

1 Genetic landscapes reveal how human genetic 2 diversity aligns with geography

3 Benjamin Marco Peter^{1,2}, Desislava Petkova^{3,4} & John Novembre^{1,5}

4 ¹ Department of Human Genetics, University of Chicago ² Max Planck Institute for Evolutionary Anthropology, Leipzig,
5 Germany

6 ³ Wellcome Trust Center for Human Genetics, University of Oxford, UK ⁴ Present Address: Procter & Gamble,
7 Brussels, Belgium ⁵ Department of Ecology & Evolution, University of Chicago

8

9 **Geographic patterns in human genetic diversity carry footprints of population history^{1,2}**
10 **and provide insights for genetic medicine and its application across human**
11 **populations^{3,4}. Summarizing and visually representing these patterns of diversity has**
12 **been a persistent goal for human geneticists⁵⁻¹⁰, and has revealed that genetic**
13 **differentiation is frequently correlated with geographic distance. However, most**
14 **analytical methods to represent population structure¹¹⁻¹⁵ do not incorporate geography**
15 **directly, and it must be considered *post hoc* alongside a visual summary. Here, we use a**
16 **recently developed spatially explicit method to estimate “effective migration” surfaces to**
17 **visualize how human genetic diversity is geographically structured (the EEMS method¹⁶).**
18 **The resulting surfaces are “rugged”, which indicates the relationship between genetic**
19 **and geographic distance is heterogenous and distorted as a rule. Most prominently,**
20 **topographic and marine features regularly align with increased genetic differentiation**
21 **(e.g. the Sahara desert, Mediterranean Sea or Himalaya at large scales; the Adriatic, inter-**
22 **island straits in near Oceania at smaller scales). In other cases, the locations of**
23 **historical migrations and boundaries of language families align with migration features.**
24 **These results provide visualizations of human genetic diversity that reveal local patterns**
25 **of differentiation in detail and emphasize that while genetic similarity generally decays**
26 **with geographic distance, there have regularly been factors that subtly distort the**
27 **underlying relationship across space observed today. The fine-scale population**
28 **structure depicted here is relevant to understanding complex processes of human**
29 **population history and may provide insights for geographic patterning in rare variants**
30 **and heritable disease risk.**

31

32 In many regions of the world, genetic diversity “mirrors” geography in the sense that genetic
33 differentiation increases with geographic distance (“isolation by distance”¹⁷⁻¹⁹); However, due to
34 the complexities of geography and history, this relationship is not one of constant
35 proportionality. The recently developed analysis method EEMS visualizes how the isolation-by-
36 distance relationship varies across geographic space¹⁶. Specifically, it uses a model based on a
37 local “effective migration” rate. For several reasons, the effective migration rates inferred by
38 EEMS do not directly represent levels of gene flow¹⁶; however they are useful for conveying
39 spatial population structure: populations in areas of high effective migration are genetically more
40 similar than other populations at the same geographic distance, and conversely, low migration
41 rates imply genetic differentiation increases rapidly with distance. In turn, a map of inferred

42 patterns of effective migration can provide a compact visualization of spatial genetic structure
43 for large, complex samples.

44
45 We apply EEMS on a combination of 27 existing single nucleotide polymorphism (SNP)
46 datasets. In total, these comprise 6066 individuals from 419 locations across Eurasia and Africa
47 (Extended Data Table 1), which we organize in seven analysis panels: an overview Afro-
48 Eurasian panel (AEA), four continental-scale panels, and two panel of Southern African
49 KhoeSan and Bantu speakers. For all analysis panels, the inferred EEMS surfaces are
50 “rugged”, with numerous high and low effective migration features (Fig 1a, Fig 2) that are
51 strongly statistically supported when compared to a uniform-migration model (Extended Data
52 Table 2). The regions of depressed effective migration often align in long, connected stretches
53 that are present in more than 95% of MCMC iterations. To facilitate discussion, we annotate
54 these stretches with dashed lines and refer to them as “troughs” of effective migration (Figs. 1a,
55 2, Extended Data Figs. 2-4). Conversely, intermediate- and high-migration areas between
56 troughs are referred to as corridors.

