

1 Genomic analysis of Andamanese provides insights
2 into ancient human migration into Asia and
3 adaptation

4
5 **Authors:** Mayukh Mondal¹⁺, Ferran Casals²⁺, Tina Xu³⁺, Giovanni M. Dall'Olio⁴, Marc Pybus¹, Mihai G.
6 Netea⁵, David Comas¹, Hafid Laayouni^{1,6}, Qibin Li^{3*}, Partha P. Majumder^{7*}, Jaume Bertranpetit^{1,8*}

7
8 **Affiliations:**

9 ¹ Institut de Biologia Evolutiva (UPF-CSIC), Universitat Pompeu Fabra, Barcelona, Catalonia, Spain

10 ² Servei de Genòmica, Universitat Pompeu Fabra, Barcelona, Catalonia, Spain

11 ³ BGI Shenzhen, Yantian District, Shenzhen, 518083, China

12 ⁴ Computational Biology, Target Sciences, GSK R&D, GlaxoSmithKline, Stevenage, Hertfordshire, United
13 Kingdom

14 ⁵ Department of Internal Medicine, Radboud University Medical Center, Nijmegen, The Netherlands

15 ⁶ Departament de Genètica i de Microbiologia, Universitat Autònoma de Barcelona, Bellaterra, Catalonia,
16 Spain

17 ⁷ National Institute of BioMedical Genomics, Kalyani, West Bengal 741251, India

18 ⁸ Leverhulme Centre for Human Evolutionary Studies, Department of Archaeology and Anthropology,
19 University of Cambridge, Cambridge, United Kingdom.

20

21 + Co-first authorship

22 * Co-senior authorship

23 Correspondence: jaume.bertranpetit@upf.edu or ppm1@nibmg.ac.in

24 **Abstract:** To shed light on the peopling of South Asia and the origins of the morphological adaptations
25 found there, we analyzed whole-genome sequences from ten Andamanese individuals and compared them
26 with 60 individuals from mainland Indian populations with different ethnic histories, and with publicly-
27 available data from other populations. We show that all Asian and Pacific populations share a single origin
28 and expansion out of Africa, contradicting an earlier proposal of two independent waves¹⁻⁴. We also show
29 that populations from South and Southeast Asia harbor a small proportion of ancestry from an unknown
30 extinct hominin, which is absent from Europeans and East Asians. The footprints of adaptive selection in the
31 genomes of the Andamanese show that their characteristic distinctive phenotypes (including very short
32 stature) do not reflect an ancient African origin, but instead result from strong natural selection on genes
33 related to human body size.

34

35 **Main Text:**

36 The origin of the Andamanese people (Andaman Islands, Bay of Bengal, India) has been considered to be
37 different from other Asian populations, because of their very distinctive so-called ‘Negrito’ morphology, and
38 the unclassifiable language that they speak⁵⁻⁷. It has been suggested that they are a living relic of a first Out-
39 of-Africa (OOA) wave of modern humans using the southern exit route, who did not subsequently mix with
40 other populations^{1,2} (since there have been multiple OOA events in human evolution, here ‘OOA’ refers to
41 the Out-of-Africa event(s) for fully modern humans only). A common origin of the Andaman (and other)
42 ‘Negrito’ populations, Melanesians and Australians, was initially proposed based on morphological
43 characteristics^{1,2} and subsequently supported by some genetic studies⁴. Previous analysis of genome-wide
44 genotyping data from several Indian populations showed that the Andamanese are one of two main reference
45 populations for estimating ancestries of Indian populations⁸. However, the lack of whole-genome sequence
46 data from the Andamanese has limited understanding of both their ancestry and the specificity of the
47 adaptations that may have resulted in their distinctive morphological features. Whether their distinctive
48 ‘Negrito’ morphological features (small body size, dark skin, curly hair, etc.) are ancestral or derived may
49 potentially be inferred by analyzing footprints of selection in their genomes. It matches known adaptations
50 due to insularity in many groups of large animals, which may explain their fast evolution in body size, a
51 feature that is shared by some extinct hominin populations⁹ as well as present-day humans¹⁰.

52 Seventy individuals from India were sequenced at ~15x coverage (Supplementary Note), including 60
53 individuals from mainland India and 10 from the Jarawa (JAR) and Onge (ONG) populations in the

54 Andaman islands (Supplementary Table 1, Supplementary Figure 1). The demographically small and
55 historically isolated Andamanese population show higher relatedness among individuals as well as higher
56 inbreeding coefficients and longer runs of homozygosity than all continental Indian populations examined
57 (Supplementary Figure 2, 3 and 4). In agreement with previous studies^{8,11}, Principal Component Analysis
58 (PCA) showed that the Andamanese constitute a genetically distinct cluster compared with the mainland
59 Indian populations (Supplementary Figure 5). Interestingly, the Jarawa and the Onge cluster tightly together,
60 indicative of their genomic homogeneity, and show a lack of recent admixture (Figure 1a), which is known
61 to have taken place in Andaman during the last century¹², but did not affect the individuals sampled.

