

Indian genetic heritage in Southeast Asian populations

Piya Changmai^{1*}, Kitipong Jaisamut¹, Jatupol Kampuansai^{2,3}, Wibhu Kutanan⁴, N. Ezgi Altınışık^{1,#a}, Olga Flegontova¹, Angkhana Inta^{2,3}, Eren Yüncü¹, Worrawit Boonthai⁵, Horolma Pamjav⁶, David Reich^{7,8,9}, Pavel Flegontov^{1,7*}

¹ Department of Biology and Ecology, Faculty of Science, University of Ostrava, Ostrava, Czech Republic

² Department of Biology, Faculty of Science, Chiang Mai University, Chiang Mai, Thailand

³ Research Center in Bioresources for Agriculture, Industry and Medicine, Chiang Mai University, Chiang Mai, Thailand

⁴ Department of Biology, Faculty of Science, Khon Kaen University, Khon Kaen, Thailand

⁵ Faculty of Sociology and Anthropology, Thammasat University, Pathum thani, Thailand

⁶ Hungarian Institute for Forensic Sciences, Institute of Forensic Genetics, Budapest, Hungary

⁷ Department of Human Evolutionary Biology, Harvard University, Cambridge, MA, USA

⁸ Department of Genetics, Harvard Medical School, Boston, MA, USA

⁹ Broad Institute of MIT and Harvard, Cambridge, MA, USA

^{#a} Current Address: Human-G Laboratory, Department of Anthropology, Hacettepe University, Ankara, Turkey

* corresponding author,

Email: piya.changmai@osu.cz (PC)

Email: pavel.flegontov@osu.cz (PF)

23 **Abstract**

24 The great ethnolinguistic diversity found today in mainland Southeast Asia (MSEA) reflects
25 multiple migration waves of people in the past. Deeply divergent East Eurasian hunter-gatherers were
26 the first anatomically modern human population known to migrate to the region. Agriculturalists from
27 South China migrated to the region and admixed with the local hunter-gatherers during the Neolithic
28 period. During the Bronze and Iron Ages, the genetic makeup of people in MSEA changed again,
29 indicating an additional influx of populations from South China. Maritime trading between MSEA and
30 India was established at the latest 300 BCE, and the formation of early states in Southeast Asia during
31 the first millennium CE was strongly influenced by Indian culture, and this cultural influence is still
32 prominent today. Several ancient Indian-influenced states were located in present-day Thailand, and
33 various populations in the country are likely to be descendants of people from those states. To
34 systematically explore Indian genetic heritage in MSEA, we generated genome-wide SNP data (the
35 HumanOrigins array) for 119 present-day individuals belonging to 10 ethnic groups from Thailand and
36 co-analyzed them with published data from MSEA using the PCA, ADMIXTURE, f_3 -statistics, qpAdm,
37 and qpGraph methods. We found South Asian low-level admixture in various MSEA populations which
38 are probably descendants of people from the ancient Indian-influenced states, but failed to find a South
39 Asian genetic component in present-day hunter-gatherer groups and relatively isolated groups from
40 highlands in Northern Thailand. Our results also support close genetic affinity between Kra-Dai-
41 speaking (also known as Tai-Kadai) and Austronesian-speaking populations, which fits a linguistic
42 hypothesis suggesting cladality of the two language families.

43

44

45

46 **Author Summary**

47 Mainland Southeast Asia is a region with great ethnolinguistic diversity and complex population
48 history. We studied genetic population history of present-day mainland Southeast Asian populations
49 using genome-wide SNP data (the HumanOrigins array). We generated new data for 10 present-day
50 ethnic groups from Thailand, which we further combined with published data from mainland and island
51 Southeast Asians and worldwide populations. We revealed South Asian genetic admixture in various
52 mainland Southeast Asian ethnic groups which are highly influenced by Indian culture, but failed to
53 find it in groups who remained culturally isolated until recently. Our finding suggests that a massive
54 migration of Indian people in the past was responsible for the spread of Indian culture in mainland
55 Southeast Asia. We also found support for a close genetic affinity between Kra-Dai- and Austronesian-
56 speaking populations, which fits a linguistic hypothesis suggesting cladality of the two language
57 families.

58

59 **Introduction**

60 Mainland Southeast Asia (MSEA) is a region with high ethnolinguistic diversity and complex
61 population history. Hundreds of indigenous languages belonging to five language families
62 (Austroasiatic, Austronesian, Hmong-Mien, Kra-Dai, and Sino-Tibetan) are spoken in MSEA [1].
63 Anatomically modern humans migrated to MSEA around 50000 – 60000 years ago [2]. Previous
64 archaeogenetic studies indicate that the earliest MSEA individuals with genome-wide data available
65 belong to the deeply diverged East Eurasian hunter-gatherers [3]. Andamanese hunter-gatherers (Onge
66 and Jarawa) and MSEA Negritos are present-day populations related to the deeply diverged East
67 Eurasian hunter-gatherers [3-4]. Neolithic populations in MSEA were established by admixture
68 between these local hunter-gatherers and agriculturalists who migrated from South China around 4000

69 years ago [3-4]. The genetic makeup of MSEA Neolithic individuals is similar to present-day
70 Austroasiatic-speaking populations [3-4]. That pair of seminal studies also detected additional waves of
71 migrations from South China to MSEA during the Bronze and Iron Ages. Early states in MSEA during
72 the first millennium CE, such as Pyu city-states, Funan, Dvaravati, Langkasuka, and Champa were
73 established with a substantial influence from Indian culture [5]. There was evidence of Indian trading in
74 MSEA and of glass bead manufacturing by MSEA locals using Indian techniques during the Iron Age
75 [2]. Thailand is a country in the middle of MSEA, and many ancient Indianized states were located in
76 its territory [5]. In Thailand 51 indigenous languages from five language families are attested.

77 We generated genome-wide SNP genotyping data for ten populations from Thailand: six
78 Austroasiatic-speaking populations (Khmer, Kuy, Lawa, Maniq, Mon, and Nyahkur), one Hmong-
79 Mien-speaking population (Hmong), one Kra-Dai-speaking population (Tai Lue), and two Sino-
80 Tibetan-speaking populations (Akha and Sgaw Karen). Maniq, a MSEA Negrito group, are present-day
81 hunter-gatherers. We combined our data with published MSEA and world-wide data. Our results
82 revealed South Asian admixture in MSEA populations which are heavily influenced by Indian culture
83 or which can be traced back to ancient Indianized states, and we failed to detect South Asian admixture
84 in relatively isolated "hill tribes" (a term commonly used in Thailand for minority ethnic groups
85 residing mainly in the northern and western highland region of the country) or in hunter-gatherers. The
86 ubiquitous South Asian admixture in MSEA populations suggests a massive migration of South Asian
87 populations to MSEA, which correlates with the spread of Indian culture across MSEA in the past. We
88 also found Atayal-related ancestry in most Kra-Dai-speaking population in MSEA and South China,
89 and that ancestry is absent in other MSEA groups apart from those with a clear history of Austronesian
90 influence. The results suggest a link between Kra-Dai and Austronesian populations as previously
91 suggested by linguistic studies proposing the existence of the Austro-Tai language macrofamily [6].

92

93 **Results**

94 **Overview of the genetic make-up of ESEA populations**

95 Using the HumanOrigins SNP array [7], we generated genome-wide genotyping data (~574,131
96 autosomal sites) for 10 present-day human populations from Thailand (Fig 1). We merged our data with
97 published data for ancient and present-day worldwide populations (S1 table). To get an overview of
98 population structure, we performed principal component analysis (PCA) (Fig 2). South Asian (SAS)
99 populations lie on a well-known North-South cline [8]. Central Asian and Siberian populations lie
100 between the European (EUR) - SAS and East-Southeast Asian (ESEA) clines. Mainland SEA Negritos
101 occupied the space between the ESEA cline and the Andamanese (Onge). Munda populations, Austro-
102 Asiatic-speaking populations from India which were shown in a previous study [9] to be a genetic
103 mixture of South Asian and Southeast Asian populations, lie between the SAS and ESEA clines, as
104 expected (Fig 2). Populations from East and Southeast Asia form a well-defined cluster, but positions
105 of some populations such as Sherpa, Burmese, Mon, Thai, Cambodian, Cham, Ede, Malay, Khmer,
106 Nyahkur, and Kuy are shifted towards the SAS cline (Fig 2).

107 Next, we performed a model-based clustering analysis using the ADMIXTURE approach. At 12
108 hypothetical ancestral populations, Burmese, Mon, Thai, Cambodian, Cham, Ede, Malay, Khmer, Kuy,
109 and Nyahkur demonstrated an ancestry component (5% on average) which peaked in Indian
110 populations (Fig 3). Due to the diversity of Thai (from Thailand) according to the PCA and
111 ADMIXTURE analyses, we separated Thai into three groups labelled Thai1, Thai2, and Thai3. Their
112 average SAS ancestry component proportions according to the ADMIXTURE analysis were as follows:
113 15, 7, and 3%, respectively (Figs 2 and 3).

114 Outgroup f_3 -statistics are used for measuring shared genetic drift between a pair of test
115 populations relative to the outgroup population. We further explored hypothetical SAS admixture in
116 MSEA by inspecting a biplot of outgroup f_3 -tests (Fig 4 and S1 Fig). We used Mbuti as an outgroup. On

117 the y-axis, statistics $f_3(\text{Mbuti}; A, \text{test group})$ are shown, where A are East Asian surrogates (Han or Dai)
118 and "test groups" are other ESEA populations. On the x-axis statistics $f_3(\text{Mbuti}; B, \text{test group})$ are
119 shown, where B are South Asian populations (Brahmin Tiwari or Coorghi). In the coordinates formed
120 by statistics $f_3(\text{Mbuti}; \text{Han}, \text{test group})$ and $f_3(\text{Mbuti}; \text{Brahmin Tiwari}, \text{test group})$ (Fig 4), most ESEA
121 populations demonstrate a linear relationship between the genetic drift shared with Han and the drift
122 shared with Brahmin Tiwari. However, positions of Burmese, Thai, Mon, Cham, Nyah Kur, Cambodian,
123 Khmer, Malay, Giarai, and Ede are shifted from that main ESEA trend line. This shift can be interpreted as an
124 elevated shared drift between the SAS group and the test population, as compared to other ESEA populations.
125 Similar results were generated when we replaced Han and Brahmin Tiwari with Dai and Coorghi, respectively
126 (S1 Fig).

127

128 **Fitting admixture models using the qpWave, qpAdm, and qpGraph approaches**

129 We also tested specific admixture models using the qpWave [10] and qpAdm methods [11, 12].
130 Previous studies indicate that deeply diverged East Eurasian hunter-gatherers (associated with the
131 Hoabinhian archaeological culture), which are related to present-day Andamanese hunter-gatherers
132 (Onge), were the first known anatomically modern humans who occupied MSEA [3, 4]. MSEA
133 populations in the Neolithic period can be modelled as a mixture of local Hoabinhians and populations
134 who migrated from East Asia [3, 4]. In our analysis, we used Atayal, Dai, and Lahu as ESEA
135 surrogates. These populations speak languages which belong to three different language families:
136 Austronesian, Kra-Dai, and Sino-Tibetan (the Tibeto-Burman branch), respectively. Onge was used as a
137 surrogate for the deeply diverged East Eurasian hunter-gatherers. 55 populations composed of at least 5
138 individuals were used as South Asian surrogates. Outgroups ("right populations") for all qpWave and
139 qpAdm analyses were the following present-day populations: Mbuti (Africans), Palestinians, Iranians

140 (Middle Easterners), Armenians (Caucasians), Papuans [7], Nganasans, Kets, Koryaks (Siberians),
141 Karitiana (Native Americans), Irish, and Sardinians (Europeans).

142 We first explored cladality of population pairs using qpWave (Fig 5, S2 Table). In other words,
143 we tested if one stream of ancestry from an ESEA surrogate is sufficient to model a Southeast Asian
144 target population. We used a cut-off p-value of 0.05. We further tested 2-way and 3-way admixture
145 models using qpAdm. We applied three criteria for defining plausible admixture models: a) all simpler
146 models should be rejected according to the chosen p-value cutoff; b) the current model should not be
147 rejected according to the chosen p-value cutoff; c) inferred admixture proportions ± 2 standard errors
148 should lie between 0 and 1 for all ancestry components. If a model meets all three criteria, we consider
149 the model as "fitting" or "passing" (S2 Table), although we caution that the only secure interpretation of
150 qpWave or qpAdm tests is in terms of model rejection, and not model fit [13]. For testing 2-way and 3-
151 way admixture, we constructed models "ESEA + Onge" and "ESEA + Onge + South Asian",
152 respectively (Fig 5, S2 Table).

153 Next, we tested more explicit demographic models using qpGraph. We first constructed two
154 skeleton graphs using different SAS surrogates, Coorghi (Fig 6A) and Palliyar (Fig 6B). The worst
155 residuals for the skeleton graphs were 2.43 and 2.24 SE intervals, respectively. Skeleton graph
156 construction is explained in Materials and methods. We then exhaustively mapped target ESEA
157 populations on all possible edges (except for edge0 in S2 Fig) on the skeleton graphs. We modeled the
158 target populations as unadmixed (33 models per target population per skeleton graph), 2-way admixed
159 (528 models per target population per skeleton graph), and 3-way admixed (5,456 models per target
160 population per skeleton graph). We compared models with different numbers of admixture sources
161 using a log-likelihood difference cut-off of 10 log-units or a worst residual difference cut-off of 0.5 SE
162 intervals (see exploration of appropriate cut-offs on simulated genetic data in Ning et al., 2020 preprint
163 [13]). For models with the same number of admixture sources, we used a log-likelihood difference cut-

164 off of 3 log-units [14]. We also avoided models with trifurcations, i.e., when drift length on any
 165 "backbone" edge equals zero. Below we discuss best models found for the studied populations grouped
 166 by language family. The summary of qpWave, qpAdm, and qpGraph results is presented in Table 1.
 167 Full results are shown in S2 Table (qpWave and qpAdm) and S3 Table (qpGraph). S3 Table shows all
 168 qpGraph models satisfying the log-likelihood difference criteria. Edge number codes for the Coorghi
 169 and Palliyar skeleton graphs are illustrated in S2 Fig.