57
58 In the broad overview Afro-Eurasia panel (Fig. 1; $n=4,697$ samples; 370 locales; $F_{ST} = 0.071$)
59 we see that troughs often align with topographical obstacles to migration, such as deserts
60 (Sahara), seas (Mediterranean, Red, Black, Caspian, East China Seas) and mountain ranges
61 (Ural, Himalayas, Caucasus). Among the main features are several large regions that have
62 mostly high effective migration, such as Europe, East Asia, Sub-Saharan Africa and Siberia.
63 Several large-scale corridors are inferred that represent long-range genetic similarity, for
64 example: India is connected by two corridors to Europe (a southern one through Anatolia and
65 Persia ‘SC’, and a northern one through the Eurasian Steppe ‘NC’); East Asia (EA) is connected
66 to Siberia and to southeast Asia and Oceania. The island populations of the Andaman islands
67 (Onge) and New Guinea show troughs nearly contiguously around them – possibly reflecting a
68 history of relative isolation^{20,21}.

69
70 Analyses on a finer geographic scale highlights subtler features (e.g. compare Europe in Fig. 1
71 vs Fig. 2a), and reveal that levels of differentiation differ both on both a local and continental
72 scale (Extended Data Table 2). At these finer scales we continue to see troughs that align with
73 landscape features, though increasingly we see troughs and corridors that coincide with contact
74 zones of language groups and proposed areas of human migrations. For example, in Europe
75 (Fig. 2b) we observe troughs (NS, CE) roughly between where Northern Slavic speaking
76 peoples currently reside relative to west Germanic speakers, and relative to the linguistically
77 complex Caucasus region. In India (Fig. 2e), troughs demarcate regions with samples of
78 Austroasiatic and Dravidian speakers, as well as central India (CI) relative to Northwestern India
79 (Sindhi, Punjabi) and Pakistan. In Southeast Asia (Fig. 2h), troughs align with several straits in
80 the Malay archipelago, but we also observe a corridor from Taiwan through Luzon to the Lower
81 Sunda Islands (LSI), and further to Melanesia, perhaps reflecting the Austronesian expansion.
82 In Africa (Fig. 2g), a trough aligns with the Sahara desert and extends south-eastward into a
83 geographic region that is a complex linguistic contact zone with Afro-Asiatic speakers (North),
84 Nilo-Saharan and Niger-Congo speakers, and linguistic isolates, the Hadza and Sandawe
85 (South). Notably, the contiguity of the South-East extension of this trough is sensitive to the

86 inclusion of the Hadza and Sandawe (Extended Data Fig. 9). In Sub-Saharan Africa we also find
87 corridors perhaps reflecting the Bantu expansion from West- into Southern and Eastern Africa,
88 where contact with Nilo-Saharan speakers resulted in complex local structure. In Southern
89 Africa, the structure seen in separate EEMS maps prepared for Bantu and Khoe-San speakers
90 (Fig. 2k/l) appear distinct from each other, illustrating that in some cases, different language
91 groups can maintain independent genetic structure in the same geographic region.

92
93 Different data visualization methods invariably emphasize different aspects of the data.
94 Contrasting with EEMS, we find that the widely used, non-spatial principal component analysis
95 (PCA) highlights large-scale geography, and the PCA-biplots typically reflect the strongest
96 gradients of diversity in a panel. For example, PCA highlights differentiation along an Out-of-
97 Africa differentiation axis in the AEA panel (Fig. 1b), the circum-Mediterranean and circum-
98 Saharan distribution of diversity in Western Eurasia and Africa, respectively, and gradients from
99 Europe into East Asia and South Asia in the Central/Eastern Eurasian panel (Fig. 2). On the
100 other hand, EEMS emphasizes local features, in particular troughs between adjacent groups
101 that are often imperceptible in the PCA-biplots. This is likely due to geographical information
102 allowing EEMS to discern subtle structure from effects of uneven sampling¹⁶. PCA easily
103 identifies outlier or admixed individuals (e.g. in Africa) that are not made apparent in EEMS
104 except when exploring model fit to find populations that are fit poorly (Extended Data Fig. 8).
105 Locally differentiated populations such as the Sardinians, Basques, and Finnish strongly shape
106 the PCA results (compare Fig. 2d to e.g. ref ¹⁷), whereas they are typically placed in low-
107 migration regions in EEMS. We also compare the model fit of EEMS and low-rank PCA (using
108 the first 2, 10 and 100 components) to the observed genetic distances as a means of assessing
109 how well each low-dimensional approach conveys structure in full genetic data. EEMS performs
110 better for small-scale panels, but PCA provides a better fit on the larger-scale AEA and CEA
111 panels (Extended Data Figure 5). We hypothesize EEMS tends to represent local genetic
112 differences relatively well, and this is supported by an analysis where we stratify the residuals of
113 genetic distances (Extended Data Fig. 6): In most panels EEMS fits best in the lowest
114 percentiles (corresponding to local differences), and the fit quality tends to decrease for larger
115 genetic distances.