62 Using several approaches, we investigated whether the Andamanese were descendants of the same OOA
63 event that resulted in the peopling of mainland India, or whether some part of their origins can be traced to
64 an earlier and independent OOA wave, as has been proposed for Aboriginal Australians⁴. First, the D-
65 statistic (Dstat) analysis¹³ (Supplementary Figure 6) showed that Andamanese share more alleles with all
66 OOA populations than with sub-Saharan Africans, suggesting that Andamanese shared a common and
67 similar ancestry with all other OOA populations. Second, TreeMix analysis¹⁴ also supports Africans as an
68 outgroup to all OOA populations (Figure 1b), with a closer relationship of Andamanese with Asians and
69 continental Indians than with Pacific populations. Third, relative cross coalescent analysis by MSMC¹⁵
70 displayed a much earlier split for Andamanese and Africans than for Andamanese and any other OOA
71 population, which are themselves very similar (Figure 1c). Estimation of historical effective population sizes
72 by MSMC suggests a similar bottleneck event for Andamanese and all other OOA populations at around
73 50,000 years ago (Supplementary Figure 7). All of these results suggest that the Andamanese shared a
74 common ancestry with all the other OOA populations, indicative of a commonality of all Asian and Pacific
75 populations and consistent with a single main OOA migration.

76 Dstat analysis (Supplementary Figure 8) revealed that the Andamanese shared more alleles with East Asian,
77 Papuan, and mainland Indian tribal populations than with Europeans, indicating that Europeans are an out-
78 group for all Asian populations. Both TreeMix (Figure 1b) and Dstat outgroup analysis (Supplementary
79 Table 2) supported this inference. Relative cross-coalescent analysis (Figure 1c) also showed a similar result:
80 the separation between Andamanese and Europeans predates the separation of Andamanese from Asians.
81 Analysis using available ancient European genome sequences from La Braña, Loschbour, and Stuttgart¹⁶⁻¹⁸
82 supported our results (Supplementary Figure 8-10 and Supplementary Table 3), showing Europeans as the
83 most distinct branch of all Eurasian and Pacific populations, even when considering the extinct Basal

84 Eurasian component of Europeans^{18,19}. Mitochondrial DNA analysis also supports a single origin for Asian
85 populations (Supplementary Table 4).

86 The analysis of the contribution of extinct hominin populations to the current genetic pool also suggests a
87 single origin for modern Asians, including Andamanese. Andamanese genomes have a similar amount of
88 Neanderthal^{13,20} introgression to other OOA populations (~2-4%), suggesting that the Neanderthal admixture
89 took place at a very early stage, before the OOA populations separated from each other (Supplementary
90 Figure 12). On the other hand, Papuans harbor a much higher proportion of Denisovan²¹ ancestry than any
91 other OOA population examined here (Supplementary Figure 13); all other Asian populations examined
92 (including the Andamanese) have only slightly more Denisovan ancestry than Europeans (Supplementary
93 Figure 14), as previously suggested²⁰. Besides that, no other difference in ancient contributions was observed
94 between the Andamanese and other Southern or Eastern Asian or Pacific populations.

95 We found that Andamanese, mainland Indian and Papuan populations carry ~2-3% fewer African alleles
96 than Europeans (Figure 2a) or East Asians (Figure 2b), as do Australians (similar yet higher value, see
97 below), a very intriguing result. We performed extensive simulations to show that this deficiency of African
98 alleles in the Andamanese cannot be explained by the Andamanese having low effective population size;
99 thus is not caused by private variants produced by specific mutations in their genome (no Admixture model,
100 Supplementary Table 5), or by later admixture between Europe or Asia and Africa (i.e. it cannot be due to a
101 “back to Africa” event; Supplementary Note and Supplementary Table 5), or by admixing with the initial
102 OOA modern humans settling in Eurasia. In contrast, it could be caused by mixture with a population that
103 diverged at least 300 kya (Supplementary Figure 15). In fact, an introgression from any hominin population
104 that can cause a bias in the Dstat calculations (Supplementary Note) would generate a false two-wave of
105 OOA (for modern humans) signal for the South Asian and Pacific populations, which is not observed. This
106 reduction in African ancestry for South Asian populations likewise cannot have originated from
107 Neanderthals or Denisovans, as these two populations have similar amounts of well-recognized ancestry in
108 Andamanese and East Asians. An alternative hypothesis is that this 2-3% reduction of African ancestry
109 originated from admixture with other hominin population(s) in Southeast Asia, such as *Homo erectus*²² or an
110 unknown extinct archaic population. A three-population model²³ confirms it (Supplementary Note and
111 Supplementary Figure 16). By calculating Dstat values for 50kb regions with a sliding window, we infer that
112 this unknown population diverged from Neanderthals and Denisova before they diverged from each other, as
113 seen initially by TreeMix (Supplementary Figure 17). To further identify specific DNA regions derived from