170

171 **Table 1. A summary of qpAdm and qpGraph admixture modelling results for the groups of**
 172 **interest.** Labels of groups genotyped in this study are italicized.

population	n	language family	country	qpAdm, best model	qpGraph, best model
Cambodian	9	Austroasiatic	Cambodia	ESEA + NEG + SAS	Atayal + Mlabri + SAS
Htin	10	Austroasiatic	Thailand	ESEA + NEG	Mlabri
<i>Khmer</i>	10	Austroasiatic	Thailand	ESEA + NEG or ESEA + NEG + SAS	Mlabri + SAS
<i>Kuy</i>	10	Austroasiatic	Thailand	ESEA + NEG or ESEA + NEG + SAS	Mlabri + SAS
<i>Lawa</i>	10	Austroasiatic	Thailand	ESEA + NEG	Tibetan + Mlabri
<i>Maniq</i>	9	Austroasiatic	Thailand	ESEA + NEG	Atayal + NEG
Mlabri	10	Austroasiatic	Thailand	ESEA + NEG	included in the skeleton graphs
<i>Mon</i>	10	Austroasiatic	Thailand	ESEA + NEG + SAS	before Tibetan + Mlabri (ESEA source) + SAS
<i>Nyahkur</i>	9	Austroasiatic	Thailand	ESEA + NEG + SAS	Mlabri + SAS
Cham	10	Austronesian	Vietnam	ESEA + NEG + SAS	Atayal + Mlabri + SAS (western source)
Ede	9	Austronesian	Vietnam	ESEA + NEG + SAS	Mlabri + SAS
Giarai	11	Austronesian	Vietnam	ESEA + NEG + SAS	Mlabri + SAS
Malay	5	Austronesian	Singapore	ESEA + NEG or ESEA + NEG + SAS	Atayal + Mlabri + SAS
<i>Hmong</i>	10	Hmong-Mien	Thailand	ESEA + NEG	before Atayal + Tibetan
<i>Tai Lue</i>	9	Kra-Dai	Thailand	ESEA	before Dai/Mlabri + Mlabri
Thai1	7	Kra-Dai	Thailand	ESEA + NEG + SAS	before Tibetan + Mlabri (ESEA source) + SAS
Thai2	2	Kra-Dai	Thailand	ESEA + NEG or ESEA + NEG + SAS	before Atayal + SAS
Thai3	1	Kra-Dai	Thailand	ESEA or ESEA + NEG	Atayal + Mlabri
<i>Akha</i>	31	Sino-Tibetan	Thailand	ESEA	Tibetan + Mlabri (ESEA source)
Burmese	6	Sino-Tibetan	Myanmar	ESEA + NEG + SAS	Tibetan + Mlabri + SAS
<i>Sgaw Karen</i>	10	Sino-Tibetan	Thailand	ESEA + NEG	Tibetan + Mlabri

173

174

175 **Sino-Tibetan (Tibeto-Burman branch).** We studied Akha, Sgaw Karen from Thailand, and Burmese
 176 (15) from Myanmar. All three groups harbor ancestry from a Tibetan-related source (S3 Table). Akha
 177 was modeled as one stream of ancestry when Lahu was used as an ESEA surrogate in qpWave (S2
 178 Table). Sgaw Karen requires an extra ancestry from the Onge surrogate in qpAdm analysis (S2 Table).
 179 The result agrees with qpGraph analysis where Sgaw Karen was modeled as a mixture of a Tibetan-

180 related and a Mlabri-related source (S3 Table). Mlabri harbor a substantial proportion of deeply
181 diverged East Eurasian ancestry (Fig 6). Additional gene flow from deep sources (edge7 and edge8) to
182 Karen on the Coorghi skeleton decreased the worst residual by ~ 0.5 SE intervals, but the inferred
183 admixture proportion was close to zero; therefore, these additional edges could be an artifact. Both
184 qpAdm and qpGraph analyses indicated South Asian ancestry in Burmese: e.g. $\sim 12\%$ inferred by
185 qpAdm (S2-3 Table). Burmese harbor ancestry from Tibetan-related + Mlabri-related + South Asian
186 sources according to a best-fitting graph model (S3 Table).

187

188 **Hmong-Mien.** We analyzed Hmong from Thailand. We were not able to model Hmong as cladal with
189 any of our three standard ESEA surrogates (Atayal, Dai, and Lahu). Then we tried to use Miao, a
190 Hmong-Mien-speaking population from China, as an ESEA surrogate. We successfully modeled
191 Hmong as Miao + Onge (S2 Table). The Hmong groups from Thailand and from Vietnam [16] are
192 cladal according to qpWave (S2 Table). Our qpGraph result showed a low level of Tibetan-related
193 ancestry ($\sim 2\%$) in Hmong (S3 Table).

194

195 **Austronesian.** There are four Austronesian-speaking populations included in this study: Cham, Ede
196 (Rade), and Giarai (Jarai) from Vietnam [16], and Malay from Singapore [15]. qpAdm and qpGraph
197 results revealed South Asian ancestry in all four Austronesian groups: 11.6%, 7.5%, 7.4%, and 2.1% in
198 Cham, Ede, Giarai, and Malay, respectively, as inferred by qpAdm (S2-3 Table). Atayal is an
199 Austronesian-speaking group from Taiwan, the homeland of Austronesian languages [17]. We failed to
200 detect Atayal-related ancestry in Ede and Giarai (S3 Table), while the ancestry is present in Cham and
201 Malay. We found Mlabri-related ancestry in all four Austronesian-speaking populations (S3 Table)

202

203 **Austroasiatic.** We studied Htin [4], Khmer, Kuy, Lawa, Maniq, Mlabri [4], Mon, Maniq, and Nyahkur
204 from Thailand, and Cambodians from Cambodia [7]. Maniq, a present-day hunter-gatherer Negrito
205 group from Southern Thailand, has a major ancestry component derived from a deeply diverged East
206 Eurasian group with ~74 % admixture proportion inferred by qpAdm (S2-3 Table). The ESEA source
207 for Maniq is Atayal-related (S3 Table). Htin was modeled as a sister group of Mlabri by qpGraph (S3
208 Table). Both groups were modelled by qpAdm as having ESEA and Onge-related ancestry (S2-3
209 Table). Lawa was modeled as Mlabri-related + Tibetan-related ancestry (S3 Table). We detected South
210 Asian admixture in five Austroasiatic-speaking groups in our study (Cambodian, Khmer, Kuy, Mon,
211 and Nyahkur): 9.4%, 4.6%, 4.3%, 11.6%, and 7%, respectively, as inferred by qpAdm. Khmer, Kuy,
212 and Nyahkur showed similar genetic makeups (S2-3 Table). We observed Atayal-related ancestry in
213 Cambodian (S3 Table) and Tibetan-related ancestry in Mon, and these ancestry sources are rare in other
214 Austroasiatic speaking populations.

215

216 **Kra-Dai.** We tested Kra-Dai-speaking populations from China (Dong, Dong Hunan, Gelao, Li,
217 Maonan, Mulam, and Zhuang from Wang et al., 2020 preprint [18]) and Vietnam (Boy, Colao, Lachi,
218 Nung, Tay, and Thai from Liu et al., 2020 [16]), and Thailand (Tai Lue from this study, Thai1, Thai2,
219 and Thai3 from Lazaridis et al., 2014 [19]). Most of the Kra-Dai-speaking populations from China and
220 Vietnam harbor Tibetan-related and Atayal-related ancestry (S3 Table). The Thai3 from Thailand was
221 modelled as getting ~56% of its ancestry from a sister group of Atayal (S3 Table). Thai2 harbors
222 ancestry from a source diverging before Atayal (S3 Table). Atayal-related ancestry is missing in Thai1
223 (S3 Table), but we found a source diverging before Tibetan Chokhopani when we mapped the Thai1
224 population on the Coorghii skeleton (S3 Table). We observed South Asian ancestry in Thai1 and Thai2,
225 but that ancestry is missing in Tai Lue and Thai3 (S2-3 Table). qpAdm inferred South Asian admixture
226 proportions in Thai1 and Thai2 at 17% and 5%, respectively.

227

228 **Discussion**

229 Indian culture was long established in MSEA, which also influenced early states formation in
230 the region during the first millennium CE [5]. Previous studies reported South Asian admixture in few
231 populations from Southeast Asia [20-22]. Some studies used the same or similar populations as those in
232 the current study but did not focus on South Asian admixture [16, 19, 23]. In this study, we thoroughly
233 analyzed South Asian admixture in present-day Southeast Asia. We also investigated the genetic
234 markup of populations in the region. Our results were consistent across various methods used in this
235 study (ADMIXTURE, f_3 -statistics, qpAdm, qpGraph). There were just one or a few admixture graph
236 models which fitted the data significantly better than ca. 6000 other models we tested per target
237 population. qpAdm and qpGraph results were in agreement: adding a South Asian-related admixture
238 edge never improved qpGraph model fits significantly when a 3-way model with South Asian
239 admixture was rejected by qpAdm. We discuss the results by language family below.

240

241 **Sino-Tibetan (Tibeto-Burman branch)**

242 Using qpWave, we were not able to reject cladality of Akha and Lahu, two Sino-Tibetan-
243 speaking populations (S2 Table). Sgaw Karen required an extra stream of ancestry from an Onge-
244 related population (S2 Table). The Onge-related ancestry in Sgaw Karen can be explained by admixture
245 with an Austroasiatic-speaking population, which harbors high genetic ancestry from Hoabinhians[3,
246 4]. Our best-fitting admixture graph model for Sgaw Karen includes genetic contribution from a
247 Mlabri-related group, which fits this explanation (S3 Table). The high worst residual of the best-fitting
248 graph including Sgaw Karen probably reflects absence of an important ancestry source on our skeleton
249 graph. Our qpAdm and qpGraph results consistently demonstrated that Burmese from Myanmar harbor

250 ancestry from South Asian populations. All three Sino-Tibetan-speaking populations tested (Akha,
251 Karen, and Burmese) have Tibetan ancestry according to the best-fitting qpGraph models (S3 Table).

252

253 **Hmong-Mien**

254 The best-fitting qpAdm model for Hmong was Miao + Onge, with a minimal admixture
255 proportion from the latter source. Cladality with Miao, another Hmong-Mien speaking population, was
256 rejected (S2 Table). qpGraph modeling also indicated a low-level gene flow (~2%) from a sister group
257 of Tibetan Chokhopani (S3 Table). The main ESEA ancestry for Hmong is a source diverging before
258 Atayal (S3 Table).

259

260 **Austronesian**

261 Malay from Singapore was modeled by qpGraph as a 3-way admixture involving sister groups
262 of Atayal, Mlabri, and South Asian populations (S3 Table). Malay is an Austronesian language. It is not
263 surprising that the Malay harbor some ancestry from a source related to Atayal, an Austronesian-
264 speaking population from Taiwan. A previous study reported admixture from an Austroasiatic-speaking
265 population in Austronesian populations from Indonesia [4]. We also detected the same signal in Malay,
266 which is represented by ancestry from a sister group of Mlabri (S3 Table). Our results generated relying
267 on various approaches indicate South Asian admixture in Malay and also in three other Austronesian-
268 speaking populations from Vietnam, i.e., Cham, Ede, and Giarai (S2-3 Table). Y-haplogroups of West
269 Eurasian origin (R1a-M420 and R2-M479) were reported in Ede and Giarai by Machold et al. 2019
270 [24], and Y-haplogroups R-M17 and R-M124 were reported in Cham by He et al. 2012 [25]. Using
271 qpGraph, we were able to confirm the Atayal-related ancestry in Cham, but that gene flow signal was
272 not supported in the case of Ede and Giarai (S3 Table). The results are consistent with a previous study

273 by Liu et al. 2020 [16], which supports the spread of Austronesian language by cultural diffusion in
274 Ede and Giarai.

275

276 **Austroasiatic**

277 Htin can be modeled by qpGraph as a sister group of Mlabri (S3 Table). Both Mlabri and Htin
278 languages belong to the Khmuic branch of the Austroasiatic family. A previous study showed that
279 Mlabri has a genetic profile similar to early Neolithic individuals from mainland Southeast Asia [4].
280 The qpGraph best-fitting models for Maniq, a mainland Negrito group, incorporate 2-way admixture
281 between an Atayal-related source and an Onge-related source, with predominant genetic contribution
282 from the latter source. Even though Maniq speak an Austroasiatic language, a better surrogate for their
283 ESEA source was Atayal, an Austronesian-speaking population (S3 Table). Maniq may harbor Atayal-
284 related ancestry from Austronesian-speaking populations in Southern Thailand (where they reside) or
285 from Malaysia nearby. Using qpGraph, we could model Lawa as a 2-way admixture between a sister
286 group of Tibetan Chokhopani and Mlabri-related ancestry, with predominant contribution from the
287 latter source (S3 Table). The Austroasiatic-speaking Lawa likely got Tibetan-related ancestry via Sgaw
288 Karen. Around 1850, Sgaw Karen started migrating from present-day Myanmar to the region that was
289 once exclusively occupied by Lawa [26]. There are villages where both Lawa and Sgaw Karen live
290 alongside each other [27], and intermarriage between the two groups became more common recently
291 [28]. A previous study [22] also observed genetic interaction between Karen and Lawa. We detected a
292 minor South Asian admixture component (~4-5%) in Kuy using both qpAdm and qpGraph methods
293 (S2-3 Table). Kutanan et al. 2019 [29] reported the presence of a West Eurasian Y-haplogroup
294 R1a1a1b2a1b (R-Y6) in Kuy.