116
117 Overall, the maps we present provide a compact summary of the complex relationship of genes
118 and geography in human populations. In contrast to methods that identify short bursts of gene
119 flow (“admixture”) between diverged populations^{22–24}, EEMS models local migration between
120 nearby groups to represent heterogeneous isolation-by-distance patterns. This leads to the first
121 of a few limitations that must be considered in interpretation. In some cases, isolation-by-
122 distance may not be the most appropriate model, and human population may overlap spatially
123 while maintaining differentiation. This can happen either for large periods of time (e.g. Southern
124 Africa Bantu vs. KhoeSan speakers) or due to recent migration, displacement or admixture
125 events (e.g. in Central Asia). In cases of population structure due to “outliers”, we found that
126 running EEMS at the highest resolutions that are computationally feasible results in easier
127 interpretable plots as the degrees of freedoms of the surface are high enough that these
128 samples can be placed in regions of isolation.

129

130 Second, the maps inferred here represent a model of gene flow that predicts genetic diversity in
131 humans sampled today – a fuller representation would represent genetic structure dynamically
132 through time. This is especially relevant as ancient DNA data have recently suggested human
133 population structure can be surprisingly dynamic (e.g. ref. ²⁵).

134

135 Third, the effective migration rates and their scales needs be interpreted with care. In each of
136 our maps the overall levels of differentiation are consistently low across all populations, and
137 EEMS draws attention to where differentiation is slightly elevated or depressed relative to
138 expectations from geographic distance. Low effective migration between a pair of populations
139 does not imply an absence of migration nor large levels of absolute differentiation; conversely,
140 high levels of effective migration do not imply ongoing gene flow. The emergence of migration
141 features in the EEMS maps that align with known topography, past historical migrations, and/or
142 linguistic/cultural distributions does not prove a causal connection and does not constitute a
143 formal statistical test. Formally testing the influence of specific features and environmental
144 variables on migration rates remain important future tasks that will require extending EEMS or
145 using different frameworks²⁶.

