Mongolic and Tungusic languages. In a time transect of 89 Mongolians, we reveal how Yamnaya steppe pastoralist spread from the west by 3300-2900 BCE in association with the Afanasievo culture, although we also document a boy buried in an Afanasievo barrow with ancestry entirely from local Mongolian hunter-gatherers, representing a unique case of someone of entirely non-Yamnaya ancestry interred in this way. The second spread of Yamnaya-derived ancestry came via groups that harbored about a third of their ancestry from European farmers, which nearly completely displaced unmixed Yamnaya-related lineages in Mongolia in the second millennium BCE, but did not replace Afanasievo lineages in western China where Afanasievo ancestry persisted, plausibly acting as the source of the early-splitting Tocharian branch of Indo-European languages. Analyzing 20 Yellow River Basin farmers dating to ~3000 BCE, we document a population that was a plausible vector for the spread of Sino-Tibetan languages both to the Tibetan Plateau and to the central plain where they mixed with southern agriculturalists to form the ancestors of Han Chinese. We show that the individuals in a time transect of 52 ancient Taiwan individuals spanning at least 1400 BCE to 600 CE were consistent with being nearly direct descendants of Yangtze Valley first farmers who likely spread Austronesian, Tai-Kadai and Austroasiatic languages across Southeast and South Asia and mixing with the people they encountered, contributing to a four-fold reduction of genetic differentiation during the emergence of complex societies. We finally report data from Jomon hunter-gatherers from Japan who harbored one of the earliest splitting branches of East Eurasian variation, and show an affinity among Jomon, Amur River Basin, ancient Taiwan, and Austronesian-speakers, as expected for ancestry if they all had contributions from a Late Pleistocene coastal route migration to East Asia.**

**Main text**

East Asia, one of the oldest centers of animal and plant domestication, today harbors more than a fifth of the world's human population, with present-day groups speaking

126    languages representing eleven major families: Sino-Tibetan, Tai-Kadai, Austronesian,

127    Austroasiatic, Hmong-Mien, Indo-European, Altaic (Mongolic, Turkic, and

128    Tungusic), Koreanic, Japonic, Yukgahiric, and Chukotko-Kanchatkan[1]. The past

129    10,000 years have been a period of profound economic and cultural change in East

130    Asia, but our current understanding of the genetic diversity, major mixture events, and

131    population movements and turnovers during the transition from foraging to

132    agriculture remains poor due to minimal sampling of the diversity of present-day

133    people on the Tibetan Plateau and southern China[2]. A particular limitation has been a

134    deficiency in ancient DNA data, which has been a powerful tool for discerning the

135    deep history of populations in Western and Central Eurasia[3-8].

136

137    We genotyped 383 present-day individuals from 46 populations indigenous to China

138    (n=337) and Nepal (n=46) using the Affymetrix Human Origins array (Table S1 and

139    Supplementary Information section 1). We also report genome-wide data from 191

140    ancient East Asians, many from cultural contexts for which there is no published

141    ancient DNA data. From Mongolia we report 89 individuals from 52 sites dating

142    between ~6000 BCE to ~1000 CE. From China we report 20 individuals from the

143    ~3000 BCE Neolithic site of Wuzhuangguoliang. From Japan we report 7 Jomon

144    hunter-gatherers from 3500-1500 BCE. From the Russian Far East we report 23

145    individuals: 18 from the Neolithic Boisman-2 cemetery at ~5000 BCE, 1 from the

146    Iron Age Yankovsky culture at ~1000 BCE, 3 from the Medieval Heishui Mohe and

147    Bohai Mohe culture at ~1000 CE; and 1 historic period hunter-gatherer from Sakhalin

148    Island. From archaeological sites in Eastern Taiwan—the Bilhun site at Hanben on the

149    main island and the Gongguan site on Green Island—we report 52 individuals from

150    the Late Neolithic through the Iron Age spanning at least 1400 BCE - 600 CE.

151

152    For all but the Chinese samples we enriched the ancient DNA for a targeted set of

153    about 1.2 million single nucleotide polymorphisms (SNPs)[4,9], while for the

154    Wuzhuangguoliang samples from China we used exome capture (18 individuals) or

155    shotgun sequencing (2 individuals) (Figure 1, Supplementary Data files 1 and 2 and

156    Supplementary Information section 1). We performed quality control to test for

157    contamination by other human sequences, assessed by the rate of cytosine to thymine

158    substitution in the terminal nucleotide and polymorphism in mitochondrial DNA

159    sequences[10] as well as X chromosome sequences in males, and restricted analysis to

160    individuals with minimal contamination[11] (Online Table 1). We detected close kinship

161    between individuals at the same site, including a Boisman nuclear family with 2

162    parents and 4 children (Table S2). We merged the new data with previously reported

163    data: 4 Jomon individuals, 8 Amur River Basin Neolithic individuals from the Devil's

164    Gate site, 72 individuals from the Neolithic to the Iron Age in Southeast Asia, and 8

165    from Nepal[7,12-20]. We assembled 123 radiocarbon dates using bone from the

166    individuals, of which 94 are newly reported (Online Table 3), and clustered

167    individuals based on time period and cultural associations, then further by genetic

168    cluster which in the Mongolian samples we designated by number (our group names

169    thus have the format "<Country>_<Time Period>_<Genetic Cluster>_<Cultural

170    Association If Any>") (Supplementary Note, Table S1 and Online Table 1). We

171    merged the data with previously reported data (Online Table 4).

172

173    We carried out Principal Component Analysis (PCA) using smartpca[21], projecting the

174    ancient samples onto axes computed using present-day people. The analysis shows

175    that population structure in East Asia is correlated with geographic and linguistic

176    categories, albeit with important exceptions. Groups in Northwest China, Nepal, and

177    Siberia deviate towards West Eurasians in the PCA (Supplementary Information

178    section 2, Figure 2), reflecting multiple episodes of West Eurasian-related admixture

179    that we estimate occurred 5 to 70 generations ago based on the decay of linkage

180    disequilibrium[22] (Table S3 and Table S4). East Asians with minimal proportions of

181    West Eurasian-related ancestry fall along a gradient with three clusters at their poles.

182    The "Amur Basin Cluster" correlates geographically with ancient and present-day

183    populations living in the Amur River Basin, and linguistically with present-day

184    indigenous people speaking Tungusic languages and the Nivkh. The "Tibetan Plateau

185    Cluster" is most strongly represented in ancient Chokhopani, Mebrak, and Samzdong

186    individuals from Nepal[15] and in present-day people speaking Tibetan-Burman

187    languages and living on the Tibetan Plateau. The "Southeast Asian Cluster" is

188    maximized in ancient Taiwan groups and present-day people in Southeast Asia and

189    southern parts of China speaking Austroasiatic, Tai-Kadai and Austronesian

190    languages (Figure S1, Figure S2). Han are intermediate among these clusters, with

191    northern Han projecting close to the Neolithic Wuzhuangguoliang individuals from

192    northern China (Figure 2). We observe two genetic clusters within Mongolia: one falls

193    closer to ancient individuals from the Amur Basin Cluster ('East' based on their

194   geography), and the second clusters toward ancient individuals of the Afanasievo

195   culture ( 'West'), while a few individuals take intermediate positions between the two

196   (Supplementary Information section 2).

197

198   The three most ancient individuals of the Mongolia 'East' cluster are from the

199   Kherlen River region of eastern Mongolia (Tamsag-Bulag culture) and date to 6000-

200   4300 BCE (this places them in the Early Neolithic period, which in Northeast Asia is

201   defined by the use of pottery and not by agriculture[23]). These individuals are

202   genetically similar to previously reported Neolithic individuals from the cis-Baikal

203   region and have minimal evidence of West Eurasian-related admixture as shown in

204   PCA (Figure 2), $f_4$-statistics and *qpAdm* (Table S5, Online Table 5, labeled as

205   Mongolia_East_N). The other seven Neolithic hunter-gatherers from northern

206   Mongolia (labeled as Mongolia_North_N) can be modeled as having 5.4% ± 1.1%

207   ancestry from a source related to previously reported West Siberian Hunter-gatherers

208   (WSHG)[8] (Online Table 5), consistent with the PCA where they are part of an east-

209   west Neolithic admixture cline in Eurasia with increasing proximity to West Eurasians

210   in groups further west. Because of this ancestry complexity, we use the

211   Mongolia_East_N individuals without significant evidence of West Eurasian-related

212   admixture as reference points for modeling the East Asian-related ancestry in later

213   groups (Online Table 5). The two oldest individuals from the Mongolia 'West' cluster

214   have very different ancestry: they are from the Shatar Chuluu kurgan site associated

215   with the Afanasievo culture, with one directly dated to 3316-2918 calBCE (we quote

216   a 95% confidence interval here and in what follows whenever we mention a direct

217   date), and are indistinguishable in ancestry from previously published ancient

218   Afanasievo individuals from the Altai region of present-day Russia, who in turn are

219   similar to previously reported Yamnaya culture individuals supporting findings that

220   eastward Yamnaya migration had a major impact on people of the Afansievo

221   culture[5,8]. All the later Mongolian individuals in our time transect were mixtures of

222   Mongolian Neolithic groups and more western steppe-related sources, as reflected by

223   statistics of the form $f_3$ (X, Y; Later Mongolian Groups), which resulted in

224   significantly negative Z scores (Z<−3) when Mongolia_East_N was used as X, and

225   when Yamnaya-related Steppe populations, AfontovaGora3, WSHG, or European

226   Middle/Late Neolithic or Bronze Age populations were used as Y (Table S6).

227

228    To quantify the admixture history of the later Mongolians, we again used *qpAdm*. A

229    large number of groups could be modeled as simple two-way admixtures of

230    Mongolia_East_N as one source (in proportions of 65-100%) and WSHG as the other

231    source (in proportions of 0-35%), with negligible contribution from Yamnaya-related

232    sources as confirmed by including Russia_Afanasievo and Russia_Sintashta groups in

233    the outgroup set (Figure 3). The groups that fit this model were not only the two