295 In this study, we generated new data for Austroasiatic-speaking Khmer from Thailand. Khmer is
296 the official language of Cambodia, and Khmer is the majority population of Cambodia [1]. Our

297 admixture graph modeling showed that Khmer from Thailand and Cambodians harbor two ancestry
298 sources in common: a Mlabri-related source and South Asian ancestry (S3 Table). West Eurasian Y-
299 haplogroups R1a1a1b2a2a (R-Z2123) and R1a1 were reported in Khmer [29] and Cambodians [30],
300 respectively. The best-fitting model for Cambodians includes additional ancestry from an Atayal-
301 related (i.e., Austronesian) source (S3 Table). Cambodians likely got this ancestry via Cham due to the
302 long-lasting interaction between the ancient Cambodian and Cham Kingdoms [5]. Cham is also the
303 largest ethnic minority in Cambodia today [1].

304 Mon and Nyahkur languages belong to the Monic branch of the Austroasiatic family [1]. Our
305 qpGraph modeling found Mlabri-related and South Asian ancestry in both populations. A previous Y-
306 chromosome study [29] reported various haplogroups of West Eurasian origin, such as J and R, in Mon,
307 and haplogroup J2a1 (J-L26) in Nyahkur. The higher frequencies of West Eurasian Y-haplogroups in
308 Mon correspond to the higher South Asian admixture proportion found in Mon as compared to
309 Nyahkur. Mon harbors additional ancestry from a source close to Tibetan Chokhopani (S3 Table).
310 Tibetan-related ancestry is missing in Nyahkur (S3 Table). The Nyahkur group is possibly a remnant of
311 an ancient Monic-speaking population from the Dvaravati kingdom located within present-day
312 Thailand [31]. Mon probably got the Tibetan-related ancestry via interactions with Sino-Tibetan-
313 speaking populations in Myanmar. Most of present-day Mon in Thailand are descendants of refugees
314 who migrated from Myanmar in the last few centuries [32]. There is some debate about the origin of
315 Mon in the Lamphun province, whether they are the direct descendants of people from the ancient Mon
316 state in present-day Thailand (ca. 1300 years before present), or their ancestors migrated from
317 Myanmar in the last few hundred years. Our results favor the latter possibility due to the Tibetan-
318 related genetic component found in Mon from Lamphun, which may reflect interaction with Burmese
319 or other Sino-Tibetan-speaking populations in Myanmar where the density of Sino-Tibetan-speaking

320 populations is much greater than in Thailand [1]. Furthermore, the Tibetan-related ancestry is absent in
321 Nyahkur, another Monic-speaking population from Thailand.

322

323 **Kra-Dai**

324 Atayal-related ancestry was found in most Kra-Dai-speaking populations in China and Vietnam,
325 according to our analysis (S3 Table). We also observed Atayal-related ancestry in Thai3 from Thailand
326 (S3 Table). Besides the Kra-Dai speakers, we were able to detect Atayal-related ancestry only in
327 Austronesian-speaking populations (Malay, Cham) or non-Austronesian populations which have
328 historical evidence of interactions with Austronesians such as Maniq and Cambodian (S3 Table).
329 Furthermore, when we used Atayal as an ESEA surrogate in 3-way qpAdm models (ESEA + Onge +
330 SAS), most of the models were rejected. Only the models with Thais (Thai1 and Thai2) as target
331 populations were not rejected (S2 Table). The genetic link between Austronesian-speaking and Kra-
332 Dai-speaking populations may reflect genealogical relationship of the two language families as
333 suggested by the Austro-Tai hypothesis [6]. Tai Lue is one of Dai ethnic groups originating in South
334 China [33]. The Tai Lue volunteers in our study migrated to Thailand less than a century ago from
335 Myanmar. Cladality of Tai Lue with all three ESEA surrogates was not rejected using qpWave (S2
336 Table). However, qpGraph modeling supported a more complex model for Tai Lue: 2-way admixture
337 between a source close to Dai and either a Mlabri-related or a source diverging before Atayal (S3
338 Table). The result suggests that after the migration from China, Tai Lue admixed with local MSEA
339 populations, or that the genetic makeup of the Dai group that gave rise to the Tai Lue group studied
340 here was different from the Dai groups sampled previously [34]. qpGraph modeling revealed different
341 genetic makeups for the three Thai sub-groups delineated in this study (S3 Table). Both qpAdm and
342 qpGraph methods consistently supported South Asian admixture in both Thai1 and Thai2 groups (S2-3
343 Table). Best-fitting models for Thai1 and Thai2 include different ESEA sources. The ESEA ancestry in

344 the Thai1 group can be traced to a source close to Dai and possibly an additional source that diverged
345 before Tibetans (S3 Table). The latter source may reflect admixture with a group that harbors a distinct
346 ESEA source, such as Chinese Han. Chinese were estimated to comprise at least 10% of the Thailand
347 population [35-36]. The Thai2 group was modelled having ESEA ancestry from a source close to
348 Atayal (S3 Table). We failed to detect South Asian ancestry in Thai3, in contrast to Thai1 and Thai2.
349 The best qpGraph model for the Thai3 group is a 2-way mixture between sister groups of Mlabri and
350 Atayal (S3 Table). Our results revealed a considerable diversity of the Thai. Previous studies also
351 reported differences in the genetic makeup of the Thai from different locations [20, 22, 29]. Samples of
352 all Thai individuals included in this study were obtained from the European Collection of Cell Cultures,
353 and we cannot trace the origin of the samples in that collection [19]. Systematic sample collection at
354 various locations will likely provide insight into the genetic diversity of the Thai.

355 Our study revealed substantial South Asian admixture in various populations across Southeast
356 Asia (~2-16% as inferred by qpAdm). We observed South Asian admixture in some populations (Cham,
357 Ede, Giarai, Khmer, Kuy, Nyahkur, and Thai) for whom the admixture was not reported before [16, 19,
358 23]. Most populations harboring South Asian admixture were heavily influenced by Indian culture in
359 the past or are related to descendants of ancient Indianized states in Southeast Asia. In contrast, we
360 failed to detect South Asian admixture in most "hill tribes" and in present-day hunter-gatherer groups
361 from Thailand. Consequently, the spread of Indian influence in the region can be explained by
362 extensive movement of people from India rather than by cultural diffusion only. The distance from the
363 coast may affect South Asian gene flow as central and southern Thai harbor South Asian ancestry, but
364 the ancestry is missing in northern Thai, who reside a long distance from the sea [22]. In this study, we
365 also observed genetic diversity in Thai, but the exact location of the Thai individuals analyzed here is
366 unknown. We detected subtle differences in populations with similar ethnolinguistic backgrounds, such
367 as Khmer from Thailand and a Khmer-speaking population (Cambodian) from Cambodia. We observed

368 Atayal-related ancestry (~3-38% as inferred by qpGraph) in most Kra-Dai-speaking populations from
369 China, Vietnam and in one group from Thailand. The results suggest a genetic connection between
370 Austronesian and Kra-Dai-speaking populations.

371

372 **Materials and methods**

373 **Sampling**

374 Sample collection and DNA extraction for all new Thailand populations in this study apart from
375 Akha was described in previous studies [23, 37-41]. Saliva samples were obtained from volunteers with
376 signed informed consent from four Akha villages in the Chiang Rai province, Thailand. The study was
377 approved by the Ethic Committee of Khon Kaen University. We performed DNA extraction as
378 described elsewhere [42]. See a list of individuals for whom genetic data is reported in this study in S4
379 Table.

380

381 **Dataset preparation**

382 Diploid genome-wide SNP data was generated using the HumanOrigins SNP array [7]. We
383 merged the new data with published ancient and present-day world-wide populations (S1 Table) using
384 PLINK v. 1.90b6.10 (<https://www.cog-genomics.org/plink/>). We first combined present-day
385 populations and applied a per site missing data threshold of 5% to create a dataset of 574,131
386 autosomal SNPs. We then added data from ancient populations. The Upper Paleolithic individual from
387 Goyet had the highest missing data percentage per individual (30%). We used the dataset for all
388 analyses except for ADMIXTURE.

389

390 **PCA**

391 The principal component analysis (PCA) was performed using PLINK v. 1.90b6.10
392 (<https://www.cog-genomics.org/plink/>) on selected populations (S1 Table) from the following regions:
393 Central, East, Southeast, and South Asia, Andamanese Islands, Siberia, and Europe.

394

395 **ADMIXTURE**

396 We performed LD filtering using PLINK v. 1.90b6.10 with the following settings: window size
397 = 50 SNPs, window step = 5 SNPs, r^2 threshold = 0.5 (the PLINK option "--indep-pairwise 50 5 0.5").
398 LD filtering produced a set of 270,700 unlinked SNPs. We carried out clustering analysis using
399 ADMIXTURE v. 1.3 (<https://dalexander.github.io/admixture/download.html>), testing from 8 to 13
400 hypothetical ancestral populations (K) with tenfold cross-validation. We performed five iterations for
401 each value of K. We selected K = 12 for presentation according to the highest model likelihood. We
402 further ran up to 30 iterations for K = 12 and ranked them by model likelihood.

403

404 **Outgroup f_3 -statistics**

405 We computed f_3 -statistics [7] using qp3Pop v. 420, a software from the ADMIXTOOLS package
406 (<https://github.com/DReichLab/AdmixTools>). We ran $f_3(\text{Mbuti}; X, \text{test group})$, where X are East Asian
407 surrogates (Han or Dai) or South Asian (Brahmin Tiwari or Coorgi) surrogates. The test groups are
408 various ESEA populations.

409

410

411

412 **qpWave and qpAdm**

413 We used qpWave v. 410 and qpAdm v. 810 from the ADMIXTOOLS package. We used the
414 following populations as outgroups ("right populations") for all qpWave and qpAdm analyses: Mbuti
415 (Africans), Palestinians, Iranians (Middle Easterners), Armenians (Caucasians), Papuans [7], Nganasan,
416 Kets, Koryaks (Siberians), Karitiana (Native Americans), Irish, and Sardinians (Europeans). We used
417 Atayal, Dai, and Lahu as ESEA surrogates. We used Onge as a surrogate for the deeply diverged East
418 Eurasian hunter-gatherers. We used 55 different populations as alternative South Asian surrogates (S2
419 Table).

420 We tested a pair of a test population and an ESEA surrogate using qpWave. We used a cut-off p-
421 value of 0.05 for qpWave modeling. We performed 2-way and 3-way admixture modeling using
422 qpAdm. 2-way admixture was modeled as "target population = ESEA surrogate + Onge", and 3-way
423 admixture was modeled as "target population = ESEA surrogate + Onge + SAS surrogate". We applied
424 three criteria for defining plausible admixture models: a) all simpler models should be rejected
425 according to the chosen p-value cutoff; b) the current model should not be rejected according to the
426 chosen p-value cutoff; c) inferred admixture proportions ± 2 standard errors should lie between 0 and 1
427 for all ancestry components.

428

429 **qpGraph**

430 We used qpGraph v. 6412 from the ADMIXTOOLS package with the following settings:
431 outpop: NULL, blgsize: 0.05, lsqmode: NO, diag: 0.0001, hires: YES, initmix: 1000, precision: 0.0001,
432 zthresh: 0, terse: NO, useallsnps: NO. We used the following criteria to select the best-fitting model.
433 Models with different numbers of admixture sources were compared using a log-likelihood difference

434 cut-off of 10 log-units or a worst residual difference cut-off of 0.5 SE intervals [13]. We used a log-
435 likelihood difference cut-off of 3 log-units for models with the same number of parameters [14].

436 We started with the following five populations: Denisovan (archaic human), Altai Neanderthal
437 (archaic human), Mbuti (African), Atayal (East Asian), and Goyet (ancient West European hunter-
438 gatherer). A best-fitting model is illustrated in S3 Fig. We fixed Neanderthal-related (node nA in S3
439 Fig) admixture proportion in non-Africans at 3%. Goyet requires extra admixture from this
440 Neanderthal-related source. When this admixture edge was missing, the worst f -statistic residual
441 increased from 2.13 to 4.56. We further mapped additional populations on the graph, one at a time. We
442 mapped a new population on all possible edges on the graph as unadmixed, 2-way, and 3-way admixed.
443 We mapped Onge on the 5-population graph (S3 Fig) and then Dai on the 6-population skeleton graph
444 (S4 Fig). Best-fitting graphs including Onge and Dai are shown in S4 Fig and S5 Fig, respectively.

445 We further mapped an ancient Iranian herder individual from Ganj Dareh (I1947 [8]). A best-
446 fitting model for this individual is a 2-way mixture between a putative West Eurasian source and a
447 basal Eurasian source (S6 Fig). Basal Eurasian admixture in ancient Iranians was reported in a previous
448 study [43]. Mlabri can be modeled as ESEA + Onge-related sources (S7 Fig), which is consistent with a
449 previous study [4].

450 We mapped South Asian populations, Coorghi or Palliyar, on the graph in S7 Fig. Both
451 populations can be modeled as a 2-way mixture between ancient Iranian-related and deep-branching
452 East Eurasian sources (S8A and B Fig). The positions of the deep East Eurasian source for Coorghi and
453 Palliyar are slightly different, but both are among the deepest East Eurasian branches.

454 Next, we added an ancient individual, Tibetan Chokhopani from Nepal (S1 Table), as the last
455 population on the skeleton graphs. The best-fitting model for this individual was an unadmixed branch
456 in the ESEA clade before the divergence of Atayal (Fig 6A and B). The total numbers of SNPs used for

457 fitting the skeleton graphs with Coorghi and Palliyar were 311,259 and 317,327, and the worst absolute
458 f -statistic residuals were are 2.43 and 2.24 SE, respectively.

459 We mapped present-day target populations on all possible edges (except for edge0 in S2 Fig) on
460 the skeleton graphs as unadmixed, 2-way admixed, and 3-way admixed. In total, we tested 6,017
461 models per target population per skeleton graph.

462

463 **Data Availability**

464 Genome-wide genotyping data generated for this study will be made publicly available when
465 the manuscript is published.