114 this hominin population, we implemented Sstar²⁴ on these putative fragments, and detected ~15Mb per
115 individual (average region length 65kb) from this hominin population that behaves either as a sister group to
116 Neanderthal and Denisova or even diverged earlier (Supplementary Figures 18 and 19). For Aboriginal
117 Australians, the deficit of African alleles is even higher (~6-7%; Figure 2), suggesting that this reduction
118 might be caused by admixture with some unknown ancient hominin population; this result needs to be
119 confirmed with additional Australian data. Rasmussen et al.⁴ suggested that Aboriginal Australians are the
120 descendants of admixture of the first OOA with later OOA populations. We failed to detect this first OOA
121 event either by Dstat (Supplementary Tables 6 and 7) or relative cross-coalescent analysis by MSMC
122 (Supplementary Figure 20). Our simulations suggest that the bias in Dstat calculation, which was interpreted
123 as the product of the first OOA population admixture with Aboriginal Australians, can instead be explained
124 by ancient hominin admixture with Aboriginal Australians.

125 To explain the genetic structure of mainland India, it has been suggested⁸ that all populations have arisen
126 from admixture between two components: (1) Ancestral North Indian (ANI) and (2) Ancestral South Indian
127 (ASI), which is genetically related to Andamanese. However, although ADMIXTURE analysis (Figure 1a)
128 showed that the Irula (ILA) and Birhor (BIR) tribal populations have high amounts of this ASI component,
129 also present in all the other non-tribal populations of Southern India examined (shown also in^{11,25}), TreeMix
130 analysis (Figure 1b) suggested that Andamanese are not directly related to this South Indian component.
131 Rather, the Andamanese are slightly closer to East Asians than to these two tribal Indian populations. Also,
132 the Andamanese do not share direct ancestry with the Australian and Papuan sequences tested (Figure 1b), as
133 has been traditionally assumed because of morphological similarities between these populations¹.

134 Since we have shown that the Andamanese and other modern Asian populations have a common origin, we
135 hypothesized that the distinct phenotype of the Andamanese should have originated by recent adaptation to
136 their environment. To detect positive selection we used the Hierarchical Boosting (HB) method, a machine-
137 learning classification framework that exploits the combined ability of some selection tests to uncover
138 features expected under the hard sweep model, while controlling for population-specific demography,
139 achieving higher power than single tests and a low rate of false positive results²⁶. We found some 1,000
140 genomic regions to have significant footprints of positive selection among the Andamanese (212 regions,
141 encompassing 107 genes, under the complete hard sweep model; and 805 regions, encompassing 509 genes,
142 under the incomplete hard sweep model). Among them, we found a significant excess of genes related to
143 body morphology, with signals in 11 of the 107 genes related to height (according to the Genetics

144 Association Database, GAD²⁷) for complete selective sweeps (Yates Chi Square=5.70, P=0.02) and 48 out of
145 509 for incomplete sweeps (Yates Chi Square=22.59, P<0.0001). Other regions under positive selection
146 included genes related to obesity or body shape and composition. It is interesting to note that these results
147 point to selective pressure on body size, likely related to low stature (in fact, the very low stature of
148 Andamanese can be recognized by the individual genotypes at height-related SNPs; see Supplementary
149 Figure 21); it could therefore represent insular dwarfism, a well-known adaptation of large animals to a
150 restricted environment that predicts a derived state for the morphology of the Andamanese. These results
151 thus provide insights into the biological bases of such adaptations, also described recently in Sardinia⁹.