234    Neolithic groups (0-5% WSHG), but also the Early Bronze Age people from the

235    Afanasievo Kurgak govi site (15%), the Ulgii group (28%), the main grouping of

236    individuals from the Middle Bronze Age Munkhkhairkhan culture (33%), Late Bronze

237    Age burials of the Ulaanzuukh type (6%), a combined group from the Center-West

238    region (27%), the Mongun Taiga type from Khukh tolgoi (35%), and people of the

239    Iron Age Slab Grave culture (9%). A striking finding in light of previous

240    archaeological and genetic data is that the male child from Kurgak govi (individual

241    I13957, skeletal code AT_629) has no evidence of Yamnaya-related ancestry despite

242    his association with Afanasievo material culture (for example, he was buried in a

243    barrow in the form of circular platform edged by vertical stone slabs, in stretched

244    position on the back on the bottom of deep rectangular pit and with a typical

245    Afanasievo egg-shaped vessel (Supplementary Note); his late Afanasievo chronology

246    is confirmed by a direct radiocarbon date of 2858-2505 BCE[24]). This is the first

247    known case of an individual buried with Afanasievo cultural traditions who is not

248    overwhelmingly Yamnaya-related, and he also shows genetic continuity with an

249    individual buried at the same site Kurgak govi 2 in a square barrow (individual I6361,

250    skeletal code AT_635, direct radiocarbon date 2618-2487 BCE). We label this second

251    individuals as having an Ulgii cultural association, although a different archaeological

252    assessment associates this individual to the Afanasievo or Chemurchek cultures[25], so

253    it is possible that this provides a second example of Afanasievo material culture being

254    adopted by individuals without any Yamnaya ancestry. The legacy of the Yamnaya-

255    era spread into Mongolia continued in two individuals from the Chemurchek culture

256    whose ancestry can be only modeled by using Afanasievo as one of the sources

257    (49.0%±2.6%, Online Table 5). This model fits even when ancient European farmers

258    are included in the outgroups, showning that if the long-distance transfer of West

259    European megalithic cultural traditions to people of the Chemurchek culture that has

260    been suggested in the archaeological literature occurred,[26] it must have been through

261    spread of ideas rather than through movement of people.

7

262

263 Beginning in the Middle Bronze Age in Mongolia, there is no compelling evidence

264 for a persistence of the Yamnaya-derivd lineages originally spread into the region

265 with Afanasievo. Instead in the Late Bronze Age and Iron Age and afterward we have

266 data from multiple Mongolian groups whose Yamnaya-related ancestry can only be

267 modeled as deriving not from the initial Afanasievo migration but instead from a later

268 eastward spread into Mongolia related to people of the Middle to Late Bronze Age

269 Sintashta and Andronovo horizons who were themselves a mixture of ~2/3 Yamnaya-

270 related and 1/3 European farmer-related ancestry[5,7,8]. The Sintashta-related ancestry is

271 detected in proportions of 5% to 57% in individuals from the

272 Mongolia_LBA_6_Khovsgol (a culturally mixed group from the literature[14]),

273 Mongolia_LBA_3_MongunTaiga, Mongolia_LBA_5_CenterWest,

274 Mongolia_EIA_4_Sagly, Mongolia_EIA_6_Pazyryk, and Mongolia_Mongol groups,

275 with the most substantial proportions of Sintashta-related ancestry always coming

276 from western Mongolia (Figure 3, Online Table 5). For all these groups, the *qpAdm*

277 ancestry models pass when Afanasievo is included in the outgroups while models

278 with Afanasievo treated as the source with Sintashta more distantly related outgroups

279 are all rejected (Figure 3, Online Table 5). Starting from the Early Iron Age, we

280 finally detect evidence of gene flow in Mongolia from groups related to Han Chinese.

281 Specifically, when Han are included in the outgroups, our models of mixtures in

282 different proportions of Mongolia_East_N, Russia_Afanasievo, Russia_Sintashta, and

283 WSHG continue to work for all Bronze Age and Neolithic groups, but fail for an

284 Early Iron Age individual from Tsengel sum (Mongolia_EIA_5), and for Xiongnu and

285 Mongols. When we include Han Chinese as a possible source, we estimate ancestry

286 proportions of 20-40% in Xiongnu and Mongols (Online Table 5).

287

288 While the Afanasievo-derived lineages are consistent with having largely disappeared

289 in Mongolia by the Late Bronze Age when our data showed that later groups with

290 Steppe pastoralist ancestry made an impact, we confirm and strengthen previous

291 ancient DNA analysis suggesting that the legacy of this expansion persisted in

292 western China into the Iron Age Shirenzigou culture (410-190 BCE)[27]. The only

293 parsimonious model for this group that fits according to our criteria is a 3-way

294 mixture of groups related to Mongolia_N_East, Russia_Afanasievo, WSHG. The only

295 other remotely plausible model (although not formally a good fit) also requires

296    Russia_Afanasievo as a source (Figure 3, Online Table 5). The findings of the original

297    study that reported evidence that the Afanasievo spread was the source of Steppe

298    ancestry in the Iron Age Shirenzigou have been questioned with the proposal of

299    alternative models that use ancient Kazakh Steppe Herders from the site of Botai,

300    Wusun, Saka and ancient Tibetans from the site of Mebrak[15] in present-day Nepal as

301    major sources for Steppe and East Asian-related ancestry[28]. However, when we fit

302    these models with Russia_Afanasievo and Mongolian_East_N added to the outgroups,

303    the proposed models are rejected (P-values between $10^{-7}$ and $10^{-2}$), except in a model

304    involving a single low coverage Saka individual from Kazakhstan as a source

305    (P=0.17, likely reflecting the limited power to reject models with this low coverage).

306    Repeating the modeling using other ancient Nepalese with very similar genetic

307    ancestry to that in Mebrak results in uniformly poor fits (Online Table 5). Thus,

308    ancestry typical of the Afanasievo culture and Mongolian Neolithic contributed to the

309    Shirenzigou individuals, supporting the theory that the Tocharian languages of the

310    Tarim Basin—from the second-oldest-known branch of the Indo-European language

311    family—spread eastward through the migration of Yamnaya steppe pastoralists to the

312    Altai Mountains and Mongolia in the guise of the Afansievo culture, from where they

313    spread further to Xinjiang[5,7,8,27,29,30]. These results are significant for theories of Indo-

314    European language diversification, as they increase the evidence in favor of the

315    hypothesis the branch time of the second-oldest branch in the Indo-European language

316    tree occurred at the end of the fourth millennium BCE[27,29,30].

317

318    The individuals from the ~5000 BCE Neolithic Boisman culture and the ~1000 BCE

319    Iron Age Yankovsky culture together with the previously published ~6000 BCE data

320    from Devil's Gate cave[19] are genetically very similar, documenting a continuous

321    presence of this ancestry profile in the Amur River Basin stretching back at least to

322    eight thousand years ago (Figure 2 and Figure S2). The genetic continuity is also

323    evident in the prevailing Y chromosomal haplogroup C2b-F1396 and mitochondrial

324    haplogroups D4 and C5 of the Boisman individuals, which are predominant lineages

325    in present-day Tungusic, Mongolic, and some Turkic-speakers. The Neolithic

326    Boisman individuals shared an affinity with Jomon as suggested by their intermediate

327    positions between Mongolia_East_N and Jomon in the PCA and confirmed by the

328    significantly positive statistic $f_4$ (Mongolia_East_N, Boisman; Mbuti, Jomon).

329    Statistics such as $f_4$ (Native American, Mbuti; Test East Asian,

330     Boisman/Mongolia_East_N) show that Native Americans share more alleles with

331     Boisman and Mongolia_East_N than they do with the great majority of other East

332     Asians in our dataset (Table S5). It is unlikely that these statistics are explained by

333     back-flow from Native Americans since Boisman and other East Asians share alleles

334     at an equal rate with the ~24,000-year-old Ancient North Eurasian MA1 who was

335     from a population that contributed about 1/3 of all Native American ancestry[31]. A

336     plausible explanation for this observation is that the Boisman/Mongolia Neolithic

337     ancestry was linked (deeply) to the source of the East Asian-related ancestry in Native

338     Americans[3,31]. We can also model published data from Neolithic and Early Bronze

339     Age individuals around Lake Baikal[7] as sharing substantial ancestry (77-94%) with

340     the lineage represented by Mongolia_East_N, revealing that this type of ancestry was

341     once spread over a wide region spanning across Lake Baikal, eastern Mongolia, and

342     the Amur River Basin (Table S7). Some present-day populations around the Amur

343     River Basin harbor large fractions of ancestry consistent with deriving from more

344     southern East Asian populations related to Han Chinese (but not necessarily Han

345     themselves) in proportions of 13-50%. We can show that this admixture occurred at

346     least by the Early Medieval period because one Heishui_Mohe individual (I3358,

347     directly dated to 1050-1220 CE) is estimated to have harbored more than 50%

348     ancestry from Han or related groups (Table S8).

349

350     The Tibetan Plateau, with an average elevation of more than 4,000 meters, is one of

351     the most extreme environments in which humans live. Archaeological evidence

352     suggests two main phases for modern human peopling of the Tibetan Plateau. The

353     first can be traced back to at least ~160,000 years ago probably by Denisovans[32] and

354     then to 40,000-30,000 years ago as reflected in abundant blade tool assemblages[33].