466

467 **Acknowledgments**

468 We thank all volunteers in Thailand who donated the samples for our study. We thank Phangard
469 Neamrat, Jaeronchai Chuaychu, and Prateep Panyadee for assisting with sample collection. This work
470 was supported by the Czech Ministry of Education, Youth and Sports: 1) Inter-Excellence program,
471 project #LTAUSA18153; 2) Large Infrastructures for Research, Experimental Development and
472 Innovations project "IT4Innovations National Supercomputing Center – LM2015070". E.Y., O.F., P.C.,
473 and P.F., were also supported by the Institutional Development Program of the University of Ostrava.
474 J.K. acknowledges partial support provided by Chiang Mai University, Thailand.

475

476 **References**

477 1. Eberhard DM, Simons GF. Ethnologue: Languages of Asia. 23rd ed. Fennig CD, editor. SIL In-
478 ternational, Global Publishing; 2020.

- 479 2. Higham C. Early Mainland Southeast Asia: From First Humans to Angkor. River Books Press
480 Dist A C; 2014.
- 481 3. McColl H, Racimo F, Vinner L, Demeter F, Gakuhari T, Moreno-Mayar JV, et al. The prehistor-
482 ic peopling of Southeast Asia. *Science*. 2018;361(6397): 88– 92.
- 483 4. Lipson M, Cheronet O, Mallick S, Rohland N, Oxenham M, Pietrusewsky M, et al. Ancient ge-
484 nomes document multiple waves of migration in Southeast Asian prehistory. *Science*.
485 2018;361(6397): 92-95.
- 486 5. Cœdès G. The Indianized states of Southeast Asia. Walter F. Vella, editor. University of Hawaii
487 Press; 1968.
- 488 6. Ostapirat W. Kra-dai and Austronesian: notes on phonological correspondences and vocabulary
489 distribution. In: Sagart L, Blench R, Sanchez-Mazas A, editors. *The peopling of East Asia: put-*
490 *ting together archaeology, linguistics and genetics*. Routledge; 2005. pp. 109-133.
- 491 7. Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, et al. Ancient admixture in hu-
492 man history. *Genetics*. 2012;192(3): 1065-1093.
- 493 8. Narasimhan VM, Patterson N, Moorjani P, Rohland N, Bernardos R, Mallick S, et al. The for-
494 mation of human populations in South and Central Asia. *Science*. 2019;365(6457): eaat7487.
- 495 9. Tätte K, Pagani L, Pathak AK, Köks S, Ho Duy B, Ho XD, et al. The genetic legacy of conti-
496 nental scale admixture in Indian Austroasiatic speakers. *Sci Rep*. 2019;9(1): 3818.
- 497 10. Reich D, Patterson N, Campbell D, Tandon A, Mazieres S, Ray N, et al. Reconstructing native
498 American population history. *Nature*. 2012;488(7411): 370-4.
- 499 11. Haak W, Lazaridis I, Patterson N, Rohland N, Mallick S, Llamas B, et al. Massive migration
500 from the steppe was a source for Indo-European languages in Europe. *Nature*. 2015;522(7555):
501 207-11.
- 502 12. Harney É, Patterson N, Reich D, Wakeley J. Assessing the performance of qpAdm: A statistical
503 tool for studying population admixture. *BioRxiv [Preprint]*. 2020 bioRxiv 2020.04.09.032664

- 504 [posted 2020 April 10]. Available from:
505 <https://www.biorxiv.org/content/10.1101/2020.04.09.032664v1> doi:
506 <https://doi.org/10.1101/2020.04.09.032664>
- 507 13. Ning C, Fernandes D, Changmai P, Flegontova O, Yuncu E, Maier R, et al. The genomic for-
508 mation of first American ancestors in East and Northeast Asia. *BioRxiv* [Preprint]. 2020
509 bioRxiv 2020.10.12.336628 [posted 2020 October 12]. Available from:
510 <https://www.biorxiv.org/content/10.1101/2020.10.12.336628v1>
511 doi: <https://doi.org/10.1101/2020.10.12.336628>
- 512 14. Flegontov P, Altınışık NE, Changmai P, Rohland N, Mallick S, Adamski N, et al. Palaeo-
513 Eskimo genetic ancestry and the peopling of Chukotka and North America. *Nature*.
514 2019;570(7760): 236–40.
- 515 15. Skoglund P, Posth C, Sirak K, Spriggs M, Valentin F, Bedford S, et al. Genomic insights into
516 the peopling of the Southwest Pacific. *Nature*. 2016;538(7626): 510–3.
- 517 16. Liu D, Duong NT, Ton ND, Van Phong N, Pakendorf B, Van Hai N, et al. Extensive ethnolin-
518 guistic diversity in Vietnam reflects multiple sources of genetic diversity. *Mol Biol Evol*.
519 2020;37(9): 2503–19.
- 520 17. Gray RD, Drummond AJ, Greenhill SJ. Language phylogenies reveal expansion pulses and
521 pauses in Pacific settlement. *Science*. 2009;323(5913): 479-83
- 522 18. Wang C-C, Yeh H-Y, Popov AN, Zhang H-Q, Matsumura H, Sirak K, et al. The Genomic for-
523 mation of human populations in East Asia. *BioRxiv* [Preprint]. 2020
524 bioRxiv 2020.03.25.004606 [posted 2020 March 25]. Available from:
525 <https://www.biorxiv.org/content/10.1101/2020.03.25.004606v1>
526 doi: <https://doi.org/10.1101/2020.03.25.004606>
- 527 19. Lazaridis I, Patterson N, Mittnik A, Renaud G, Mallick S, Kirsanow K, et al. Ancient human
528 genomes suggest three ancestral populations for present-day Europeans. *Nature*.

- 529 2014;513(7518): 409–13.
- 530 20. Vongpaisarnsin K, Listman JB, Malison RT, Gelernter J. Ancestry informative markers for dis-
531 tinguishing between Thai populations based on genome-wide association datasets. *Legal Medi-*
532 *cine*. 2015;17(4): 245-50.
- 533 21. Mörseburg A, Pagani L, Ricaut FX, Yngvadottir B, Harney E, Castillo C, et al. Multi-layered
534 population structure in Island Southeast Asians. *European Journal of Human Genetics*.
535 2016;24(11): 1605-11.
- 536 22. Kutanan W, Liu D, Kampuansai J, Srikummool M, Srithawong S, Shoocongdej R, et al. Recon-
537 structing the human genetic history of mainland Southeast Asia: insights from genome-wide da-
538 ta from Thailand and Laos. *BioRxiv [Preprint]*. 2020bioRxiv 2020.12.24.424294 [posted 2020
539 December 24]. Available from: <https://www.biorxiv.org/content/10.1101/2020.12.24.424294v1>
540 doi: <https://doi.org/10.1101/2020.12.24.424294>
- 541 23. Kutanan W, Ghirotto S, Bertorelle G, Srithawong S, Srithongdaeng K, Pontham N, et al. Geog-
542 raphy has more influence than language on maternal genetic structure of various northeastern
543 Thai ethnicities. *Journal of human genetics*. 2014;59(9): 512-20.
- 544 24. Macholdt E, Arias L, Duong NT, Ton ND, Van Phong N, Schröder R, et al. The paternal and
545 maternal genetic history of Vietnamese populations. *Eur J Hum Genet*. 2020;28(5): 636–45.
- 546 25. He JD, Peng MS, Quang HH, Dang KP, Trieu AV, Wu SF, et al. Patrilineal perspective on the
547 Austronesian diffusion in Mainland Southeast Asia. *PLoS One*. 2012;7(5): e36437.
- 548 26. Kunstadter P. Subsistence agricultural economies of Lua’ and Karen hill farmers, Mae Sariang
549 District, Northwestern Thailand. In: Kunstadter P, Chapman EC, Sabhasri S, editors. *Farmers in*
550 *the forest*. University of Hawai’i Press; 1978. pp. 118–208.
- 551 27. Kauffmann HE. Some social and religious institutions of the Lawā (Northwest Thailand). *J Si-*
552 *am Soc*. 1972;60(1): 237–306.
- 553 28. Nahhas RW. Sociolinguistic Survey of Lawa in Thailand. *SIL Electron Surv Reports*. 2011;222.

- 554 29. Kutanan W, Kampuansai J, Srikumool M, Brunelli A, Ghirotto S, Arias L, et al. Contrasting
555 paternal and maternal genetic histories of Thai and Lao populations. *Mol Biol Evol.* 2019 Jul
556 1;36(7): 1490–506.
- 557 30. Black ML, Dufall K, Wise C, Sullivan S, Bittles AH. Genetic ancestries in northwest Cambodia.
558 *Ann Hum Biol.* 2006;33(5–6): 620–7.
- 559 31. Hla NP. The Major Role of the Mons in Southeast Asia. *J Siam Soc.* 1992;79: 13-21.
- 560 32. Schliesinger J. Ethnic groups of Thailand: Non-Tai-speaking peoples. White Lotus Press; 2000.
- 561 33. Baba Y. The Meaning of ‘History’ or ‘Past’ in the Context of the Tai-Lue Cultural Revival
562 Movement. In: 13th International Conference on Thai Studies “Globalized Thailand? Connec-
563 tivity, Conflict and Conundrums of Thai Studies”. 2017. pp. 58–67.
- 564 34. Li JZ, Absher DM, Tang H, Southwick AM, Casto AM, Ramachandran S, et al. Worldwide
565 Human Relationships Inferred from Genome-Wide Patterns of Variation. *Science.*
566 2008;319(5866): 1100–4.
- 567 35. Rae I, Witzel M. The overseas Chinese of South East Asia: history, culture, business. Palgrave
568 Macmillan UK; 2008.
- 569 36. West BA. Encyclopedia of the peoples of Asia and Oceania. Facts on File, Incorporated; 2009
- 570 37. Kampuansai J, Bertorelle G, Castri L, Nakbunlung S, Seielstad M, Kangwanpong D. Mitochon-
571 drial DNA variation of Tai speaking peoples in Northern Thailand. *Sci Asia.* 2007;33: 443-8.
- 572 38. HUGO Pan-Asian SNP Consortium. Mapping human genetic diversity in Asia. *Science.*
573 2009;326(5959): 1541-5.
- 574 39. Besaggio D, Fuselli S, Srikumool M, Kampuansai J, Castrì L, Tyler-Smith C, et al. Genetic
575 variation in northern Thailand hill tribes: origins and relationships with social structure and lin-
576 guistic differences. *BMC Evolutionary Biology.* 2007;7(S2): S12.
- 577 40. Lithanatudom P, Wipasa J, Inti P, Chawansuntati K, Svasti S, Fucharoen S, et al. Hemoglobin E
578 prevalence among ethnic groups residing in malaria-endemic areas of Northern Thailand and its

579 lack of association with *Plasmodium falciparum* invasion in vitro. PLoS One. 2016;11(1):
580 e0148079–e0148079.

581 41. Kutanan W, Kumpansai J, Changmai P, Flegontov P, Schröder R, Macholdt E, et al. Con-
582 trasting maternal and paternal genetic variation of hunter-gatherer groups in Thailand. Scientific
583 reports. 2018;8(1): 1-9.

584 42. Flegontov P, Changmai P, Zidkova A, Logacheva MD, Altınışık NE, Flegontova O, et al. Ge-
585 nomic study of the Ket: a Paleo-Eskimo-related ethnic group with significant ancient North
586 Eurasian ancestry. Scientific reports. 2016;6(1): 1-2.

587 43. Lazaridis I, Nadel D, Rollefson G, Merrett DC, Rohland N, Mallick S, et al. Genomic insights
588 into the origin of farming in the ancient Near East. Nature. 2016;536(7617):419-24.

589

590 **Fig 1. Locations of populations for whom genome-wide data was generated in this study.** Colors
591 represent language families: pink, Sino-Tibetan; green, Hmong-Mien; red, Austroasiatic; and purple,
592 Kra-Dai. The map was created using R package “rworldmap” ([https://cran.r-](https://cran.r-project.org/web/packages/rworldmap/)
593 [project.org/web/packages/rworldmap/](https://cran.r-project.org/web/packages/rworldmap/)).

594 **Fig 2. A principal component analysis (PCA) plot of present-day Eurasian populations.** PCA was
595 performed using PLINK. **Left panel:** An overview of the PC1 vs. PC2 space for all populations. The
596 legend at the bottom of the plot lists abbreviations of meta-populations: CAS, Central Asians; ESEA,
597 East and Southeast Asians; NEGM = Mainland Negritos; SAS, South Asians; EUR, Europeans; Munda,
598 Austroasiatic-speaking populations (the Munda branch) from India; Onge, Onge (Andamanese hunter-
599 gatherers); and SIB, Siberians. **Right panel:** A zoomed in view on the rectangle in the left panel.

600 **Fig 3. Results of an ADMIXTURE analysis.** The plot represents results for 12 hypothetical ancestral
601 populations. Abbreviations of meta-populations are shown above the plot: AFR, Africans; EUR,
602 Europeans; CAU, Caucasians; PA, Papuans and Australians; Onge, Onge (Andamanese hunter-

603 gatherers), SAS, South Asians; Munda, Austroasiatic-speaking populations (the Munda branch) from
604 India; SAM, Native South Americans; SIB, Siberian; CAS, Central Asians; ESEA, East and Southeast
605 Asians; and NEGM, Mainland Negritos.

606 **Fig 4. A biplot showing results of outgroup f_3 -tests.** The biplot of $f_3(\text{Mbuti}; \text{Brahmin Tiwari}, X)$ vs.
607 $f_3(\text{Mbuti}; \text{Han}, X)$ illustrates the amount of genetic drift shared between test ESEA populations and
608 Brahmin Tiwari or Han. The trend line represents a ratio of shared genetic drifts that is common for
609 most ESEA populations. The positions of few ESEA populations deviated from the trend line, which
610 indicates elevated shared drift between the Indian reference population and the test population, as
611 compared to most ESEA populations.