146

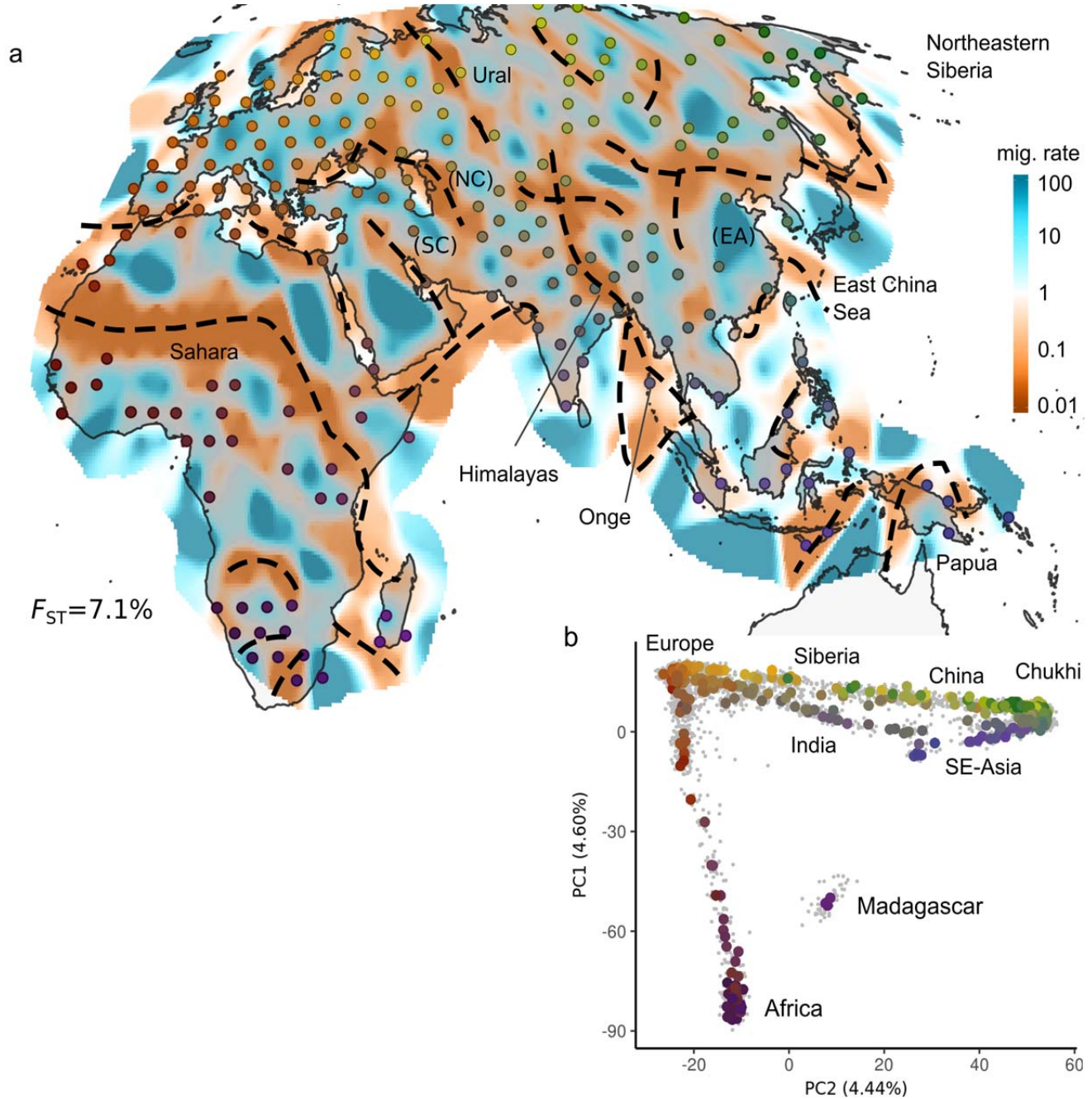
147

148 Finally, ascertainment decisions of which samples to include will affect the outcome of any
149 analysis. When there is a feature inferred in a region with few samples, the exact positioning of
150 the inferred change on the map will be imprecise (e.g. the trough presumably associated with
151 the English Channel in Fig 2b). The maps of posterior variance (Extended Data Figures 2 and
152 4) partly convey where there is uncertainty in positioning, but caution is still warranted as the
153 modelling assumptions will introduce further uncertainty. In other cases, the presence or
154 absence of a particular group may open or close corridors, sometime depending on resolution.
155 Examples of this are the Kusundas, a Nepali people with both Tibetan and Indian ancestry
156 causing a corridor through the Himalayas, the Kalmyk, a mongolian people in Southern Russia
157 that are linked by a corridor to Mongolia in the CEA, but not the AEA panel, and this corridor
158 disappears if the Kalmyks are excluded from the analysis and the Hadza and Sandawe, which
159 cause inference of a trough in Eastern Africa when included.