355     However, it is only in the last ~3,600 years that there is evidence for continuous

356     permanent occupation of this region with the advent of agriculture[34]. We grouped 17

357     present-day populations from the highlands into three categories based on genetic

358     clustering patterns (Figure S3): "Core Tibetans" who are closely related to the ancient

359     Nepal individuals such as Chokopani with a minimal amount of admixture with

360     groups related to West Eurasians and lowland East Asians in the last dozens of

361     generations, "northern Tibetans" who are admixed between lineages related to Core

362     Tibetans and West Eurasians, and "Tibeto-Yi Corridor" populations (the eastern edge

363     of the Tibetan Plateau connecting the highlands to the lowlands) that includes not just

10

Tibetan speakers but also Qiang and Lolo-Burmese speakers who we estimate using *qpAdm*[4,35] have 30-70% Southeast Asian Cluster-related ancestry (Table S9). We computed $f_3$ (Mbuti; Core Tibetan, non-Tibetan East Asian) to search for non-Tibetans that share the most genetic drift with Tibetans. Neolithic Wuzhuangguoliang, Han and Qiang appear at the top of the list (Table S10), suggesting that Tibetans harbor ancestry from a population closely related to Wuzhuangguoliang that also contributed more to Qiang and Han than to other present-day East Asian groups. We estimate that the mixture occurred 60-80 generations ago (2240-1680 years ago assuming 28 years per generation[36] under a model of a single pulse of admixture (Table S11). This represents an average date and so only provides a lower bound on when these two populations began to mix; the start of their period of admixture could plausibly be as old as the ~3,600-year-old date for the spread of agriculture onto the Tibetan plateau. These findings are therefore consistent with archaeological evidence that expansions of farmers from the Upper and Middle Yellow River Basin influenced populations of the Tibetan Plateau from the Neolithic to the Bronze Age as they spread across the China Central plain[37,38], and with Y chromosome evidence that the shared common haplogroup Oα-F5 between Han and Tibetans coalesced to a common ancestry less than 5,800 years ago[39].

In the south, we find that the ancient Taiwan Hanben and Gongguan culture individuals dating from at least a span of 1400 BCE - 600 CE are genetically most similar to present-day Austronesian speakers and ancient Lapita individuals from Vanuatu as shown in outgroup $f_3$-statistics and significantly positive $f_4$-statistics (Taiwan_Hanben/Gongguan, Mbuti; Ami/Atayal/Lapita, other Asians) (Table S8). The similarity to Austronesian-speakers is also evident in the Iron Age dominant paternal Y chromosome lineage O3a2c2-N6 and maternal mtDNA lineages E1a, B4a1a, F3b1, and F4b, which are widespread lineages among Austronesian-speakers[40,41]. We compared the present-day Austronesian-speaking Ami and Atayal of Taiwan with diverse Asian populations using statistics like $f_4$ (Taiwan Iron Age/Austronesian, Mbuti; Asian1, Asian2). Ancient Taiwan groups and Austronesian-speakers share significantly more alleles with Tai-Kadai speakers in southern mainland China and in Hainan Island[42] than they do with other East Asians (Table S8), consistent with the hypothesis that ancient populations related to present-day Tai-Kadai speakers are the source for the spread of agriculture to Taiwan island around

11

398    5000 years ago[43]. The Jomon share alleles at an elevated rate with ancient Taiwan

399    individuals and Ami/Atayal as measured by statistics of the form $f_4$ (Jomon, Mbuti;

400    Ancient Taiwan/Austronesian-speaker, other Asians) compared with other East Asian

401    groups, with the exception of groups in the Amur Basin Cluster (Table S8)[44].

402

403    The Han Chinese are the world's largest ethnic group. It has been hypothesized based

404    on the archaeologically documented spread of material culture and farming

405    technology, as well as the linguistic evidence of links among Sino-Tibetan languages,

406    that one of the ancestral populations of the Han might have consisted of early farmers

407    along the Upper and Middle Yellow River in northern China, some of whose

408    descendants also may have spread to the Tibetan Plateau and contributed to present-

409    day Tibeto-Burmans[45]. Archaeological and historical evidence document how during

410    the past two millennia, the Han expanded south into regions inhabited by previously

411    established agriculturalists[46]. Analysis of genome-wide variation among present-day

412    populations has revealed that the Han Chinese are characterized by a "North-South"

413    cline[47,48], which is confirmed by our analysis. The Neolithic Wuzhuangguoliang,

414    present-day Tibetans, and Amur River Basin populations, share significantly more

415    alleles with Han Chinese compared with the Southeast Asian Cluster, while the

416    Southeast Asian Cluster groups share significantly more alleles with the majority of

417    Han Chinese groups when compared with the Neolithic Wuzhuangguoliang (Table

418    S12, Table S13). These findings suggest that Han Chinese may be admixed in variable

419    proportions between groups related to Neolithic Wuzhuangguoliang and people

420    related to those of the Southeast Asian Cluster. To determine the minimum number of

421    source populations needed to explain the ancestry of the Han, we used $qpWave$[4,49] to

422    study the matrix of all possible statistics of the form $f_4$ (Han$_1$, Han$_2$; O$_1$, O$_2$), where

423    "O$_1$" and "O$_2$" are outgroups that are unlikely to have been affected by recent gene

424    flow from Han Chinese. This analysis confirms that two source populations are

425    consistent with all of the ancestry in most Han Chinese groups (with the exception of

426    some West Eurasian-related admixture that affects some northern Han Chinese in

427    proportions of 2-4% among the groups we sampled; Table S14 and Table S15).

428    Specifically, we can model almost all present-day Han Chinese as mixtures of two

429    ancestral populations, in a variety of proportions, with 77-93% related to Neolithic

430    Wuzhuangguoliang from the Yellow River basin, and the remainder from a

431    population related to ancient Taiwan that we hypothesize was closely related to the

432    rice farmers of the Yangtze River Basin. This is also consistent with our inference that

433    the Yangtze River farmer related ancestry contributed nearly all the ancestry of

434    Austronesian speakers and Tai-Kadai speakers and about 2/3 of some Austroasiatic

435    speakers[17,20] (Figure 4). A caveat is that there is a modest level of modern

436    contamination in the Wuzhuangguoliang we use as a source population for this

437    analysis (Online Table 1), but this would not bias admixture estimates by more than

438    the contamination estimate of 3-4%. The average dates of West Eurasian-related

439    admixture in northern Han Chinese populations Han_NChina and Han_Shanxi are 32-

440    45 generations ago, suggesting that mixture was continuing at the time of the Tang

441    Dynasty (618-907 CE) and Song Dynasty (960-1279 BCE) during which time there

442    are historical records of integration of Han Chinese amd western ethnic groups, but

443    this date is an average so the mixture between groups could have begun earlier.

444

445    To obtain insight into the formation of present-day Japanese archipelago populations,

446    we searched for groups that contribute most strongly to present-day Japanese through

447    admixture $f_3$-statistics. The most strongly negative signals come from mixtures of Han

448    Chinese and ancient Jomon ($f_3$(Japanese; Han Chinese, Jomon)) (Table S16). We can

449    model present-day Japanese as two-way mixtures of 84.3% Han Chinese and 15.7%

450    Jomon or 87.6% Korean and 12.4% Jomon (we cannot distinguish statistically

451    between these two sources; Table S17 and Table S18). This analysis by no means

452    suggests that the mainland ancestry in Japan was contributed directly by the Han

453    Chinese or Koreans themselves, but does suggest that it is from an ancestral

454    population related to those that contributed in large proportion to Han Chinese as well

455    as to Koreans for which we do not yet have ancient DNA data.

456

457    We used *qpGraph*[35] to explore models with population splits and gene flow, and

458    tested their fit to the data by computing $f_2$-, $f_3$- and $f_4$- statistics measuring allele

459    sharing among pairs, triples, and quadruples of populations, evaluating fit based on

460    the maximum |Z|-score comparing predicted and observed values. We further

461    constrained the models by using estimates of the relative population split times

462    between the selected pairs of populations based on the output of the MSMC

463    software[50]. While admixture graph modeling based on allele frequency correlation

464    statistics is not able to reject a model in which ancient Taiwan individuals and

465    Boisman share substantial ancestry with each other more recently than either does

466    with the ancestors of Chokopani and Core Tibetans, this model cannot be correct

467    because our MSMC analysis reveals that Core Tibetans (closely related to Chokopani)

468    and Ulchi (closely related to Boisman) share ancestry more recently in time on

469    average than either does with Ami (related to Taiwan_Hanben). This MSMC-based

470    constraint allowed us to identify a parsimonious working model for the deep history

471    of key lineages discussed in this study (Supplementary Information section 3:

472    *qpGraph* Modeling). Our fitted model (Figure 5), suggests that much of East Asian

473    ancestry today can be modelled as derived from two ancient populations: one from the

474    same lineage as the approximately ~40,000-year-old Tianyuan individual and the

475    other more closely related to Onge, with groups today having variable proportions of

476    ancestry from these two deep sources. In this model, the Mongolia_East_N and Amur

477    River Basin Boisman related lineages derive the largest proportion of their ancestry

478    from the Tianyuan-related lineage and the least proportion of ancestry from the Onge-

479    related lineage compared with other East Asians. A sister lineage of

480    Mongolia_East_N is consistent with expanding into the Tibetan Plateau and mixing

481    with the local hunter-gatherers who represent an Onge-related branch in the tree. The

482    Taiwan Hanben are well modelled as deriving about 14% of their ancestry from a

483    lineage remotely related to Onge and the rest of their ancestry from a lineage that also

484    contributed to Jomon and Boisman on the Tianyuan side, a scenario that would

485    explain the observed affinity among Jomon, Boisman and Taiwan Hanben. We

486    estimate that Jomon individuals derived 45% of their ancestry from a deep basal

487    lineage on the Onge side. These results are consistent with the scenario a Late

488    Pleistocene coastal route of human migration linking Southeast Asia, the Japanese

489    Archipelago and the Russian Far East[51]. Due to the paucity of ancient genomic data

490    from Upper Paleolithic East Asians, there are limited constraints at present for

491    reconstructing the deep branching patterns of East Asian ancestral populations, and it

492    is certain that this admixture graph is an oversimplification and that additional

493    features of deep population relationships will be revealed through future work.

494

495    At the end of the last Ice Age, there were multiple highly differentiated populations in

496    East as well as West Eurasia, and it is now clear that these groups mixed in both

497    regions, instead of one population displacing the others. In West Eurasia, there were

498    at least four divergent populations each as genetically differentiated from each other

499    as Europeans and East Asians today (average $F_{ST}=0.10$), which mixed in the

500     Neolithic, reducing heterogeneity (average $F_{ST}$=0.03) and mixed further in the Bronze

501     Age and Iron Age to produce the present-relatively low differentiation that

502     characterizes modern West Eurasia (average $F_{ST}$=0.01)[52]. In East Eurasia, our study

503     suggests an analogous process, with the differentiation characteristic of the Amur

504     River Basin groups, Neolithic Yellow River farmers, and people related to those of

505     the Taiwan Iron Age (average $F_{ST}$=0.06 in our data) collapsing through mixture to

506     today's relatively low differentiation (average $F_{ST}$=0.01-0.02) (Figure 6). A priority

507     should be to obtain ancient DNA data for the hypothesized Yangtze River population

508     (the putative source for the ancestry prevalent in the Southeast Asian Cluster of

509     present-day groups), which should, in turn, make it possible to test and further extend

510     these models, and in particular to understand if dispersals of people in Southeast Asia

511     do or do not correlate to ancient movements of people.