612 **Fig 5. An overview of admixture proportions estimated by qpAdm.** Admixture proportions were
613 inferred using qpAdm with three groups of surrogates representing three ancestries: deeply diverged
614 East Eurasian (NEG), South Asian (SAS), and East Asian (EA). Admixture proportions were averaged
615 across all models which passed our criteria for "fitting" models. The map was plotted using R package
616 "rnaturalearth" (<https://github.com/ropensci/rnaturalearth>).

617 **Fig 6. Skeleton graphs used for the qpGraph mapping method.** We used the skeleton graphs to
618 explore the genetic make-up of ESEA populations. We used different South Indian populations for two
619 skeleton graphs: Coorghi in panel **A** and Palliyar in panel **B**.

620

621 **Supporting information**

622 **S1 Fig. A biplot of $f_3(\text{Mbuti}; \text{Coorghi}, X)$ vs. $f_3(\text{Mbuti}; \text{Dai}, X)$ (A), $f_3(\text{Mbuti}; \text{Coorghi}, X)$ vs. $f_3(\text{Mbuti};$
623 $\text{Han}, X)$ (B), and $f_3(\text{Mbuti}; \text{Brahmin Tiwari}, X)$ vs. $f_3(\text{Mbuti}; \text{Dai}, X)$ (C).**

624 **S2 Fig. Skeleton graphs used for qpGraph mapping, with edges numbered.** Coorghi was used as an
625 Indian surrogate for skeleton graph **A** and Palliyar for skeleton graph **B**.

- 626 **S3 Fig. The starting skeleton graph with 5 populations.**
- 627 **S4 Fig. The best-fitting graph for Onge mapped on the 5-population skeleton graph (Fig S3).**
- 628 **S5 Fig. The best-fitting graph for Dai mapped on the 6-population skeleton graph (Fig S4).**
- 629 **S6 Fig. The best-fitting graph for an ancient Iranian herder from Ganj Dareh mapped on the 7-**
630 **population skeleton graph (Fig S5).**
- 631 **S7 Fig. The best-fitting graph for Mlabri mapped on the 8-population skeleton graph (Fig S6).**
- 632 **S8 Fig. The best-fitting graphs for Coorghhi (A) and Palliyar (B) mapped on the 9-population**
633 **skeleton graph (Fig S7).**
- 634
- 635 **S1 Table. Information on reference populations used in this study.**
- 636 **S2 Table. qpWave and qpAdm results.**
- 637 **S3 Table. All best-fitting qpGraph models.**
- 638 **S4 Table. Metadata for newly genotyped present-day individuals.**
- 639

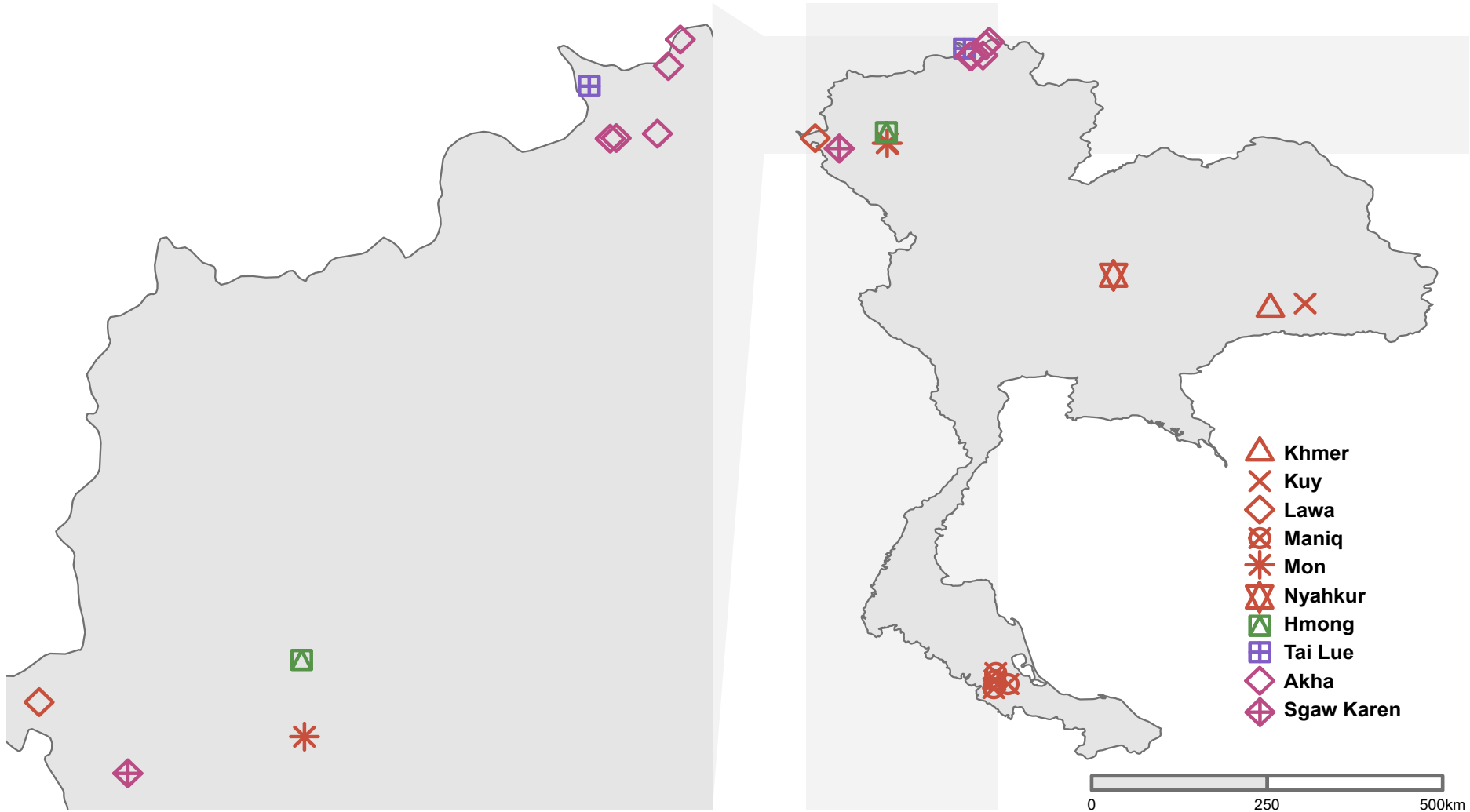


Fig.1

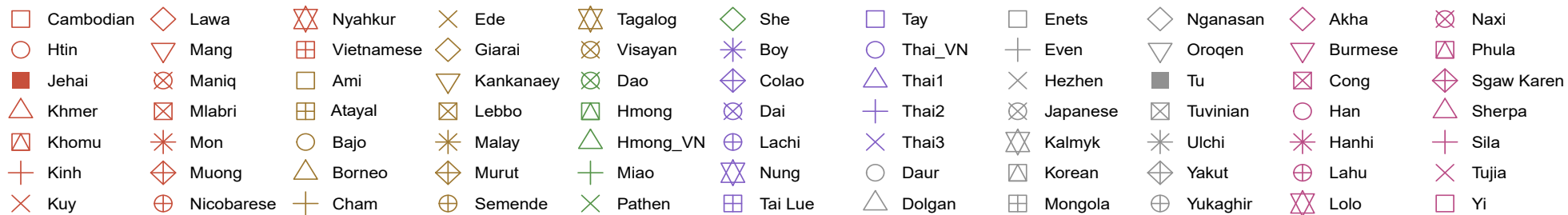
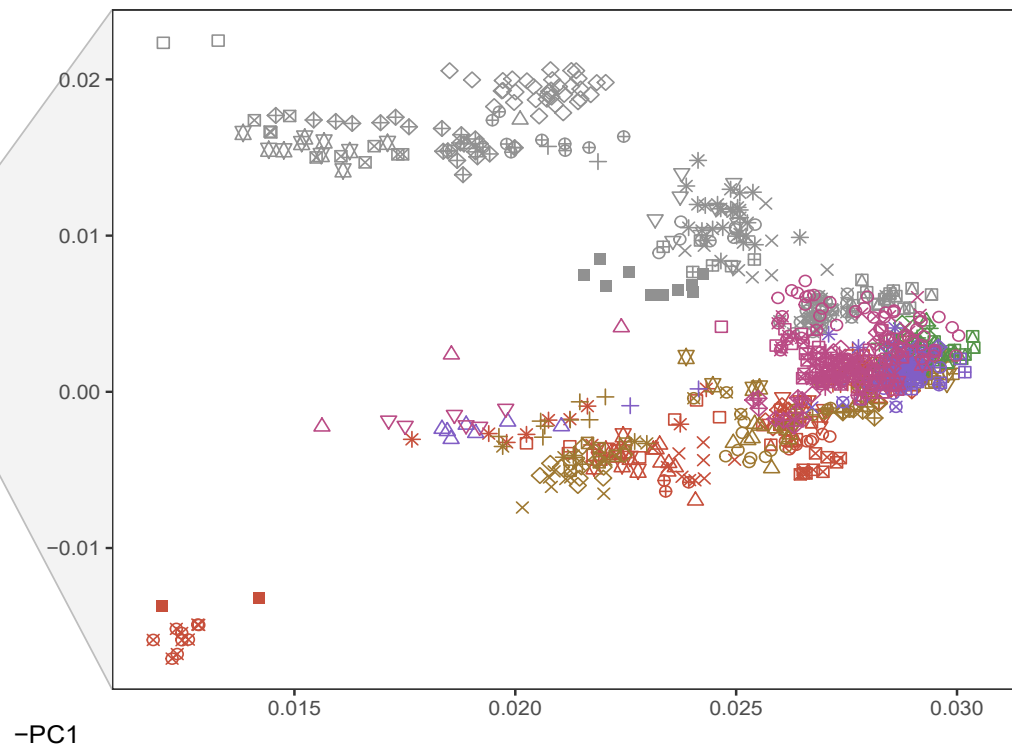
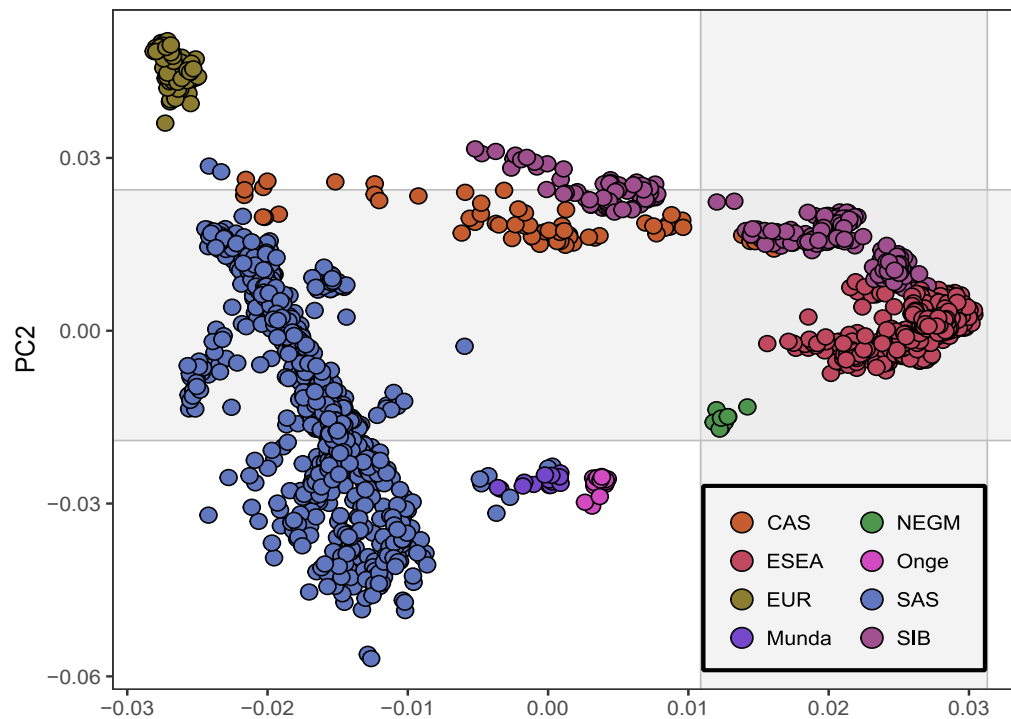