152 Our analysis supports a distinct model for the human settlement of Asia and Pacific, with two novel insights
153 (Figure 3): (i) Asian populations, including ones from the Pacific, have a single origin and OOA expansion,
154 sharing a more recent common ancestor between themselves than with Europeans; our analyses do not
155 support the hypothesis of two independent OOA events, postulated a long time ago based on physical
156 appearance¹ and apparently confirmed by genetics⁴; and, (ii) Indian mainland populations, Andamanese,
157 Papuans and Aboriginal Australians (but not East Asians) carry genomic contributions from an extinct
158 hominin population, with admixture ranging between 2-3% (higher in Australians, but this estimate needs to
159 be confirmed with new data). Our results do not indicate whether or not the introgression is derived from the
160 same hominin in all populations, but in the case of the Andamanese (Supplementary Figure 22) we have
161 shown that it comes from a new unknown hominin population, that likely separated very early in the
162 hominin tree. Also, we have shown that the hominin admixture in these populations can cause a bias in Dstat
163 calculation that can be erroneously interpreted as a first OOA migration of modern. Finally, the distinctive
164 morphology of the Andamanese (and probably of other 'Negrito' populations) has probably originated from
165 strong adaptive selection as shown by the excess of genes under selection related to height and body mass,
166 and it is not an ancestral character, but derived, leading to the possibilities of understanding the basic biology
167 of a complex adaptation in an island environment.

168 **References**

- 169 1. Coon, C. S. & Hunt, E. E. *The living races of man*. (Knopf, 1966).
- 170 2. Molnar, S. *Human Variation: Races, Types and Ethnic Groups*. (Routledge, 2015).
- 171 3. Cavalli-Sforza, L. L., Menozzi, P. & Piazza, A. *The history and geography of human genes*.
172 (Princeton University Press, 1994).
- 173 4. Rasmussen, M. *et al.* An Aboriginal Australian Genome Reveals Separate Human Dispersals into
174 Asia. *Science* **334**, 94–98 (2011).
- 175 5. Huxley, T. H. On the Geographical Distribution of the Chief Modifications of Mankind. *J. Ethnol.*
176 *Soc. London* **2**, 404–412 (1870).
- 177 6. Brown, A. R. *The Andaman Islanders: A Study in Social Anthropology*. (Cambridge University Press,
178 1922).
- 179 7. Abbi, A. Is Great Andamanese genealogically and typologically distinct from Onge and Jarawa?
180 *Lang. Sci.* **31**, 791–812 (2009).
- 181 8. Reich, D., Thangaraj, K., Patterson, N., Price, A. L. & Singh, L. Reconstructing Indian population
182 history. *Nature* **461**, 489–494 (2009).
- 183 9. Montgomery, S. H. Primate brains, the ‘island rule’ and the evolution of *Homo floresiensis*. *J. Hum.*
184 *Evol.* **65**, 750–760 (2013).
- 185 10. Zoledziewska, M. *et al.* Height-reducing variants and selection for short stature in Sardinia. *Nat.*
186 *Genet.* **47**, 1352–1356 (2015).
- 187 11. Basu, A., Sarkar-Roy, N. & Majumder, P. P. Genomic reconstruction of the history of extant
188 populations of India reveals five distinct ancestral components and a complex structure. *Proc. Natl.*
189 *Acad. Sci.* **113**, 201513197 (2016).
- 190 12. Dass, F. *The Andaman Islands*.PDF. (The Good Shepherd Convent Press, 1937).
- 191 13. Green, R. E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**, 710–722 (2010).
- 192 14. Pickrell, J. K. & Pritchard, J. K. Inference of population splits and mixtures from genome-wide allele
193 frequency data. *PLoS Genet.* **8**, e1002967 (2012).
- 194 15. Schiffels, S. & Durbin, R. Inferring human population size and separation history from multiple
195 genome sequences. *Nat. Genet.* **46**, 919–925 (2014).
- 196 16. Raghavan, M. *et al.* Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans.
197 *Nature* **505**, 87–91 (2014).
- 198 17. Olalde, I. *et al.* Derived immune and ancestral pigmentation alleles in a 7,000-year-old Mesolithic
199 European. *Nature* **507**, 225–228 (2014).
- 200 18. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for present-day
201 Europeans. *Nature* **513**, 409–413 (2014).
- 202 19. Haak, W. *et al.* Massive migration from the steppe was a source for Indo-European languages in
203 Europe. *Nature* **522**, 207–211 (2015).