15

512  **References**

513  1.  Cavalli-Sforza, L. L. The Chinese human genome diversity project. *Proc. Natl. Acad. Sci.*
514      *USA* **95**, 11501-11503 (1998).

515  2.  HUGO Pan-Asian SNP Consortium. Mapping human genetic diversity in Asia. *Science* **326**,
516      1541-1545 (2009).

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

519  4.  Haak, W., et al. Massive migration from the steppe was a source for Indo-European languages
520      in Europe. *Nature* **522**, 207–211 (2015).

521  5.  Allentoft, M.E., et al.. Population genomics of Bronze Age Eurasia. *Nature* **522**,167-172
522      (2015).

523  6.  Fu, Q., et al.. The genetic history of ice age Europe. *Nature* **534**, 200-205 (2016).

524  7.  de Barros Damgaard, P., et al.. 137 ancient human genomes from across the Eurasian steppes.
525      *Nature* **557**, 369-374 (2018).

526  8.  Narasimhan, V.M., et al. The formation of human populations in South and Central Asia.
527      *Science* **365**, eaat7487 (2019).

528  9.  Fu, Q., et al. An early modern human from Romania with a recent Neanderthal ancestor.
529      *Nature* **524**, 216–219 (2015).

530  10. Fu, Q.,et al. DNA analysis of an early modern human from Tianyuan Cave, China. *Proc. Natl*
531      *Acad. Sci. USA* **110**, 2223–2227 (2013).

532  11. Korneliussen, T. S., Albrechtsen, A., & Nielsen, R. ANGSD: Analysis of Next Generation
533      Sequencing Data. *BMC Bioinformatics* **15**, 356 (2014).

534  12. Posth, C., et al. Language continuity despite population replacement in Remote Oceania. *Nat*
535      *Ecol Evol*. **2**, 731-740 (2018).

536  13. Sikora, M., et al. The population history of northeastern Siberia since the Pleistocene. *Nature*
537      **570**, 182-188 (2019).

538  14. Jeong, C., et al. The genetic history of admixture across inner Eurasia. *Nat. Ecol. Evol*. **3**,
539      966–976 (2019).

540  15. Jeong, C., et al. Long-term genetic stability and a high-altitude East Asian origin for the
541      peoples of the high valleys of the Himalayan arc. *Proc. Natl. Acad. Sci. USA* **113**, 7485–7490
542      (2016).

543  16. Skoglund, P., et al. Genomic insights into the peopling of the Southwest Pacific. *Nature* **538**,
544      510-513 (2016).

545  17. Lipson, M., et al. Ancient genomes document multiple waves of migration in Southeast Asian
546      prehistory. *Science* **361**, 92-95 (2018).

547  18. Kanzawa-Kiriyama, H., et al. A partial nuclear genome of the Jomons who lived 3000 years
548      ago in Fukushima, Japan. *J. Hum. Genet* **62**, 213–221 (2016).

16

549  19. Siska, V., et al. Genome-wide data from two early Neolithic East Asian individuals dating to
550      7700 years ago. *Sci Adv*. **3**, e1601877 (2017).

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

552  21. Patterson, N., Price, A. L., & Reich, D. Population structure and eigenanalysis. *PLoS Genet*. **2**,
553      e190 (2006).

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

556  23. Jordan, P. & Zvelebil M. ed. *Ceramics Before Farming: The Dispersal of Pottery Among*
557      *Prehistoric Eurasian Hunter-Gatherers* (Routledge, New York, 2010).

558  24. Kovalev, A. A., & Erdenebaatar, D. Discovery of New Cultures of the Bronze Age in
559      Mongolia according to the Data obtained by the International Central Asian Archaeological
560      Expedition. In *Current Archaeological Research in Mongolia*, (eds Bemmann, J., H.
561      Parzinger, H., Pohl, E., D. Tseveendorzh, D.) 149–170 (Bonn: Vor- und Frügeschichtliche
562      Archäologie Rheinische Friedrich-Wilhelm-Universität Bonn, 2009).

563  25. Wilkins, S., et al. Dairy pastoralism sustained eastern Eurasian steppe populations for 5,000
564      years. *Nat Ecol Evol* **4**, 346–355 (2020).

565  26. Kovalev, A. The Great Migration of the Chemurchek People from France to the Altai in the
566      Early 3rd Millennium BCE . *International Journal of Eurasian Studies*. 1(11) , pp. 1-58
567      (2011).

568  27. Ning, C., et al. Ancient Genomes Reveal Yamnaya-Related Ancestry and a Potential Source of
569      Indo-European Speakers in Iron Age Tianshan. *Curr Biol*. **29**, 2526-2532.e4 (2019).

570  28. Weslowski. Eurogenes Blog: A surprising twist to the Shirenzigou nomads story
571      (2019).https://eurogenes.blogspot.com/2019/08/a-surprising-twist-to-shirenzigou.html.

572  29. Mallory, J.P. *In Search of the INDO-Europeans: Language, Archaeology and Myth* (Thames &
573      Hudson, New York, 1991).

574  30. Anthony, D. *The Horse, the Wheel, and Language: How Bronze-Age Riders from the Eurasian*
575      *Steppes Shaped the Modern World* (Princeton University Press, Princeton and Oxford, 2007).

576  31. Raghavan, M., et al. Upper Palaeolithic Siberian genome reveals dual ancestry of Native
577      Americans. *Nature* **505**, 87-91 (2014).

578  32. Chen, F., et al. A late Middle Pleistocene Denisovan mandible from the Tibetan Plateau.
579      *Nature* **569**, 409-412 (2019).

580  33. Zhang, X. L., et al. The earliest human occupation of the high-altitude Tibetan Plateau 40
581      thousand to 30 thousand years ago. *Science* **362**, 1049-1051 (2018).

582  34. Chen, F.H., et al. Agriculture facilitated permanent human occupation of the Tibetan Plateau
583      after 3600 B.P. *Science* **347**, 248-250 (2015).

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

585   36. Moorjani, M., Sankararaman, S., Fu, Q., Przeworski, M., Patterson, N., & Reich, D. A genetic
586        method for dating ancient genomes provides a direct estimate of human generation interval in
587        the last 45,000 years. *Proc. Natl. Acad. Sci. USA* **113**, 5652-5657 (2016).
588   37. Barton, L., et al. Agricultural origins and the isotopic identity of domestication in northern
589        China. *Proc. Natl. Acad. Sci. USA* **106**, 5523-5528 (2009).
590   38. Yang, X., et al. Early millet use in northern China. *Proc. Natl. Acad. Sci. USA* **109**, 3726–3730
591        (2012).
592   39. Wang, L.X., et al. Reconstruction of Y-chromosome phylogeny reveals two neolithic
593        expansions of Tibeto-Burman populations. *Mol Genet Genomics*. **293**, 1293-1300 (2018).
594   40. Wei, L.H., et al. Phylogeography of Y-chromosome haplogroup O3a2b2-N6 reveals patrilineal
595        traces of Austronesian populations on the eastern coastal regions of Asia. *PLoS One* **12**,
596        e0175080 (2017).
597   41. Ko, A.M., et al. Early Austronesians: into and out of Taiwan. *Am. J. Hum. Genet*. **94**, 426-36
598        (2014).
599   42. Lipson, M., et al. Reconstructing Austronesian population history in island Southeast Asia.
600        *Nat Commun*. **5**, 4689 (2014).
601   43. Bellwood, P. The checkered prehistory of rice movement southwards as a domesticated
602        cereal—from the Yangzi to the equator. *Rice* **4**, 93-103 (2011).
603   44. Kanzawa-Kiriyama H., et al. Late Jomon male and female genome sequences from the
604        Funadomari site in Hokkaido, Japan. *Anthropol Sci*. **127**, 83-108 (2019).
605   45. Su, B., et al. Y chromosome haplotypes reveal prehistorical migrations to the Himalayas.
606        *Hum. Genet*. **107**, 582-590 (2000).
607   46. Ge, J. X., Wu, S. D., & Chao, S. J. *Zhongguo yimin shi (The Migration History of China)*
608        (Fujian People's Publishing House, Fuzhou, 1997).
609   47. Chen, J., et al. Genetic structure of the Han Chinese population revealed by genome-wide
610        SNP variation. *Am. J. Hum. Genet*. **85**, 775-785 (2009).
611   48. Xu, S., et al. Genomic dissection of population substructure of Han Chinese and its
612        implication in association studies. *Am. J. Hum. Genet*. **85**, 762-774 (2009).
613   49. Reich, D., et al. Reconstructing Native American population history. *Nature* **488**, 370-374
614        (2012).
615   50. Schiffels, S., & Durbin, R. Inferring human population size and separation history from
616        multiple genome sequences. *Nat Genet*. **46**, 919-925 (2014).
617   51. Matsumura, H., et al. Craniometrics Reveal "Two Layers" of Prehistoric Human Dispersal in
618        Eastern Eurasia. *Sci Rep.* **9**, 1451 (2019).
619   52. Lazaridis, I., et al. Genomic insights into the origin of farming in the ancient Near East.
620        *Nature* **536**, 419–424 (2016).
621
622

623 **Methods**

624 **Ancient DNA laboratory work**

625 All samples except those from Wuzhuangguoliang were prepared in dedicated clean

626 room facilities at Harvard Medical School, Boston, USA. Online Table 2 lists

627 experimental settings for each sample and library included in the dataset. Skeletal

628 samples were surface cleaned and drilled or sandblasted and milled to produce a fine

629 powder for DNA extraction[53,54]. We then either followed the extraction protocol by