Fig.2

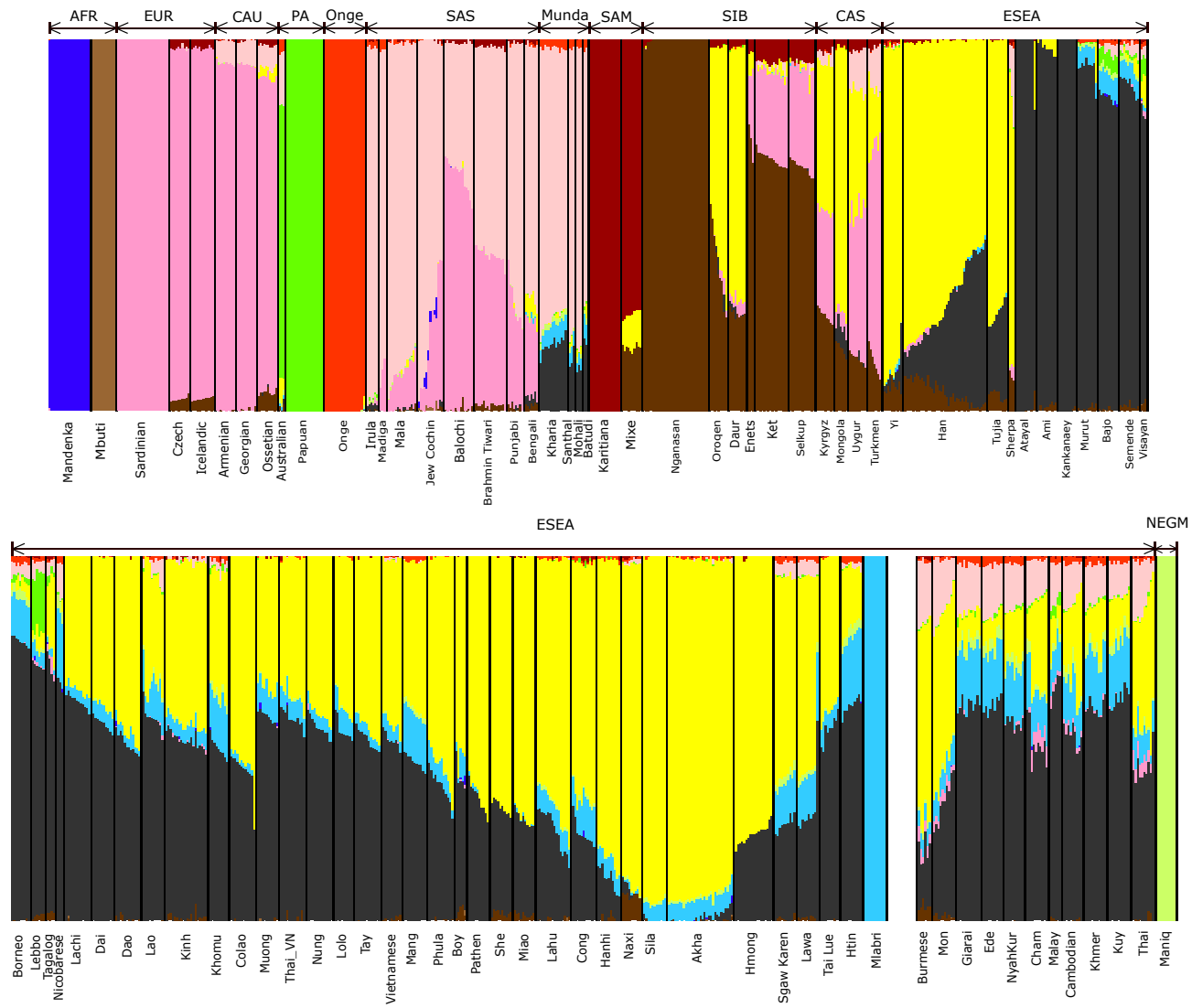


Fig.3

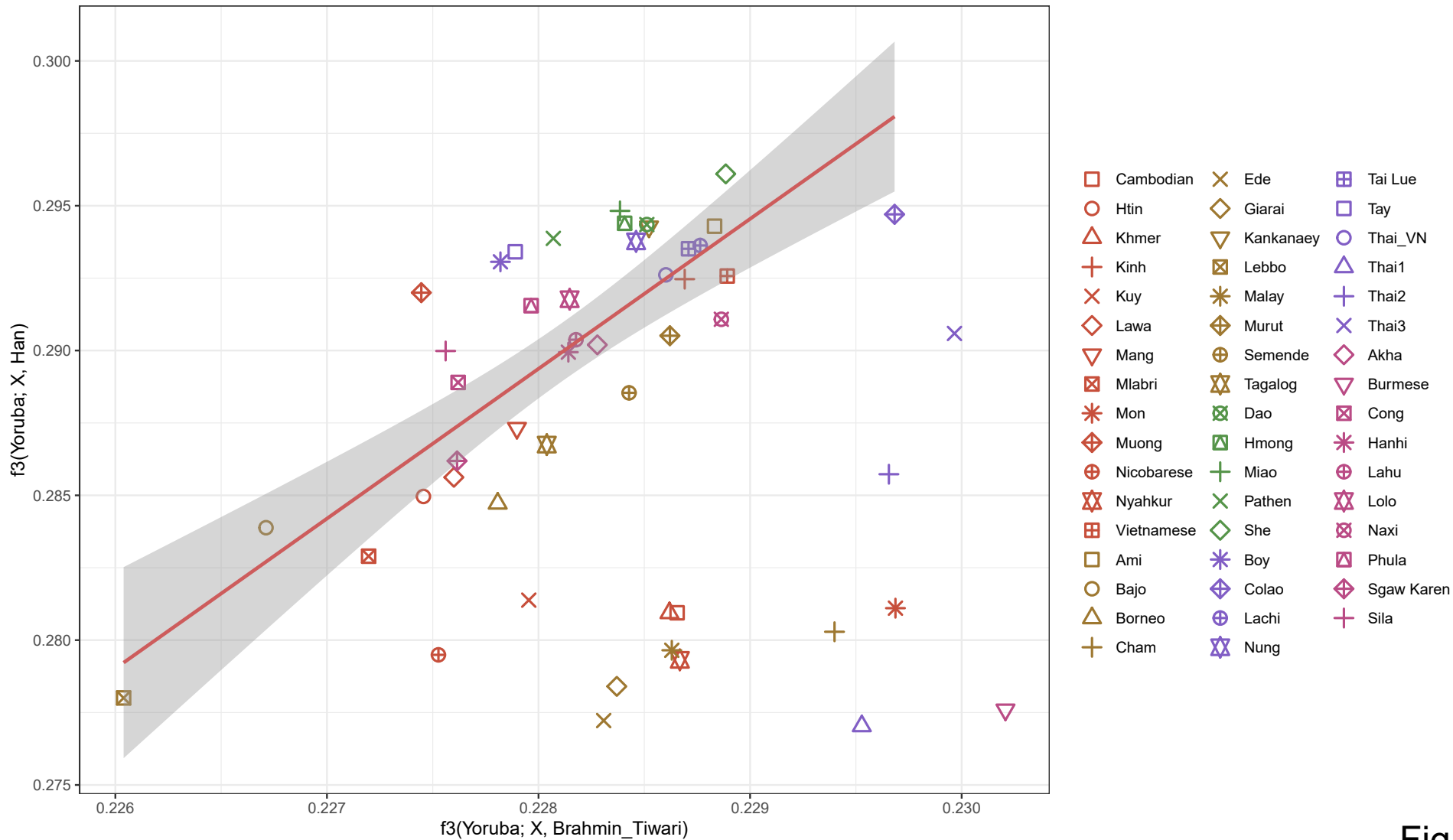


Fig.4

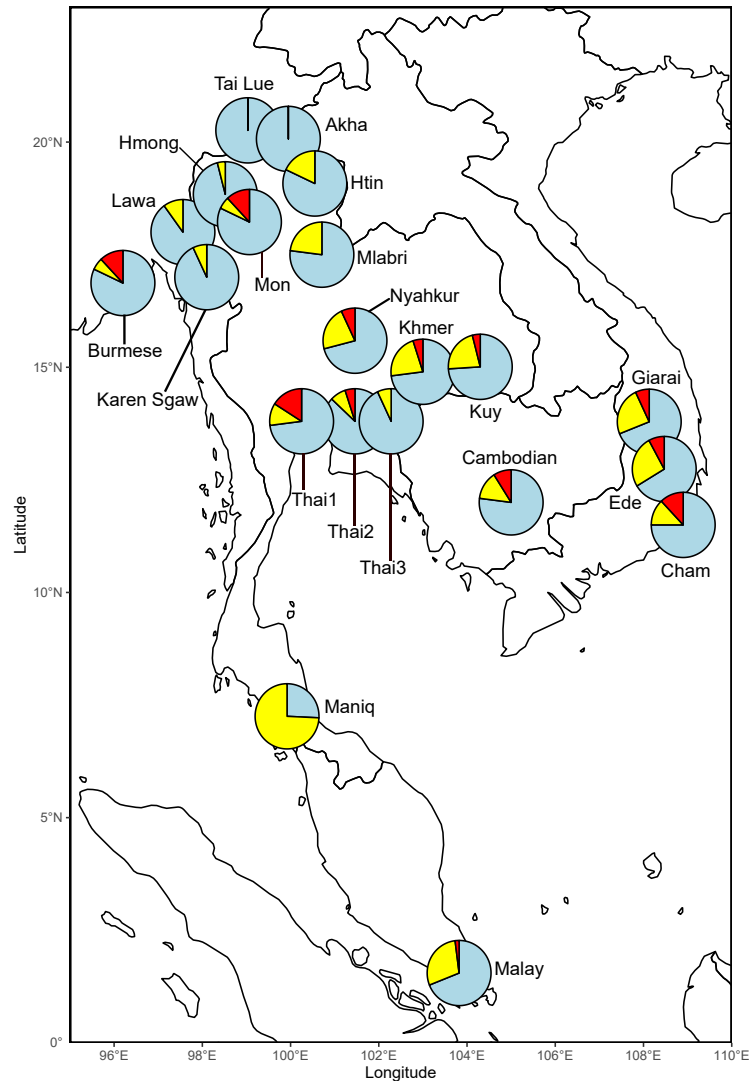


Fig.5

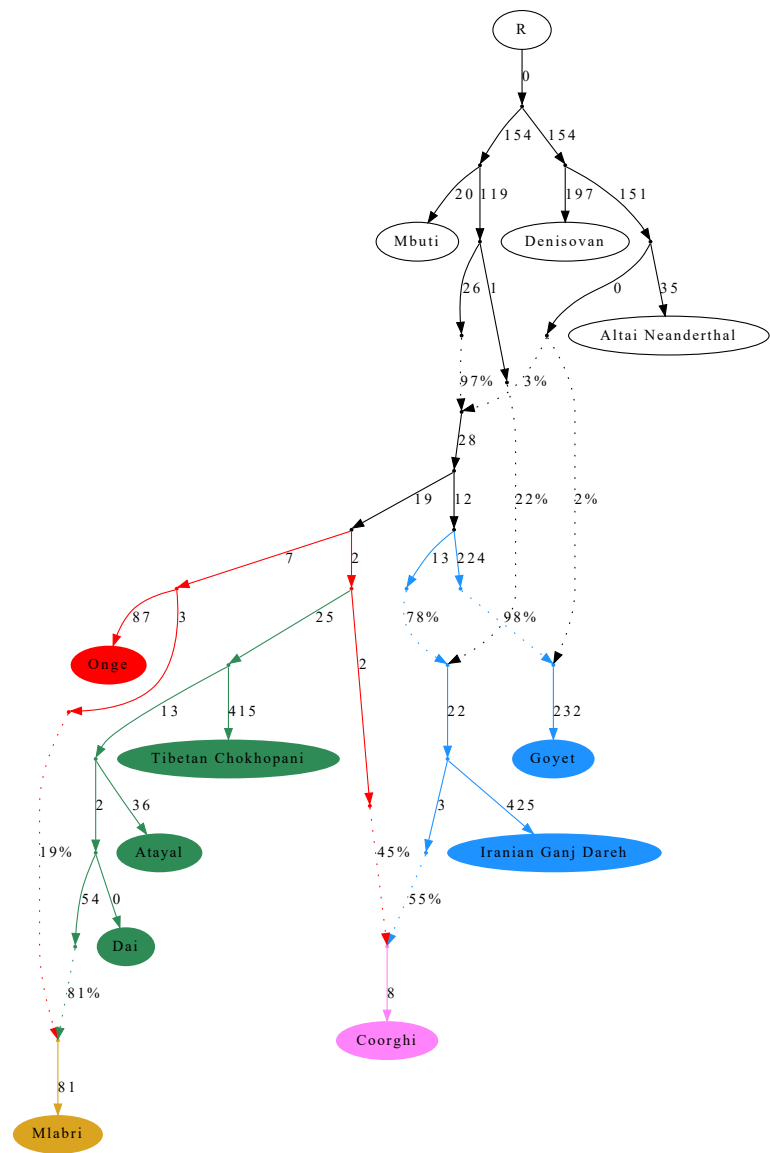
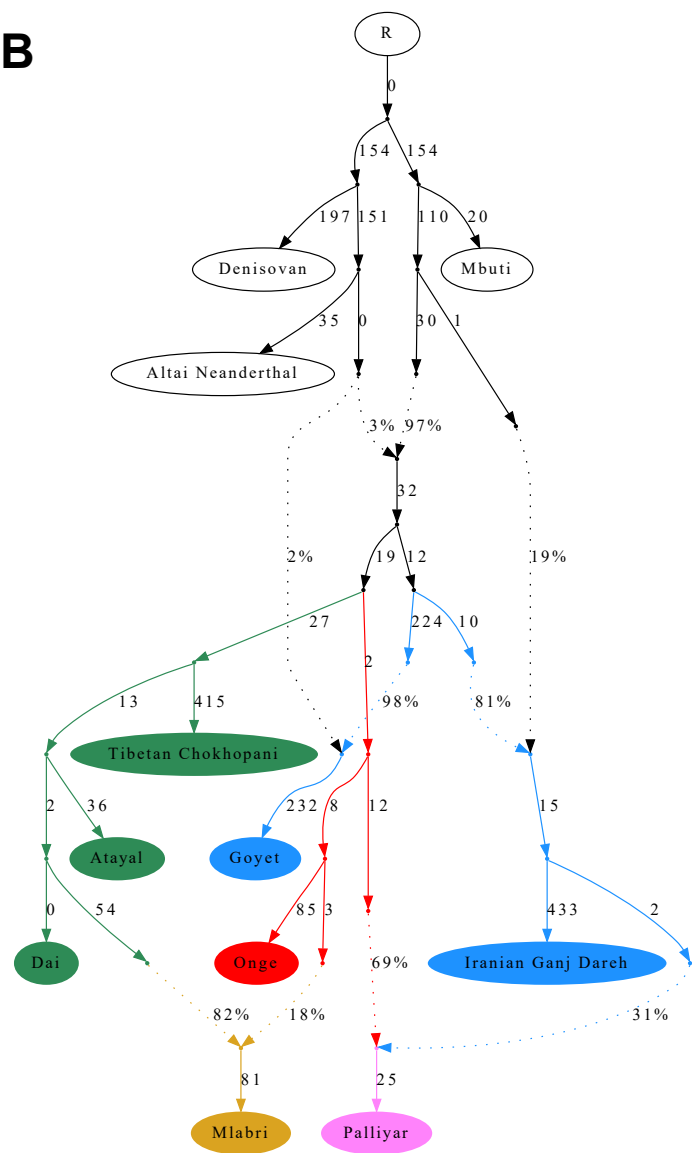
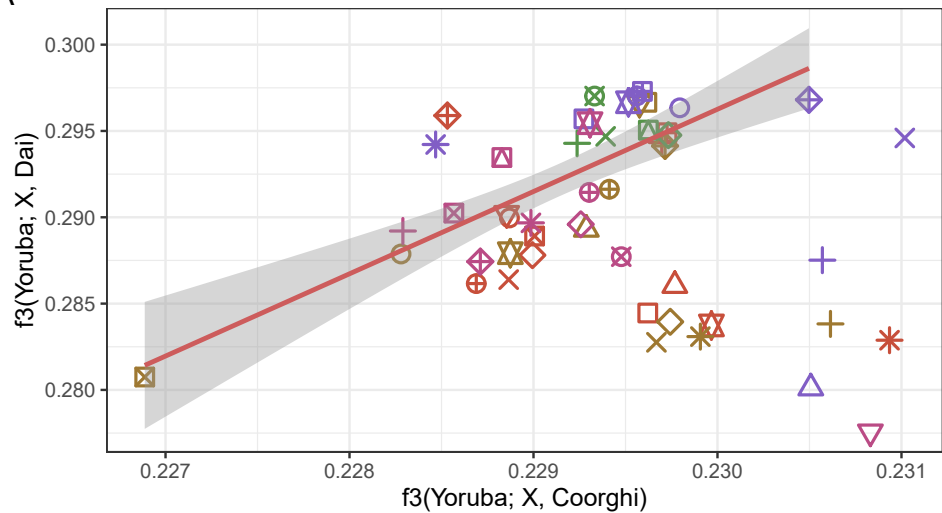
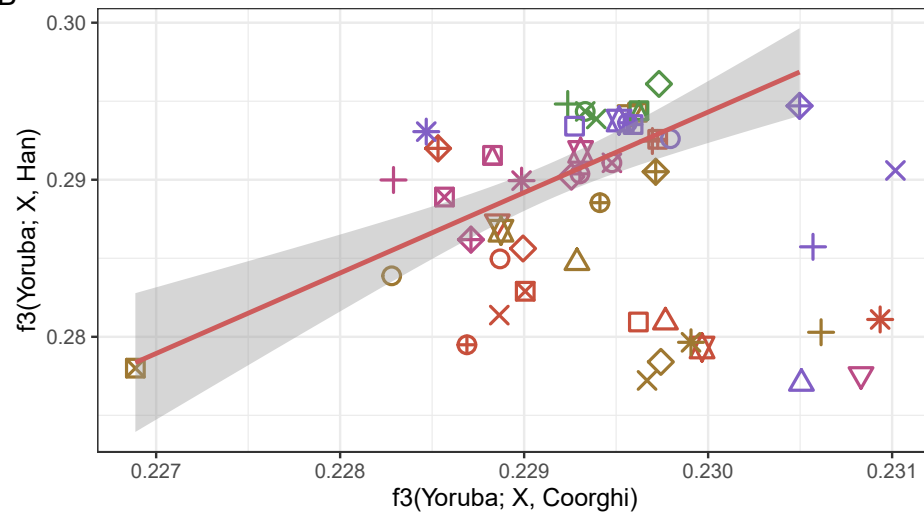
A**B**