- 204 20. Prüfer, K. *et al.* The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature*
205 **505**, 43–49 (2014).
- 206 21. Meyer, M. *et al.* A High-Coverage Genome Sequence from an Archaic Denisovan Individual. *Science*
207 **338**, 222–226 (2012).
- 208 22. Swisher, C. C. *et al.* Latest Homo erectus of Java: potential contemporaneity with Homo sapiens in
209 southeast Asia. *Science* **274**, 1870–1874 (1996).
- 210 23. Hammer, M. F., Woerner, A. E., Mendez, F. L., Watkins, J. C. & Wall, J. D. Genetic evidence for
211 archaic admixture in Africa. *Proc. Natl. Acad. Sci. U. S. A.* **108**, 15123–15128 (2011).
- 212 24. Vernot, B., Akey, J. M. & Vernot B., A. J. M. Resurrecting surviving Neandertal lineages from
213 modern human genomes. *Science* **343**, 1017–1021 (2014).
- 214 25. Juyal, G. *et al.* Population and genomic lessons from genetic analysis of two Indian populations. *Hum.*
215 *Genet.* **133**, 1273–1287 (2014).
- 216 26. Pybus, M. *et al.* Hierarchical boosting: a machine-learning framework to detect and classify hard
217 selective sweeps in human populations. *Bioinformatics* **31**, btv493 (2015).
- 218 27. Becker, K. G., Barnes, K. C., Bright, T. J. & Wang, S. A. The genetic association database. *Nat.*
219 *Genet.* **36**, 431–432 (2004).
- 220

221 URLs

222 European Nucleotide Archive (<http://www.ebi.ac.uk/ena>), Picard tools (<http://picard.sourceforge.net/>), Broad
223 ftp server (<ftp.broadinstitute.org>), 1000 Genome ancestral file
224 (http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase1/analysis_results/supporting/ancestral_alignments/).

225 Accession codes

226 The whole-genome sequences (Andamanese vcf files) have been deposited in the European Nucleotide
227 Archive, Accession ID: PRJEB11455.

228 Acknowledgements

229 Our main funding was provided by the joint Spain-India bilateral grant PRI-PIBIN-2011-0942 from the
230 Ministerio de Economía y Competitividad (Spain). Complementary funding was from grant BFU2013-
231 43726-P from the Ministerio de Economía y Competitividad (Spain), with the support of Secretaria
232 d'Universitats i Recerca, Departament d'Economia i Coneixement de la Generalitat de Catalunya (GRC
233 2014 SGR866). Jessica Nye (Universitat Pompeu Fabra) and Chris Tyler-Smith kindly corrected in depth the
234 manuscript. Thanks to R.A. Foley for discussion and inspiring input for Figure 3.

235 Author Contributions

236 MM, FC, PPM and JB conceived and designed the project. PPM provided the samples. PPM, ZH and QL
237 sequenced samples and carried out initial analyses. MM performed the remaining genetic data analyses. FC,
238 GMDO, MP, MGN, DC, HL, PPM and JM participated in and discussed analyses. MM, FC, PPM and JB
239 wrote the manuscript.

240 Author Information

241 The authors declare no competing financial interests. Correspondence should be addressed to JB
242 (jaume.bertranpetit@upf.edu) or PPM (ppm1@nibmg.ac.in).

243 **Fig. 1: Ancestry of Indian Populations.**

244 a. ADMIXTURE analysis using 10 randomly-chosen individuals from CEU, CHB and YRI taken from
245 the 1000 Genomes Project and individuals from our data set: Punjabi (PUN), Uttar Pradesh Brahmins
246 (UBR), Rajput (RAJ), Bengali (BEN), Vellalar (VLR), Irula (ILA), Birhor (BIR), Jarawa (JAR),
247 Onge (ONG) and Riang (RIA). Results are shown for five ancestral components, the optimal number.
248 Each vertical bar represents one individual, colored according to the proportion of the five ancestral
249 components.

250 b. TreeMix analysis without migration. Africans are Yoruba (YRI), Mandenka (MAD), Mbuti pygmy
251 (MBT) and San (SAN), Europeans are French (FRN) and Sardinian (SAR), East Asians are Dai
252 (DAI) and Han Chinese (HAN), Pacific population are Papuans (PAP) and Aboriginal Australians
253 (AUS), and Indians are (BIR, ILA and RIA) and Andamanese (JAR and ONG). Inferred ancestral
254 genome information from the 1000 Genomes Project was used as outgroup. The scale bar shows 10
255 times the standard error and the x axis shows the amount of drift. Drift considered non-significant is
256 indicated by a red line, so the three branches (RIA, HAN, DAI), (ONG, JAR) and (BIR, ILA) form a
257 trichotomy.

258 MSMC Relative Cross Coalescent Rate showing genetic separation between two populations. In each
259 curve one individual was from Jarawa (JAR) and the other from either a Tribal population of India
260 (ILA, BIR or RIA), ONG, or from outside India (FRN, DAI, PAP and YRI). The x-axis shows time
261 and the y-axis shows a measure of the similarity between the two populations.