630 Dabney et al[55] replacing the extender-MinElute-column assembly with the columns

631 from the Roche High Pure Viral Nucleic Acid Large Volume Kit[56] (manual

632 extraction) or, for samples prepared later, used DNA extraction protocol based on

633 silica beads instead of spin columns (and Dabney buffer) to allow for automated DNA

634 purification[57] (robotic extraction). We prepared individually barcoded double-

635 stranded libraries for most samples using a protocol that included a DNA repair step

636 with Uracil-DNA-glycosylase (UDG) treatment to cut molecules at locations

637 containing ancient DNA damage that is inefficient at the terminal positions of DNA

638 molecules (Online Table 1, UDG: "half")[58], or, without UDG pre-treatment (double

639 stranded minus). For a few samples processed later, single stranded DNA libraries[59]

640 were prepared with USER (NEB) addition in the dephosphorylation step that results

641 in inefficient uracil removal at the 5'end of the DNA molecules, and does not affect

642 deamination rates at the terminal 3' end[60]. We performed target enrichment via

643 hybridization of these libraries with previously reported protocols[10]. We either

644 enriched for the mitochondrial genome and 1.2M SNPs in two separate experiments

645 or together in a single experiment. If split over two experiments, the first enrichment

646 was for sequences aligning to mitochondrial DNA[58,61] with some baits overlapping

647 nuclear targets spiked in to screen libraries for nuclear DNA content. The second in-

648 solution enrichment was for a targeted set of 1,237,207 SNPs that comprises a merge

649 of two previously reported sets of 394,577 SNPs (390k capture)[4] and 842,630 SNPs[9].

650 We sequenced the enriched libraries on an Illumina NextSeq500 instrument for 2x76

651 cycles (and both indices) or on Hiseq X10 instruments at the Broad Institute of MIT

652 and Harvard for 2x101 cycles. We also shotgun sequenced each library for a few

653 hundred thousand reads to assess the fraction of human reads.

654

655 Ancient DNA extractions of the Wuzhuangguoliang samples were performed in the

656 clean room at Xi'an Jiaotong University and Xiamen University following the

19

657 protocol by Rohland and Hofreiter[62]. Each sample extract was converted into double-

658 stranded Illumina libraries following the manufacturer's protocol (Fast Library Prep

659 Kit, iGeneTech, Beijing, China). Sample-specific indexing barcodes were added to

660 both sides of the fragments via amplification. Nuclear DNA capture was performed

661 with AIExome Enrichment Kit V1 (iGeneTech, Beijing, China) according to the

662 manufacturer's protocol and sequenced on an Illumina NovaSeq instrument with 150

663 base pair paired-end reads. Sequences that did not perfectly match one of the expected

664 index combinations were discarded.

665

666 For the AH1-7 and AH1-17 DNA extracts, we prepared whole genome sequencing

667 libraries. The two DNA extracts were converted into barcoded Illumina sequencing

668 libraries using commercially available library kits (NEBNext® Ultra™ II DNA

669 Library Prep Kit) and Illumina-specific primers[63]. DNA libraries were not treated with

670 uracil-DNA-glycosylase (UDG) [59]. We used a MinElute Gel Extraction Kit (Qiagen,

671 Hilden, Germany) for purification. Two libraries were sequenced on a HiSeqX10

672 instrument (2×150 bp, PE) at the Novogene Sequencing Centre (Beijing, China). The

673 base calling was performed using CASAVA software.

674

675 **Bioinformatic processing**

676 For the sequencing data produced at Harvard Medical School, we used one of two

677 pipelines ("pipeline 1" or "pipeline 2"; Online Table 2). An up-to-date description of

678 both pipelines and analyses showing that the differences between them do not cause

679 systematic bias in population genetic analysis can be found in Fernandes et al[64]. For

680 both pipelines we began by de-multiplexed the data and assigning sequences to

681 samples based on the barcodes and/or indices, allowing up to one mismatch per

682 barcode or index. We trimmed adapters and restricted to fragments where the two

683 reads overlapped by at least 15 nucleotides. In pipeline 1 we merged the sequences

684 (allowing up to one mismatch) using a modified version of *Seqprep*[65] where bases in

685 the merged region are chosen based on highest quality in case of a conflict, and in

686 pipeline 2 we used custom software (https://github.com/DReichLab/ADNA-Tools).

687 For mitochondrial DNA analysis, we aligned the resulting merged sequences to the

688 RSRS reference genome[66] using *bwa* (version 0.6.1 for pipeline 1 and version 0.7.15

689 for pipeline 2)[67], and removed duplicates with the same orientation, start and stop

690 positions, and molecular barcodes. We determined mitochondrial DNA haplogroups

20

691    using *HaploGrep2*[68]. We also analyzed the sequences to generate two assessments of

692    ancient DNA authenticity. The first assessment estimated the rate of cytosine to

693    thymine substitution in the final nucleotide, which is expected to be at least 3% at

694    cytosines in libraries prepared with a partial UDG treatment protocol and at least 10%

695    for untreated libraries (minus) and single stranded libraries; all libraries we analyzed

696    met this threshold. The second assessment used *contamMix* (version 1.0.9 for pipeline

697    1 and 1.0.12 for pipeline 2)[10] to determine the fraction of mtDNA sequences in an

698    ancient sample that match the endogenous majority consensus more closely than a

699    comparison set of 311 worldwide present-day human mtDNAs (Online Table 1).

700    Computational processing of the sequence data from the whole genome was the same

701    as the mtDNA enrichment except that the human genome (hg19) was used as the

702    target reference. Due to the low coverage, diploid calling was not possible; instead,

703    we randomly selected a single sequence covering every SNP position of interest

704    ("pseudo-haploid" data) using custom software, only using nucleotides that were a

705    minimum distance from the ends of the sequences to avoid deamination artifacts

706    (https://github.com/DReichLab/adna-workflow). The coverages and numbers of SNPs

707    covered at least once on the autosomes (chromosomes 1-22) are in Online Table 1.

708

709    For the sequencing data from the Wuzhuangguoliang samples, we clipped adaptors

710    with *leehom*[69] and then further processed using *EAGER*[70], including mapping with

711    *bwa* (v0.6.1)[67] against the human genome reference GRCh37/hg19 (or just the

712    mitochondrial reference sequence), and removing duplicate reads with the same

713    orientation and start and end positions. To avoid an excess of remaining C-to-T and

714    G-to-A transitions at the ends of the sequences, we clipped three bases of the ends of

715    each read for each sample using trimBam

716    (https://genome.sph.umich.edu/wiki/BamUtil:_trimBam). We generated pseudo-

717    haploid calls by selecting a single read randomly for each individual using

718    pileupCaller (https://github.com/stschiff/sequenceTools/tree/master/srcpileupCaller).

719

720    **Accelerator Mass Spectrometry Radiocarbon Dating**

721    We generated 94 direct AMS (Accelerator Mass Spectrometry) radiocarbon ($^{14}$C)

722    dates as part of this study; 87 at Pennsylvania State University (PSU) and 7 at Poznan

723    Radiocarbon Laboratory. The methods used at both laboratories are published, and

724    here we summarize the methods from PSU. Bone collagen from petrous, phalanx, or

21

725    tooth (dentine) samples was extracted and purified using a modified Longin method

726    with ultrafiltration (>30kDa gelatin)[71]. If bone collagen was poorly preserved or

727    contaminated we hydrolyzed the collagen and purified the amino acids using solid

728    phase extraction columns (XAD amino acids)[72]. Prior to extraction we sequentially

729    sonicated all samples in ACS grade methanol, acetone, and dichloromethane (30

730    minutes each) at room temperature to remove conservants or adhesives possibly used

731    during curation. Extracted collagen or amino acid preservation was evaluated using

732    crude gelatin yields (% wt), %C, %N and C/N ratios. Stable carbon and nitrogen

733    isotopes were measured on a Thermo DeltaPlus instrument with a Costech elemental

734    analyzer at Yale University. C/N ratios between 3.14 and 3.45 indicate that all

735    radiocarbon dated samples are well preserved. All samples were combusted and

736    graphitized at PSU using methods described in Kennett et al. 2017[71]. $^{14}$C

737    measurements were made on a modified National Electronics Corporation 1.5SDH-1

738    compact accelerator mass spectrometer at either the PSUAMS facility or the Keck-

739    Carbon Cycle AMS Facility. All dates were calibrated using the IntCal13 curve[73] in

740    OxCal v 4.3.2[74] and are presented in calendar years BCE/CE .

741

742    **Y chromosomal haplogroup analysis**

743    We performed Y-haplogroup determination by examining the state of SNPs present in

744    ISOGG version 11.89 (accessed March 31, 2016) and our unpublished updated

745    phylogeny.

746

747    **X-chromosome contamination estimates**

748    We performed an X-chromosomal contamination test for the male individuals

749    following an approach introduced by Rasmussen et al[75] and implemented in the

750    *ANGSD* software suite[11]. We used the "MoM" (Methods of Moments) estimates. The

751    estimates for some males are not informative because of the limited number of X-

752    chromosomal SNPs covered by at least two sequences, and hence we only report

753    results for individuals with at least 200 SNPs covered at least twice. The estimated

754    contamination rates for the male samples are low (Online Table 1). The contamination

755    rates for all samples are quite low except those from Wuzhuangguoliang. We detected

756    3-6% contamination in the Wuzhuangguoliang samples, and restricted population

757    genetic modeling analysis only to three males with 3-4% contamination.

758

759 **Data merging**

760 We merged the data with previously published datasets genotyped on Affymetrix

761 Human Origins arrays[3,35], restricting to individuals with >95% genotyping

762 completeness. We manually curated the data using ADMIXTURE[76] and

763 EIGENSOFT[21] to identify samples that were outliers compared with other samples

764 from their own populations. We removed seven individuals from subsequent analysis;

765 the population IDs for these individuals are prefixed by the string "Ignore_" in the

766 dataset we release, so users who wish to analyze these samples are still able to do so.