Fig.6

A



B



C

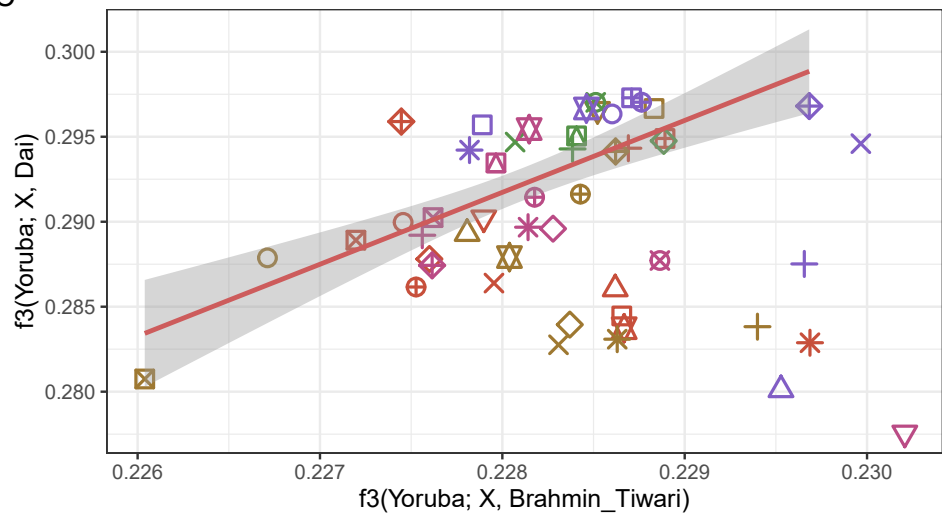


Fig.S1

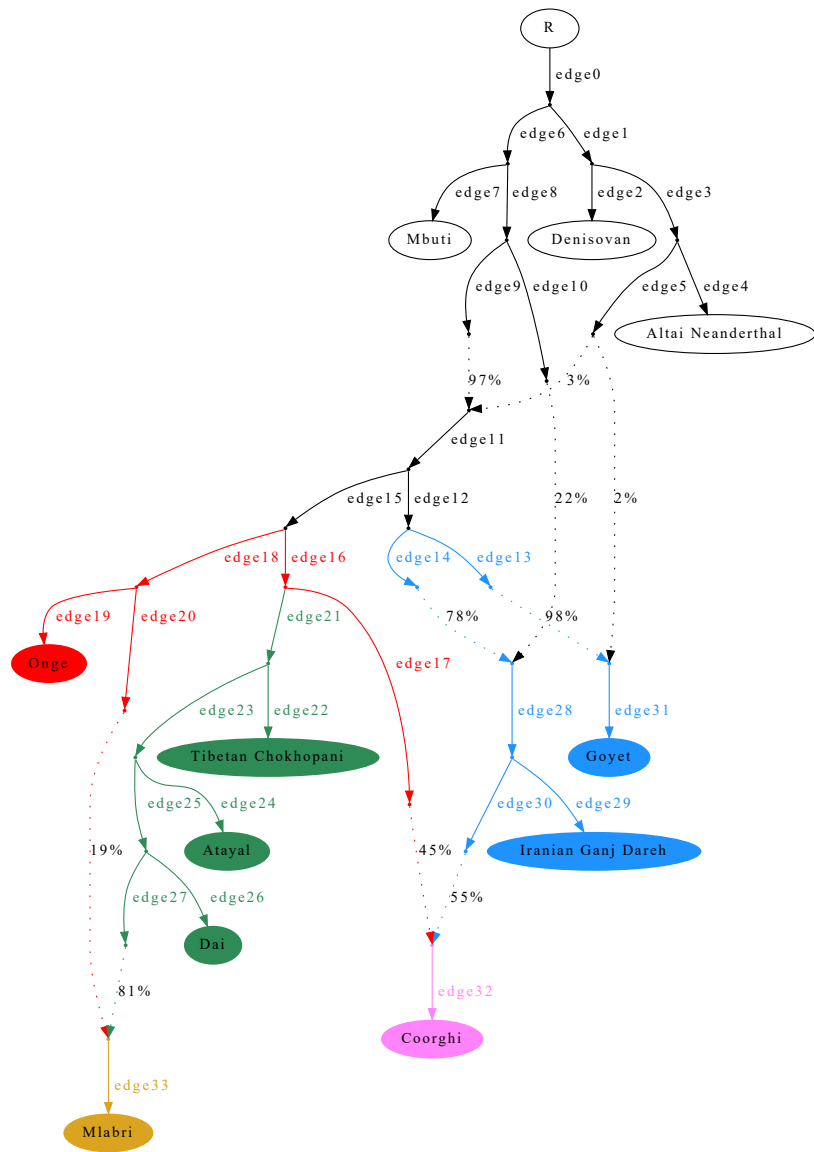
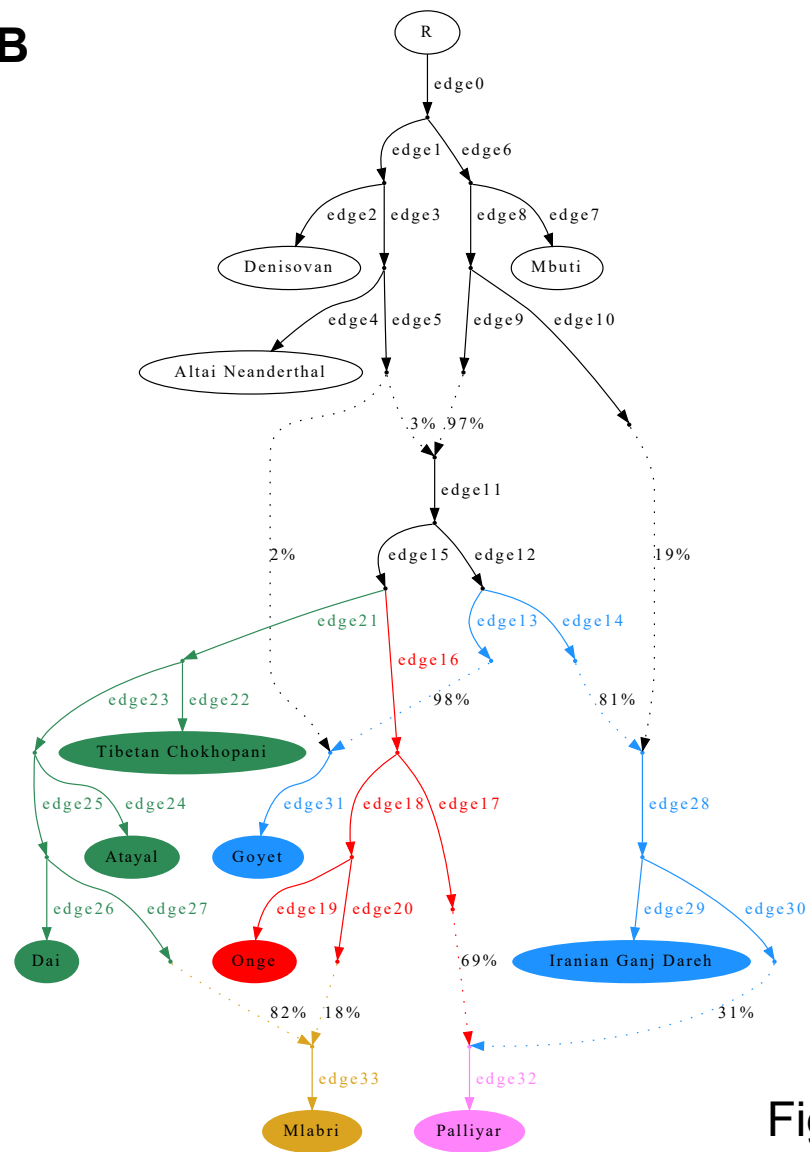
A**B**

Fig.S2

template_first4popNG2_lock3 :: Alt Den Bel Ata 0.002448 0.000018 -0.002429 0.001139 -2.134

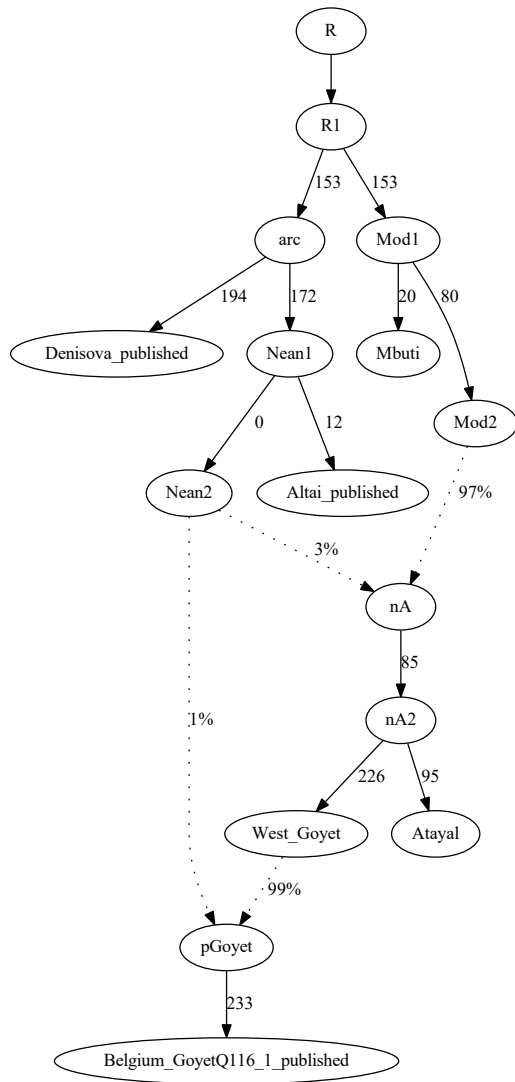


Fig.S3

template_first5popNG2_lock3_1way_nA2_Atayal_test_Onge :: Alt Den Bel Ata 0.002379 0.000018 -0.002360 0.001132 -2.086