262

263

264

265

266

267

268 **Fig. 2: Fewer African derived alleles in Indian, Andamanese, Papuan and Aboriginal Australians than**
269 **Europeans or East Asians.** Each horizontal line shows the result of D-statistics [Dstat(W,X;Y,Z)] where the
270 W population is either a) French (FRN) or b) East Asian Dai (DAI). The X population is either from India:
271 Punjabi (PUN), Uttar Pradesh Brahmins (UBR), Rajput (RAJ), Bengali (BEN), Vellalar (VLR), Irula (ILA),
272 Birhor (BIR) or Riang (RIA); Andamanese: Jarawa (JAR) or Onge (ONG); and FRN, Sardinia (SAR), DAI,
273 Han Chinese (HAN), Papuans (PAP) or Aboriginal Australians (AUS); names are shown to the right of the
274 two figures. The Y population is African (Yoruba (YRI), Mandenka (MAD), Mbuti pygmy (MBT) or San
275 (SAN). Ancestral allele information from the 1000 Genomes Project is used as outgroup (Z population).
276 Colour coding of the populations: Europeans (Pink), East Asians (Deep Yellow), African (Brown), Indo-
277 Europeans (Red), Dravidians (Black), Austro Asiatics (Blue), Andamanese (Light Green), Tibeto Burman
278 (Yellow), Pacific Islanders and Australian Aboriginals (Deep Green). A positive value means that the W and
279 Y populations share more derived alleles with each other compared with X and Y, while a negative value
280 means X and Y populations share more derived allele with each other as compared with W and Y. The
281 statistically significant results (in this case defined by a Z score more or less than ± 3) are marked with a star.
282 a. Dstat results of D(FRN(W),X;AFR(Y),Ancestral(Z)).
283 b. Dstat results of D(DAI(W),X;AFR(Y),Ancestral(Z)).

284

285

286 **Fig. 3: Model of gene flow in Asia.** Red boxes indicate extinct non-African hominins who introgressed into
287 modern humans; these introgressions are marked with dotted lines. The green box indicates populations that
288 may have admixed with the new unknown hominin; Andamanese and Indian are fully analyzed here; the
289 others will have to be further studied in the future. To properly solve the question mark trichotomy would
290 require more data.

291

292 **Methods**

293 **Samples**

294 In total, 70 samples were collected from 10 Indian populations from different geographical regions,
295 linguistic affiliations and social categories (Supplementary Table 1). The 10 populations were: Punjabi
296 (PUN), Uttar Pradesh Upper caste Brahmins (UBR), Rajput (RAJ), Bengali (BEN), Vellalar (VLR), Irula
297 (ILA), Birhor (BIR), Jarawa (JAR), Onge (ONG) and Riang (RIA). The blood and saliva samples were
298 collected with voluntary informed consent from the participants. More information on the populations is
299 found in Basu et al¹¹.

300 Additional samples were also used to understand Indian populations from a global perspective. We used the
301 1000 Genomes Phase 1 data²⁸, the Great Ape Genome Project (GAGP) data²⁹, high-coverage data from three
302 Aboriginal Australians³⁰, nine Yoruba (YRI) high-coverage data and five Utah residents with Northern and
303 Western European Ancestry (CEU)³¹. We used some Ancient genome sequences: Malta¹⁶, La Braña¹⁷,
304 Loschbour and Stuttgart¹⁸. Neanderthal²⁰ and Denisova²¹ data were used to calculate the admixture level of
305 these subspecies in Indian populations. We have used the 1000 Genomes Project ancestral file³² to identify
306 the ancestral allele.

307 **Sequencing**

308 The whole-genome sequencing was done in two different places (BGI, NIBMG) using Illumina technology.
309 50 of the 70 samples were sequenced in BGI, whereas 20 were sequenced in NIBMG (Supplementary Tables
310 1 and 8). Sequencing libraries with an insert size of ~500 bp were constructed and paired-end reads were
311 generated by HiSeq 2000. The raw sequencing reads were mapped to hg19 using BWA³³. Duplicates were
312 removed by Picard tools. We followed best practice recommendations from GATK 2.8-1³⁴ using
313 IndelRealigner and BaseRecalibrator with their default values. For IndelRealigner we used 1000 Genomes
314 Phase 1 indel interval files, and for BaseRecalibrator we used dbSNP 137. Variants were called by
315 HaplotypeCaller from GATK. After creation of the raw vcf files, we used VariantRecalibrator from GATK
316 on the autosomes using dbSNP 137, HapMap 3.3, 1000 Genomes Project Omni 2.5 and 1000 Genomes
317 Project Phase 1 SNPs with high confidence, Mills and 1000 Genomes Project gold standard indels to assign
318 a well-calibrated probability to each variants; all these files were downloaded from the Broad Institute ftp
319 site (date 11/05/2013) as described in the website of GATK. The average coverage for autosomes was ~15x
320 and the accessible genome was close to 100% (Supplementary Table 8). Though the sequencing was done in

321 two different institutes, Principal Component (PC) and ADMIXTURE analysis (Supplementary Note)
322 demonstrated a very tight clustering for samples from the same population, suggesting that influences from
323 the two sequencing centers were not detectable.