767

768 **Principal Components Analysis.** We carried out principal components analysis in

769 the *smartpca* program of EIGENSOFT[21], using default parameters and the lsqproject:

770 YES and numoutlieriter: 0 options.

771

772 **ADMIXTURE Analysis.** We carried out ADMIXTURE analysis in unsupervised

773 mode[76] after pruning for linkage disequilibrium in PLINK[77] with parameters --indep-

774 pairwise 200 25 0.4 which retained 256,427 SNPs for Human Origin Dataset. We ran

775 ADMIXTURE with default 5-fold cross-validation (--cv=5), varying the number of

776 ancestral populations between K=2 and K=18 in 100 bootstraps with different random

777 seeds.

778

779 *f*-statistics. We computed $f_3$-statistics and $f_4$-statistics using ADMIXTOOLS[35] with

780 default parameters. We computed standard errors using a block jackknife[78].

781

782 **$F_{ST}$ computation.** We estimated $F_{ST}$ using EIGENSOFT[21] with default parameters,

783 inbreed: YES, and fstonly: YES. We found that the inbreeding corrected and

784 uncorrected $F_{ST}$ were nearly identical (within ~0.001), and in this study, always

785 analyzed uncorrected $F_{ST}$.

786

787 **Admixture graph modeling.** Admixture graph modeling was carried out with the

788 *qpGraph* software as implemented in ADMIXTOOLS[35] using Mbuti as an outgroup.

789

790 **Testing for the number of streams of ancestry.** We used *qpWave*[4,35] as

791 implemented in ADMIXTOOLS to test whether a set of test populations is consistent

792 with being related via *N* streams of ancestry from a set of outgroup populations.

23

793

794 **Inferring mixture proportions without an explicit phylogeny.** We used *qpAdm*[4] as

795 implemented in ADMIXTOOLS to estimate mixture proportions for a *Test* population

796 as a combination of *N* 'reference' populations by exploiting (but not explicitly

797 modeling) shared genetic drift with a set of 'Outgroup' populations.

798

799 **Weighted linkage disequilibrium (LD) analysis**. LD decay was calculated using

800 ALDER[22] to infer admixture parameters including dates and mixture proportions.

801

802 **MSMC.** We used MSMC[50] following the procedures in Mallick et al[79] to infer cross-

803 coalescence rates and population sizes among Ami/Atayal, Tibetan, and Ulchi.

804

805 **Kinship analysis**. We used READ software[80] as well as a custom method[81] to

806 determine genetic kinship between individual pairs.

807

808 **Data availability**

809 The aligned sequences are available through the European Nucleotide Archive under

810 accession number [to be made available on publication]. Genotype data used in

811 analysis are available at https://reich.hms.harvard.edu/datasets. Any other relevant

812 data are available from the corresponding author upon reasonable request.

813

814 53. Pinhasi, R., Fernandes, D.M., Sirak, K., & Cheronet, O. Isolating the human cochlea to

815     generate bone powder for ancient DNA analysis. *Nat Protoc*. **14**, 1194-1205 (2019).

816 54. Sirak, K.A., et al., A minimally-invasive method for sampling human petrous bones from the

817     cranial base for ancient DNA analysis. *Biotechniques*. **62**, 283-289 (2017).

818 55. Dabney, J., et al. Complete mitochondrial genome sequence of a Middle Pleistocene cave bear

819     reconstructed from ultrashort DNA fragments. *Proc Natl Acad Sci U S A*. **110**, 15758-63

820     (2013).

821 56. Korlević, P. Reducing microbial and human contamination in DNA extractions from ancient

822     bones and teeth. *Biotechniques*. **59**, 87-93 (2015).

823 57. Rohland, N., Glocke, I., Aximu-Petri, A., & Meyer, M. Extraction of highly degraded DNA

824     from ancient bones, teeth and sediments for high-throughput sequencing. *Nat Protoc*. **13**,

825     2447-2461 (2018).

826    58. Rohland, N., Harney, E., Mallick, S., Nordenfelt, S. & Reich, D. Partial uracil–DNA–
827        glycosylase treatment for screening of ancient DNA. *Phil. Trans. R. Soc. Lond. B* **370**,
828        20130624 (2015).

829    59. Gansauge, M.T., & Meyer, M. Selective enrichment of damaged DNA molecules for ancient
830        genome sequencing. *Genome Res*. **24**, 1543-1549 (2014).

831    60. Meyer, M., et al., A high-coverage genome sequence from an archaic Denisovan individual.
832        *Science*. **338**, 222-226 (2012).

833    61. Maricic, T., Whitten, M., & Pääbo, S. Multiplexed DNA sequence capture of mitochondrial
834        genomes using PCR products. *PLoS One* **5**, e14004 (2010).

835    62. Rohland, N., & Hofreiter, M. Ancient DNA extraction from bones and teeth. *Nat. Protoc*. **2**,
836        1756–1762 (2007).

837    63. Meyer, M., & Kircher, M. Illumina sequencing library preparation for highly multiplexed
838        target capture and sequencing. *Cold Spring Harb. Protoc*. **6**, pdb.prot5448 (2010).

839    64. Fernandes, D.M., et al. The spread of steppe and Iranian-related ancestry in the islands of the
840        western Mediterranean. *Nat Ecol Evol*. 4, 334-345 (2020).

841    65. John, J. S. SeqPrep, https://github.com/jstjohn/SeqPrep (2011).

842    66. Behar, D.M., et al. A "Copernican" reassessment of the human mitochondrial DNA tree from
843        its root. *Am J Hum Genet.* 90, 675-84 (2012). Erratum in: *Am J Hum Genet.* 90, 936 (2012).

844    67. Li, H., & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform.
845        *Bioinformatics*, **25**, 1754-1760 (2009).

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

848    69. Renaud, G., Stenzel, U., & Kelso, J. leeHom: adaptor trimming and merging for Illumina
849        sequencing reads. *Nucleic Acids Res*. **42**, e141 (2014).

850    70. Peltzer, A., et al. EAGER: efficient ancient genome reconstruction. *Genome Biol*. **17**, 60
851        (2016).

852    71. Kennett, D. J. et al. Archaeogenomic evidence reveals prehistoric matrilineal dynasty. *Nat.
853        Commun.* 8, 14115 (2017).

854    72. Lohse, J. C., Madsen, D. B., Culleton, B. J. & Kennett, D. J. Isotope paleoecology of episodic
855        mid-to-late Holocene bison population expansions in the southern Plains, U.S.A. *Quat. Sci.
856        Rev.* 102, 14–26 (2014).

857    73. Reimer, P. J. et al. IntCal13 and Marine13 radiocarbon age calibration curves 0–50,000 years
858        cal BP. *Radiocarbon* 55, 1869–1887 (2013).

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

860    75. Rasmussen, M., et al. An Aboriginal Australian Genome Reveals Separate Human Dispersals
861        into Asia. *Science* **334**, 94–98 (2011).

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

864  77. Chang, C., et al. Second-generation PLINK: rising to the challenge of larger and richer
865      datasets. *GigaScience* **4**, 7 (2015).
866  78. Busing, F. T. A., Meijer, E., & Leeden, R. Delete-m Jackknife for Unequal m. *Statistics and*
867      *Computing* **9**, 3-8 (1999).
868  79. Mallick, S.M., et al. The Simons Genome Diversity Project: 300 genomes from 142 diverse
869      populations, *Nature* **538**, 201-206 (2016).
870  80. Monroy, K.J.M., Jakobsson, M., & Günther, T. Estimating genetic kin relationships in
871      prehistoric populations. *PLoS One* **13**, e0195491 (2018).
872  81. Kennett, D.J., et al. Archaeogenomic evidence reveals prehistoric matrilineal dynasty. *Nat.*
873      *Commun*. **8**, 14115 (2017).
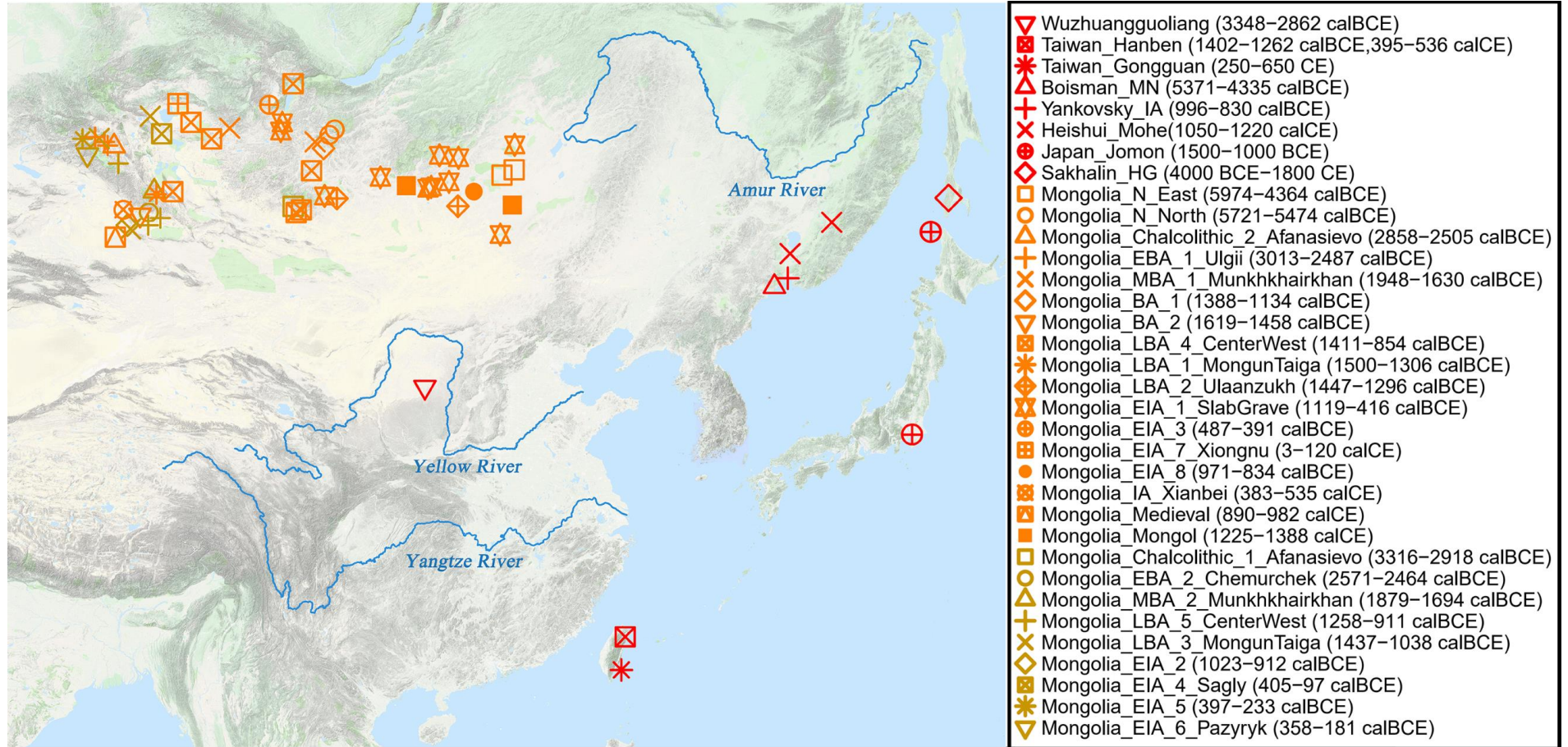
874