160

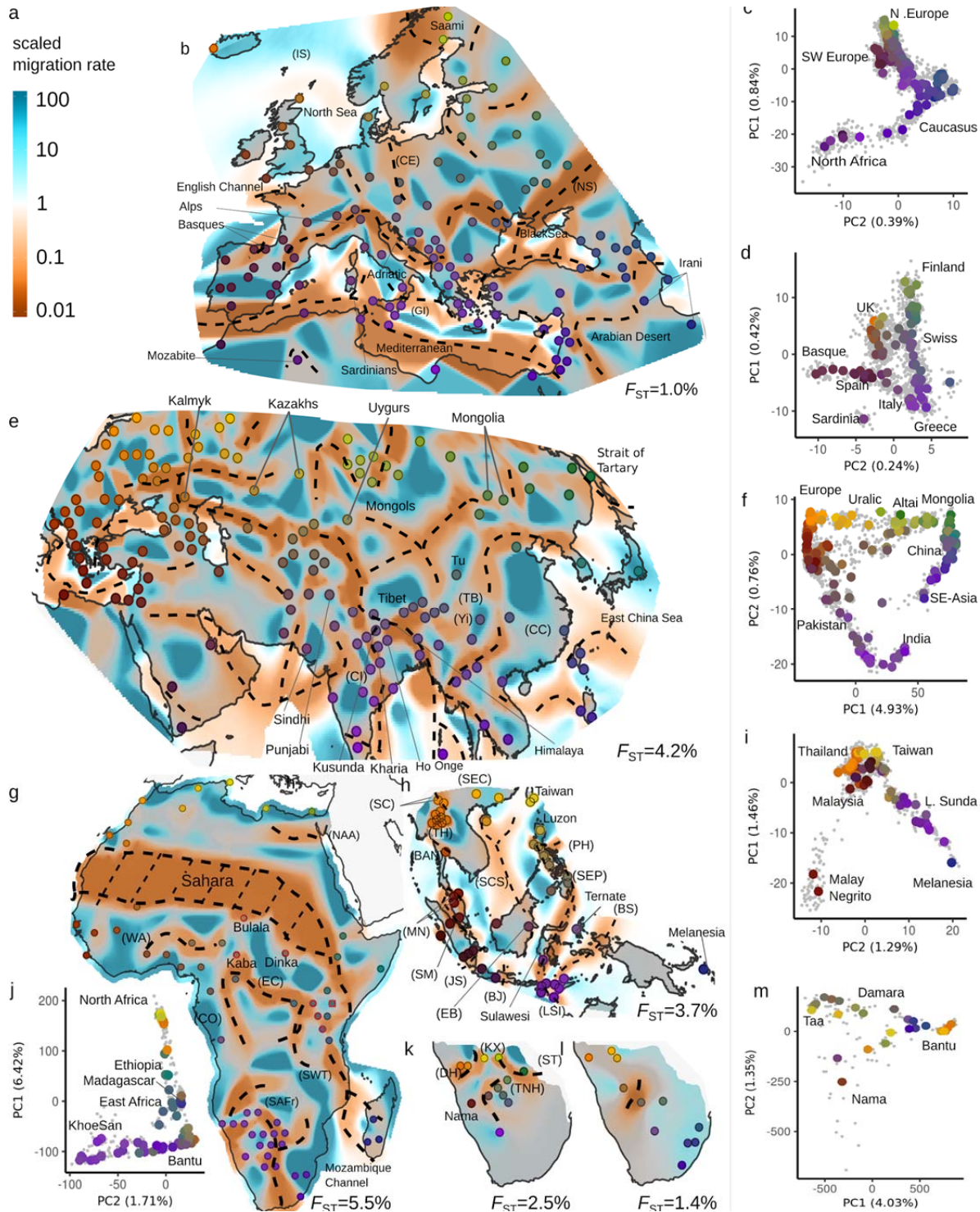
161 Nonetheless, the maps presented here provide a useful representation of human genetic
162 diversity, that complements results from geography-agnostic methods. Our results emphasize
163 the importance of geographical features on shaping human genetic history and help describe
164 fine-scale patterns of human genetic diversity²⁷. By using recent large-scale SNP data and a
165 novel analysis method, our work expands beyond previous studies of gene flow in humans²⁸⁻³⁰.
166 Our rugged migration landscapes suggest a synthesis of the clusters versus clines paradigms
167 for human structure^{7,8,31}: By revealing both sharp and diffuse features that structure human
168 genetic diversity, our results suggest that more continuous definitions of ancestry in human
169 population genetics should complement models of discrete populations with admixture. As rare
170 disease variants are commonly geographically localized^{32,33}, the maps presented here may help
171 predict regions where clustering of alleles should be expected. The maps also annotate
172 present-day population structure that ancient DNA and historical/archaeological studies can aim
173 to explain.

174
175



176
177
178
179
180
181
182
183
184
185

Figure 1: Large-scale patterns of population structure. a: EEMS posterior mean effective migration surface for Afro-Eurasia (AEA) panel. Regions and features discussed in the main text are labeled. Approximate location of troughs are annotated with dashed lines (see Extended Data Figure 2). b: PCA plot of AEA panel: Individuals are displayed as grey dots, Colored dots reflect median of sample locations; with colors reflecting geography and matching with the EEMS plot. Locations displayed in the EEMS plot reflect the position of populations after alignment to grid vertices used in the model (see methods). For exact locations, see annotated Extended Data Figure 2 and Table S1. The displayed value of F_{ST} emphasizes the low absolute level of differentiation in human SNP data.