324 **Relatedness, Inbreeding and Homozygosity Run**

325 Relatedness was calculated using KING³⁵ software with 13,679,600 autosomal bi-allelic SNPs. Inbreeding
326 was calculated by vcfTools³⁶ using the same SNPs and the default parameters. Homozygosity runs were done
327 by PLINK v1.07³⁷ software using 4,475,795 autosomal bi-allelic unlinked SNPs with the default parameters.
328 SNPs were unlinked according to the variance inflation factor (VIF) method implemented in PLINK with a
329 window size of 50 SNPs, a step size of 5, and a variance inflation factor of 2.

330 **PCA**

331 SmartPCA from the EIGENSOFT package³⁸ was used for PCA. We kept only autosomal, bi-allelic SNPs
332 that have Minor Allele Frequency (MAF) of at least 0.05. We also removed SNPs which had missing
333 information for any individual. Only 10 individuals per population from the 1000 Genomes Project data
334 were kept to avoid sample size bias.

335 **Admixture**

336 ADMIXTURE³⁹ was used to calculate admixture per individual with the same filters as the PCA analysis.
337 To explore the optimal number of ancestral populations (k), we used k= 2–6, performing ten iterations for
338 each. The best k value was estimated using the cross-validation error method implemented in
339 ADMIXTURE.

340 **MSMC**

341 Effective population size and population separation over time were calculated using MSMC¹⁵. Only
342 autosomes were used. MSMC recommendations were followed to create input files from BAM files. We
343 phased genomes using 1000 Genomes Project Phase 3 data as the reference using Shapeit⁴⁰.

344 **Dstat**

345 ADMIXTOOLS were used⁴¹ for Dstat analysis. To reduce biases (especially ascertainment bias), we called
346 variants from India and the Great Ape Genome Project (only humans) together as described above. SNP
347 information from Aboriginal Australians, Neanderthal, Denisova and other ancient samples were extracted
348 as described in Supplementary Information 5. Ancestral information was extracted from the fasta file given
349 on the 1000 Genomes Project website.

350 **TreeMix**

351 TreeMix¹⁴ was used to analyse the divergence of the populations from each other, using the data described
352 above. We used migration values from 0 to 20. The inferred ancestral genome was used to root the tree. To
353 allow for linkage disequilibrium (LD) we used the -k flag. The LD blocks were defined as 1 Mb in length,
354 which in our case corresponds to about 5,000 SNPs.

355 **Simulations**

356 For simulations, we used ms⁴² following published parameters⁴³. We added Andamanese parameters
357 determined from our inferences about Andamanese ancestry (See Supplementary Note).

358 **Dadi and a three-population model for Archaic Admixture**

359 We first built a null model without introgression of archaic hominins into the Andamanese using dadi-1.7.0⁴⁴
360 following parameters from Gravel et al⁴³. Then a three-population model for archaic admixture was
361 implemented to estimate the divergence of this unknown population from humans and the time of admixture
362 with Andamanese by simulating 2% of hominin genome introgression into Andamanese at different time
363 points.

364 **Selection**

365 This analysis, used Andamanese genomes from our data and YRI sequences from Complete Genomics³¹ and
366 merged them. After removing any SNP which has missing information for any individual, we phased the
367 Andamanese with Shapeit⁴⁵ using 1000 Genomes Project phase 1 samples as a reference⁴⁰. Then, the
368 following selection tests were performed on the data:

- 369 1. Tajima's D⁴⁶.
- 370 2. CLR⁴⁷.
- 371 3. Fay and Wu's H⁴⁸.
- 372 4. Fu & Li's D⁴⁹.
- 373 5. XP-EHH⁵⁰.
- 374 6. ΔiHH ⁵¹.
- 375 7. iHS⁵¹.
- 376 8. EHH average⁵².