## Acknowledgements

906

## Author Contributions

908    Conceptualization, C.-C.W., H.-Y.Y., A.N.P., H.M., A.M.K., L.J., H.L., J.K., R.P.,

909    and D.R.; Formal Analysis, C.-C.W., R.B., M.Ma., S.M., Z.Z., B.J.C, and D.R.;

910    Investigation, C.-C.W., K.Si., O.C., A.K., N.R., A.M.K., M.Ma., S.M., K.W., N.A.,

911    N.B., K.C., B.J.C, L.E., A.M.L., M.Mi., J.O., K.S., S.W., S.Y., F.Z., J.G., Q.D., L.K.,

912    Da.L, Do.L, R.L., W.C., R.S., L.-X.W., L.W., G.X., H.Y., M.Z., G.H., X.Y., R.H.,

913    S.S., D.J.K., L.J., H.L., J.K., R.P., and D.R.; Resources, H.-Y.Y., A.N.P., R.B., D.T.,

914    J.Z., Y.-C.L, J.-Y.L., M.Ma., S.M., Z.Z., R.C., C.-J. H., C.-C.S., Y.G.N., A.V.T.,

915    A.A.T., S.L., Z.-Y.S., X.-M.W., T.-L.Y., X.H., L.C., H.D., J.B., E.Mi., D.E., T.-O.I.,

916    E.My., H.K.-K., M.N., K.Sh., D.J.K., R.P., and D.R.; Data Curation, C.-C.W., K.Si.,

917    O.C., A.K., N.R., R.B., M.Ma., S.M., B.J.C, L.E., A.A.T., and D.R.; Writing, C.-

918    C.W., H.-Y.Y., A.N.P., H.M., A.K., and D.R.; Supervision, C.-C.W., H.-Q.Z., N.R.,

919    M.R., S.S., D.J.K., L.J., H.L., J.K., R.P., and D.R.
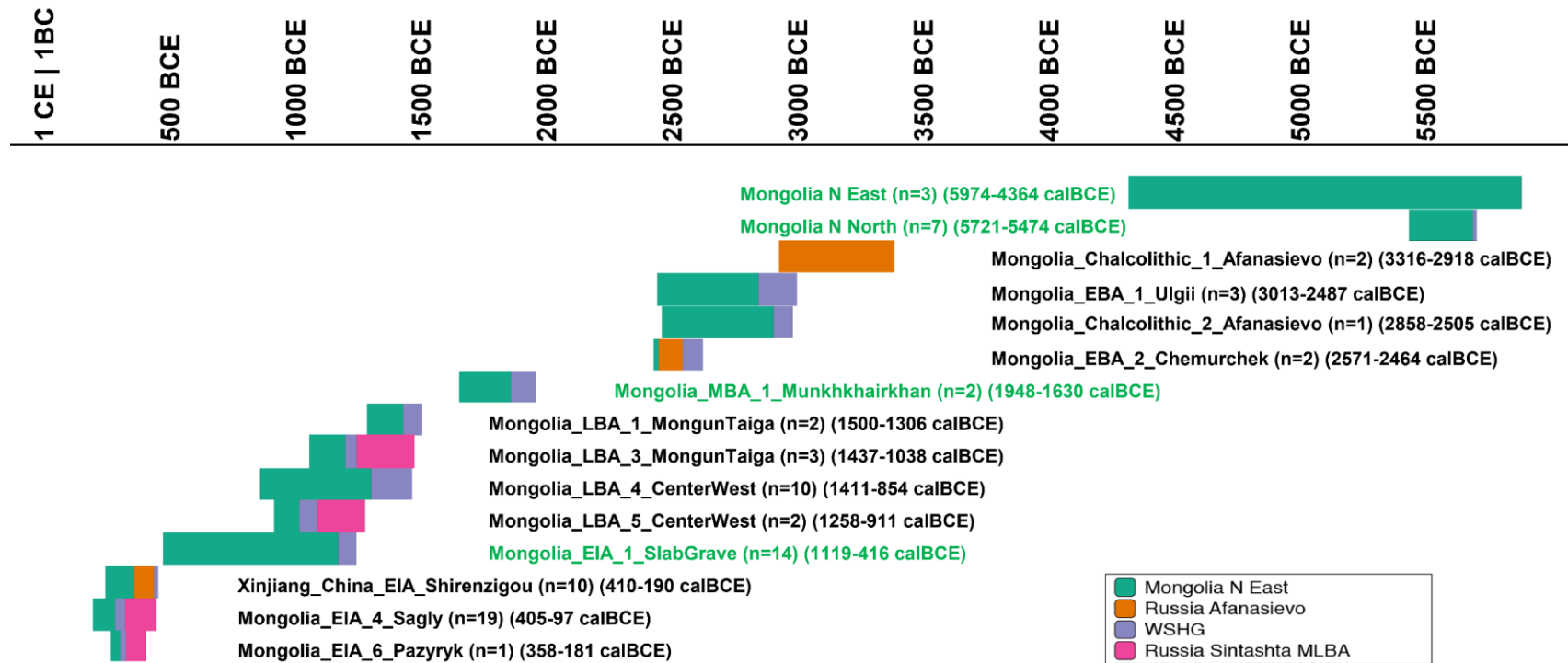
920

## Competing interests

922    The authors declare no competing interests.

27

**Figure Legends**



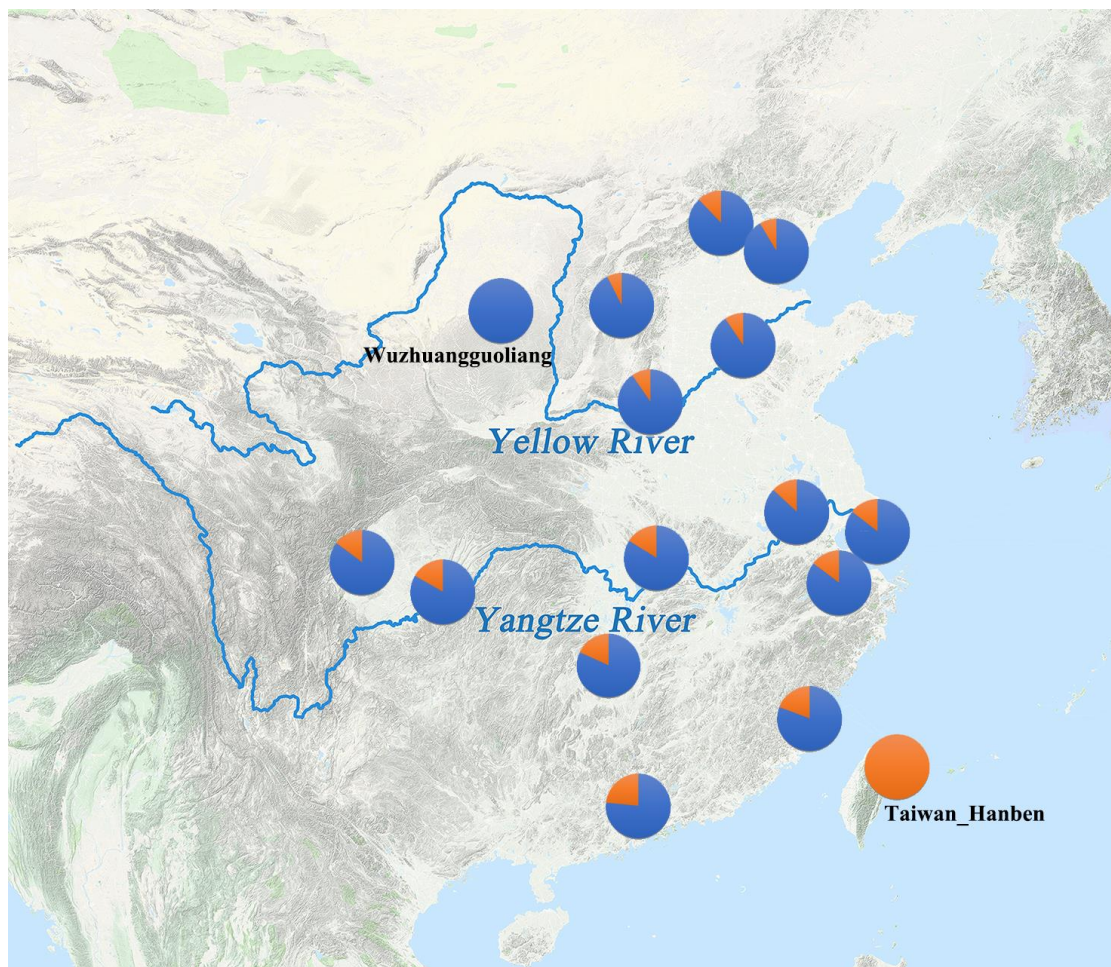| | |
|---|---|
| ▽ | Wuzhuangguoliang (3348–2862 calBCE) |
| ⊠ | Taiwan_Hanben (1402–1262 calBCE,395–536 calCE) |
| ✳ | Taiwan_Gongguan (250–650 CE) |
| △ | Boisman_MN (5371–4335 calBCE) |
| ✛ | Yankovsky_IA (996–830 calBCE) |
| ✕ | Heishui_Mohe(1050–1220 calCE) |
| ⊕ | Japan_Jomon (1500–1000 BCE) |
| ◇ | Sakhalin_HG (4000 BCE–1800 CE) |
| ☐ | Mongolia_N_East (5974–4364 calBCE) |
| ○ | Mongolia_N_North (5721–5474 calBCE) |
| △ | Mongolia_Chalcolithic_2_Afanasievo (2858–2505 calBCE) |
| ✛ | Mongolia_EBA_1_Ulgii (3013–2487 calBCE) |
| ✕ | Mongolia_MBA_1_Munkhkhairkhan (1948–1630 calBCE) |
| ◇ | Mongolia_BA_1 (1388–1134 calBCE) |
| ▽ | Mongolia_BA_2 (1619–1458 calBCE) |
| ⊠ | Mongolia_LBA_4_CenterWest (1411–854 calBCE) |
| ✳ | Mongolia_LBA_1_MongunTaiga (1500–1306 calBCE) |
| ⬖ | Mongolia_LBA_2_Ulaanzukh (1447–1296 calBCE) |
| ⋈ | Mongolia_EIA_1_SlabGrave (1119–416 calBCE) |
| ⊕ | Mongolia_EIA_3 (487–391 calBCE) |
| ⊞ | Mongolia_EIA_7_Xiongnu (3–120 calCE) |
| ● | Mongolia_EIA_8 (971–834 calBCE) |
| ⊠ | Mongolia_IA_Xianbei (383–535 calCE) |
| ◤ | Mongolia_Medieval (890–982 calCE) |
| ■ | Mongolia_Mongol (1225–1388 calCE) |
| ☐ | Mongolia_Chalcolithic_1_Afanasievo (3316–2918 calBCE) |
| ○ | Mongolia_EBA_2_Chemurchek (2571–2464 calBCE) |
| △ | Mongolia_MBA_2_Munkhkhairkhan (1879–1694 calBCE) |
| ✛ | Mongolia_LBA_5_CenterWest (1258–911 calBCE) |
| ✕ | Mongolia_LBA_3_MongunTaiga (1437–1038 calBCE) |
| ◇ | Mongolia_EIA_2 (1023–912 calBCE) |
| ⊠ | Mongolia_EIA_4_Sagly (405–97 calBCE) |
| ✳ | Mongolia_EIA_5 (397–233 calBCE) |
| ▽ | Mongolia_EIA_6_Pazyryk (358–181 calBCE) |