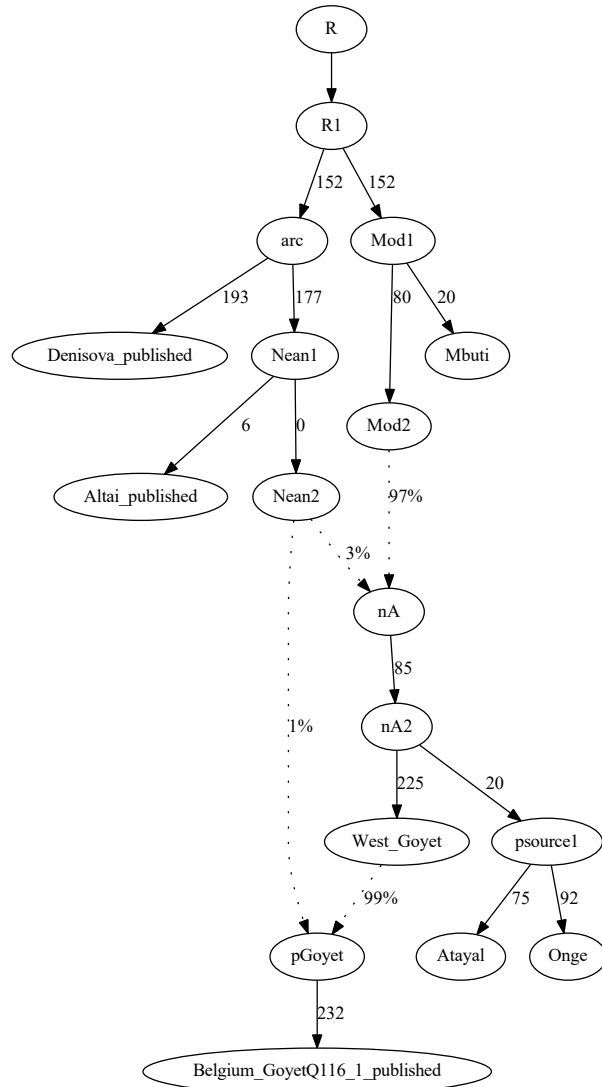


Fig.S4

template_first6popNG2_lock3_1way_East_Atayal_test_Dai:: Alt Den Bel Ata 0.002393 0.000018 -0.002374 0.001128 -2.104

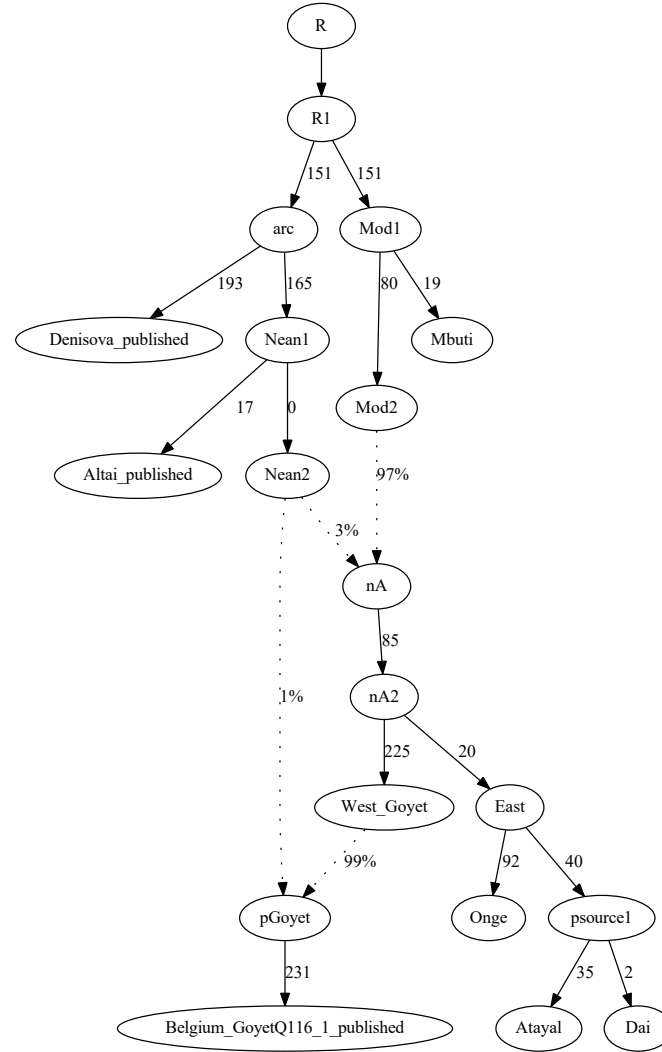


Fig.S5

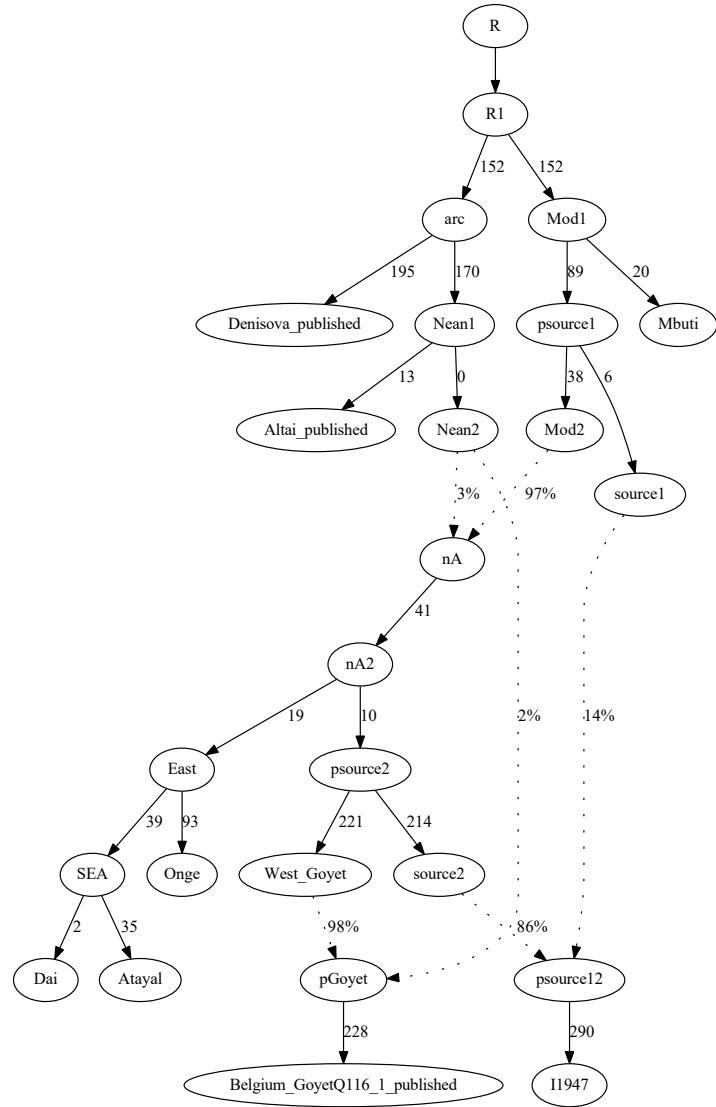


Fig.S6

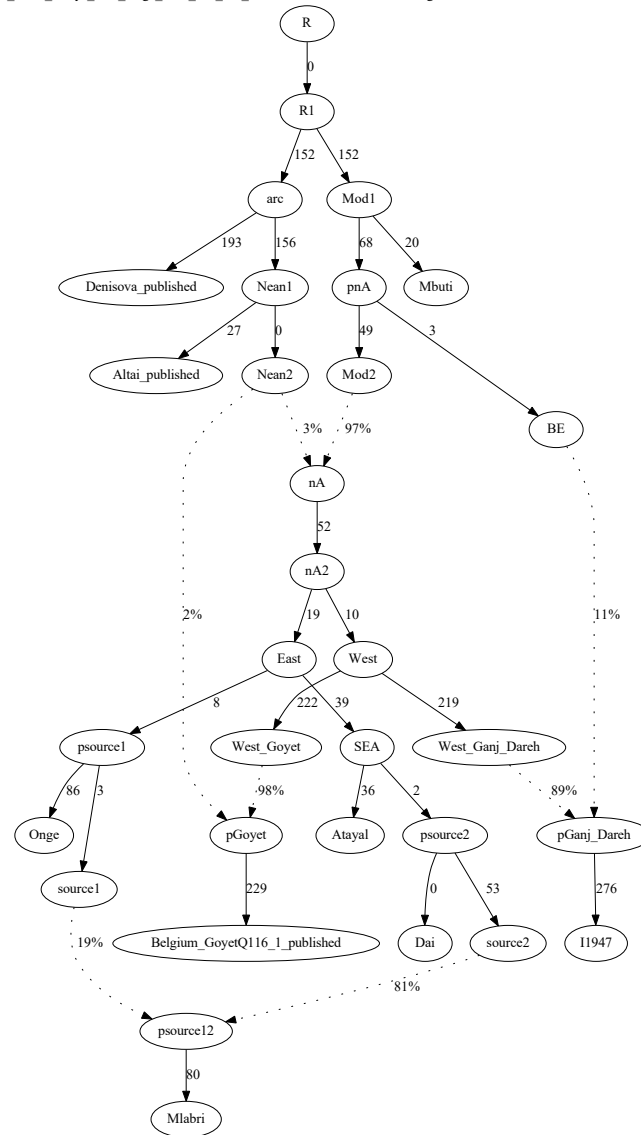
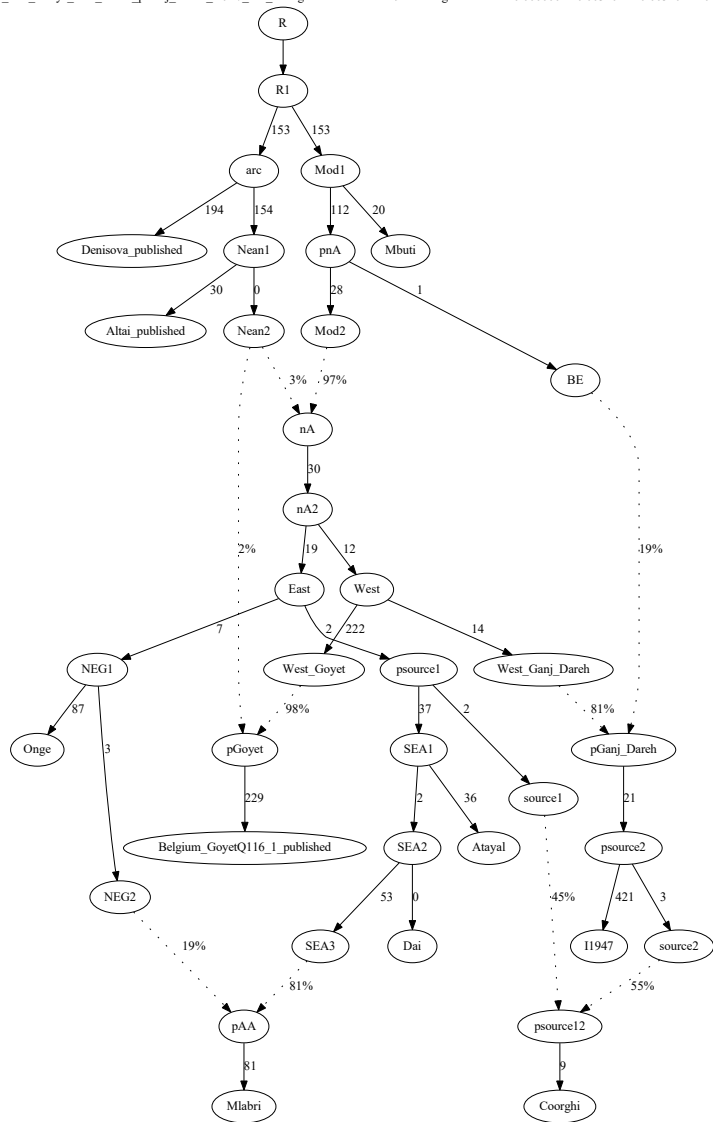


Fig.S7

A



B

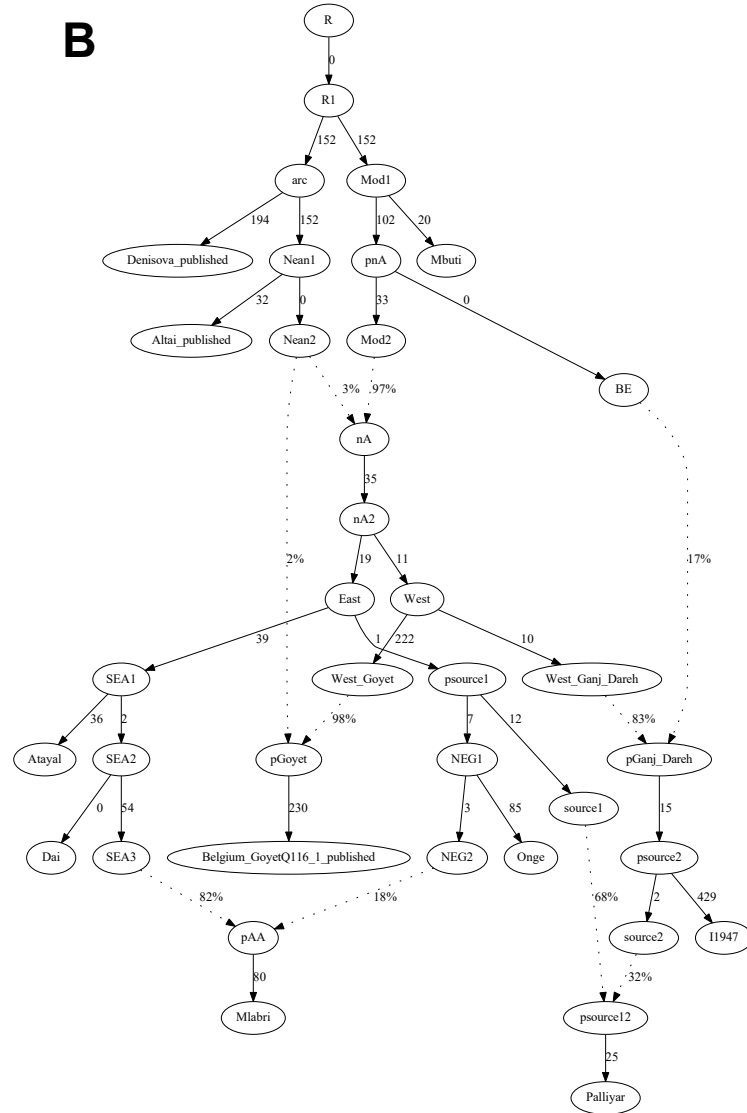


Fig.S8