377 After calculating all tests, we ran the boosting algorithm²⁶ using parameters both from the East Asian and the
378 European hierarchical boosting strategy (simulated under neutrality and under selection using cosi with
379 demographic models from Schaffner et al⁵³ for both East Asian and European demography and then
380 calculating the best strategy to detect selection). In fact, results for the hierarchical boosting strategy for non-
381 African populations are very similar (Supplementary Note). Information about body size genes was obtained
382 from the Genetics Association Database²⁷ and their functional annotation from ANNOVAR⁵⁴

383 **Dstat with sliding windows and Sstar**

384 To identify candidate introgressed regions from an unknown hominin, we calculated Dstat per individual for
385 50 kb regions with sliding windows of 5kb and retained regions where Andamanese have fewer African
386 derived alleles than Europeans or East Asians:

$$387 \quad D_{stat} = \frac{\sum(F_w - F_x)(F_y - F_z)}{\sum(F_w + F_x - 2F_w F_x)(F_y + F_z - 2F_y F_z)}$$

388 F is the allele frequency in w, x, y or z populations.

389 We ran TreeMix on the putative introgressed regions (Supplementary Note) and Sstar²⁴ to refine the
390 identification of the introgressed hominin haplotypes thus only taking regions which is positive for both
391 Dstat by sliding windows and Sstar (Supplementary Note).

392 Methods Reference

- 393 28. Abecasis, G. R. *et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**,
394 56–65 (2012).
- 395 29. Prado-Martinez, J. *et al.* Great ape genetic diversity and population history. *Nature* **499**, 471–475
396 (2013).
- 397 30. Sudmant, P. H. *et al.* Global diversity, population stratification, and selection of human copy-number
398 variation. *Science* **349**, aab3761 (2015).
- 399 31. Drmanac, R. *et al.* Human genome sequencing using unchained base reads on self-assembling DNA
400 nanoarrays. *Science* **327**, 78–81 (2010).
- 401 32. Abecasis, G. R. *et al.* A map of human genome variation from population-scale sequencing. *Nature*
402 **467**, 1061–1073 (2010).
- 403 33. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
404 *Bioinformatics* **25**, 1754–1760 (2009).
- 405 34. McKenna, A. *et al.* The genome analysis toolkit: A MapReduce framework for analyzing next-
406 generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
- 407 35. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.
408 *Bioinformatics* **26**, 2867–2873 (2010).
- 409 36. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 410 37. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage
411 analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 412 38. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association
413 studies. *Nat. Genet.* **38**, 904–909 (2006).
- 414 39. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated
415 individuals. *Genome Res.* **19**, 1655–1664 (2009).
- 416 40. Delaneau, O. *et al.* Integrating sequence and array data to create an improved 1000 Genomes Project
417 haplotype reference panel. *Nat. Commun.* **5**, 3934 (2014).
- 418 41. Patterson, N. *et al.* Ancient admixture in human history. *Genetics* **192**, 1065–1093 (2012).
- 419 42. Hudson, R. R. Generating samples under a Wright-Fisher neutral model of genetic variation.
420 *Bioinformatics* **18**, 337–338 (2002).
- 421 43. Gravel, S. *et al.* Demographic history and rare allele sharing among human populations. *Proc. Natl.*
422 *Acad. Sci. U. S. A.* **108**, 11983–11988 (2011).
- 423 44. Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H. & Bustamante, C. D. Inferring the joint
424 demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genet.*
425 **5**, (2009).
- 426 45. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for thousands of
427 genomes. *Nature Methods* **9**, 179–181 (2011).

- 428 46. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism.
429 *Genetics* **123**, 585–595 (1989).
- 430 47. Nielsen, R. *et al.* Genomic scans for selective sweeps using SNP data. *Genome Res.* **15**, 1566–1575
431 (2005).
- 432 48. Fay, J. C. & Wu, C. I. Hitchhiking under positive Darwinian selection. *Genetics* **155**, 1405–1413
433 (2000).
- 434 49. Fu, Y. X. & Li, W. H. Statistical tests of neutrality of mutations. *Genetics* **133**, 693–709 (1993).
- 435 50. Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive selection in human
436 populations. *Nature* **449**, 913–918 (2007).
- 437 51. Voight, B. F., Kudaravalli, S., Wen, X. & Pritchard, J. K. A map of recent positive selection in the
438 human genome. *PLoS Biol.* **4**, 0446–0458 (2006).
- 439 52. Sabeti, P. C. *et al.* Detecting recent positive selection in the human genome from haplotype structure.
440 *Nature* **419**, 832–837 (2002).
- 441 53. Schaffner, S. F. *et al.* Calibrating a coalescent simulation of human genome sequence variation.
442 *Genome Res.* **15**, 1576–1583 (2005).
- 443 54. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-
444 throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).