**Figure 1: Geographical locations of newly reported ancient individuals.** We use different colors for the two ancient Mongolia clusters.

Detailed information are given in Table S1, Online Table 1 and Supplemental Experimental Procedures.
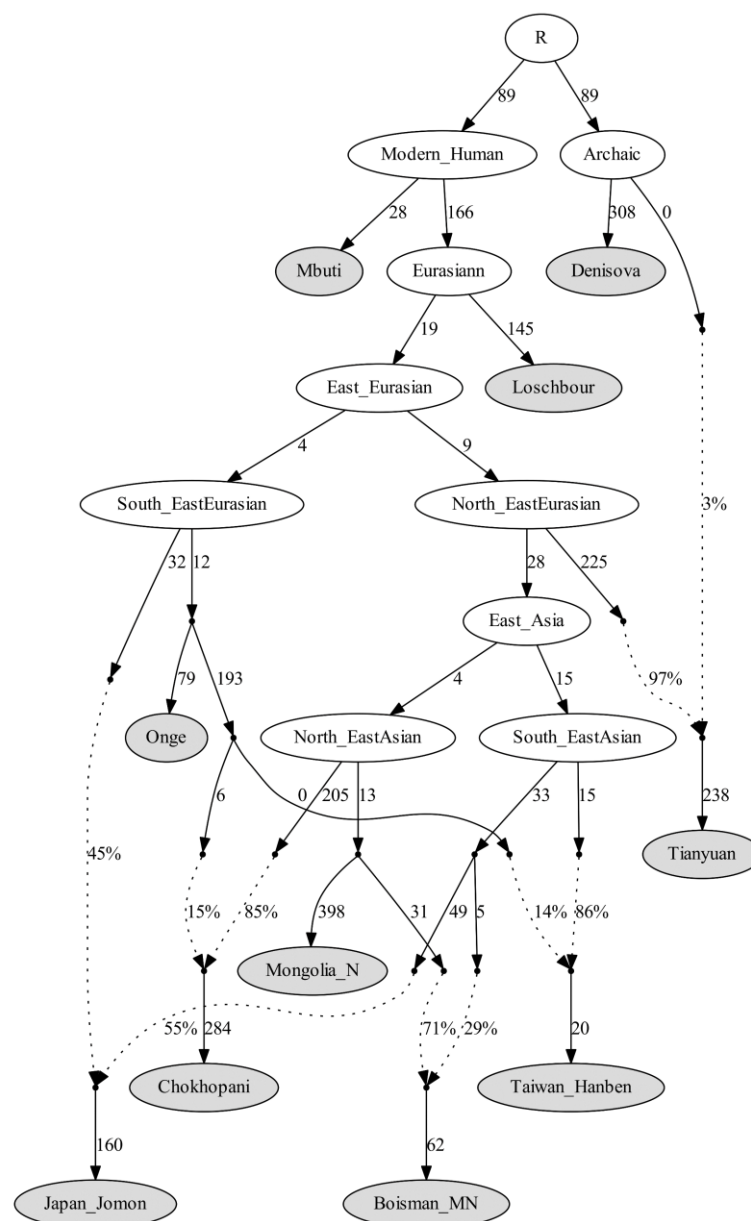
**Figure 2: Principal Component Analysis (PCA). (A)** Projection of ancient samples onto PCA dimensions 1 and 2 defined by East Asians, Europeans, Siberians and Native Americans. **(B)** Projection onto groups with the little West Eurasian mixture.

**Figure 3:** *qpAdm* **modeling of ancestry change over time in Mongolia.** We use Mongolia_East_N, Afanasievo, WSHG, and Sintashta_MLBA as sources, and for each combined archaeological and genetic grouping identify maximally parsimonious models (fewest numbers of sources) that fit with P>0.05 (Online Table 5). We plot results for groupings that give a unique parsimonious model, and include at least one individual with data that "PASS" at high quality and with a confident chronological assignment (Online Table 1). The bars show proportions of each ancestry source, and we also include time spans for the individuals in the cluster. Groupings that include more eastern individuals (longitude >102.7 degrees) are indicated in green and typically have very little Yamnaya-related admixture even at late dates.
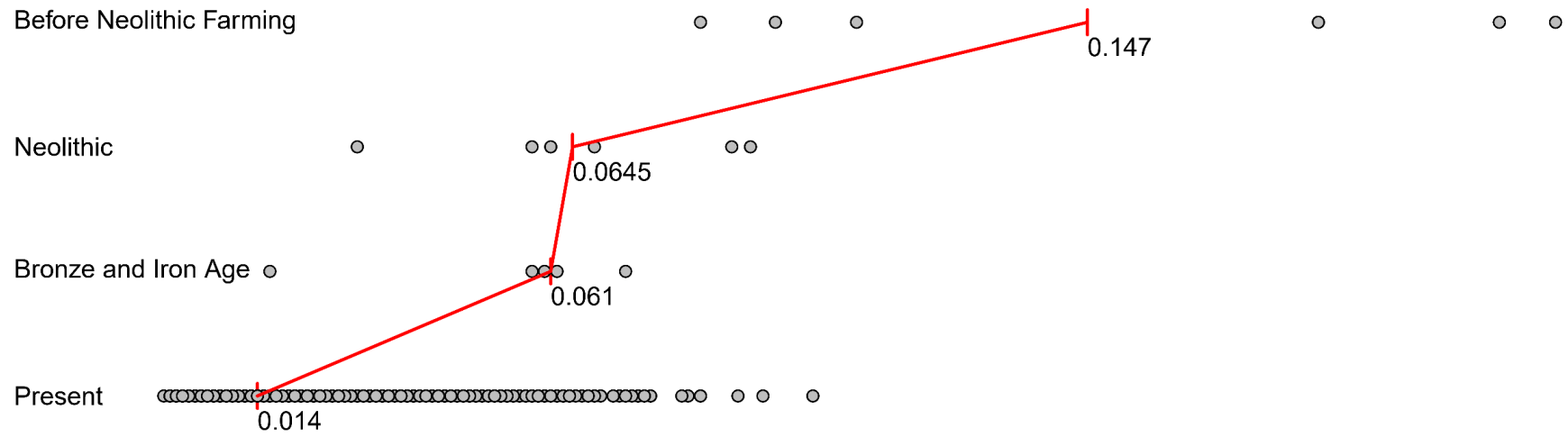
**Figure 4: *qpAdm* modeling of Han Chinese cline.** We used the ancient Wuzhuangguoliang as a proxy for Yellow River Farmers and Taiwan_Hanben as a proxy for Yangtze River Farmers related ancestry.

**Figure 5: *qpGraph* modeling of a subset of East Asians.** We used all available sites in the 1240K dataset, restricting to transversions only to replicate key results (Supplementary Information). We started with a skeleton tree that fits the data with Denisova, Mbuti, Onge, Tianyuan and Loschbour and one admixture event. We then grafted on Mongolia_East_N, Jomon, Taiwan_Hanben, Chokhopani, and Boisman in turn, adding them consecutively to all possible edges in the tree and retaining only graph solutions that provided no differences of |Z|>3 between fitted and estimated statistics. We used the MSMC relative population split time to constrain models (the maximum discrepancy for this model is |Z|=2.8). Drifts along edges are multiplied by 1000. Dashed lines represent admixture. Deep population splits are not well constrained due to a lack of data from Upper Paleolithic East Asians.

**Figure 6. Homogenization of East Asian populations through mixture.** Pairwise $F_{ST}$ distribution among populations belonging to four time slices in East Asia; the median (red) of $F_{ST}$ is shown.