186
 187
 188 **Figure 2: Regional patterns of genetic diversity.** **a:** scale bar for relative effective migration rate. Posterior effective migration
 189 surfaces for **b:** Western Eurasia (WEA) **e:** Central/Eastern Eurasia (CEA) **g:** Africa (AFR) **h:** South East Asian (SEA) **k:** Southern
 190 African KhoesSan (SAKS) **i:** Southern African Bantu (SAB) analysis panels.; In panel **g**, red circles indicate Nilo-Saharan speakers.
 191 Approximate location of troughs are shown with dashed lines (see Extended Data Figure 4). PCA plots: **c:** WEA **d:** Europeans in
 192 WEA **f:** CEA **i:** SEA **j:** AFR **m:** SAHG+SAB. Individuals are displayed as grey dots. Large dots reflect median PC position for a
 193 sample; with colors reflecting geography matched to the corresponding EEMS figure. In the EEMS plots, approximate sample
 194 locations are annotated. For exact locations, see annotated Extended Data Figure 4 and Table S1. Features discussed in the main
 text and supplement are labeled. F_{ST} values per panel emphasize the low absolute levels of differentiation.

195 References

- 196 1. Veeramah, K. R. & Hammer, M. F. The impact of whole-genome sequencing on the reconstruction of human population history.
197 *Nat. Rev. Genet.* **15**, 149–162 (2014).
- 198 2. Schraiber, J. G. & Akey, J. M. Methods and models for unravelling human evolutionary history. *Nat. Rev. Genet.* **16**, 727–740
199 (2015).
- 200 3. Bustamante, C. D., Burchard, E. G. & De la Vega, F. M. Genomics for the world. *Nature* **475**, 163–165 (2011).
- 201 4. Popejoy, A. B. & Fullerton, S. M. Genomics is failing on diversity. *Nature* **538**, 161–164 (2016).
- 202 5. Roberts, L. How to sample the world’s genetic diversity. *Science* **257**, 1204–1205 (1992).
- 203 6. Cavalli-Sforza, L. L., Menozzi, P. & Piazza, A. *The History and Geography of Human Genes*. (Princeton University Press,
204 1994).
- 205 7. Rosenberg, N. A. *et al.* Genetic structure of human populations. *Science* **298**, 2381–2385 (2002).
- 206 8. Serre, D. & Pääbo, S. Evidence for gradients of human genetic diversity within and among continents. *Genome Res.* **14**, 1679–
207 1685 (2004).
- 208 9. International HapMap Consortium. A haplotype map of the human genome. *Nature* **437**, 1299–1320 (2005).
- 209 10. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- 210 11. Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of population structure using multilocus genotype data. *Genetics* **155**,
211 945–959 (2000).
- 212 12. Pickrell, J. K. & Pritchard, J. K. Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS*
213 *Genet.* **8**, e1002967 (2012).
- 214 13. Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS Genet.* **2**, e190 (2006).
- 215 14. Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of population structure using dense haplotype data. *PLoS*
216 *Genet.* **8**, e1002453 (2012).
- 217 15. Lipson, M. *et al.* Efficient moment-based inference of admixture parameters and sources of gene flow. *Mol. Biol. Evol.* **30**,
218 1788–1802 (2013).
- 219 16. Petkova, D., Novembre, J. & Stephens, M. Visualizing spatial population structure with estimated effective migration surfaces.
220 *Nat. Genet.* **48**, 94–100 (2016).
- 221 17. Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).
- 222 18. Ramachandran, S. *et al.* Support from the relationship of genetic and geographic distance in human populations for a serial
223 founder effect originating in Africa. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15942–15947 (2005).
- 224 19. Wang, C., Zöllner, S. & Rosenberg, N. A. A Quantitative Comparison of the Similarity between Genes and Geography in
225 Worldwide Human Populations. *PLoS Genet.* **8**, e1002886 (2012).
- 226 20. Reich, D., Thangaraj, K., Patterson, N., Price, A. L. & Singh, L. Reconstructing Indian population history. *Nature* **461**, 489–494
227 (2009).
- 228 21. Pugach, I., Delfin, F., Gunnarsdóttir, E., Kayser, M. & Stoneking, M. Genome-wide data substantiate Holocene gene flow from
229 India to Australia. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 1803–1808 (2013).

- 230 22. Patterson, N. J. *et al.* Ancient Admixture in Human History. *Genetics* **112**, 145037 (2012).
- 231 23. Loh, P.-R. *et al.* Inferring admixture histories of human populations using linkage disequilibrium. *Genetics* **193**, 1233–1254
- 232 (2013).
- 233 24. Hellenthal, G. *et al.* A Genetic Atlas of Human Admixture History. *Science* **343**, 747–751 (2014).
- 234 25. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature* **513**, 409–
- 235 413 (2014).
- 236 26. Hanks, E. M. & Hooten, M. B. Circuit Theory and Model-Based Inference for Landscape Connectivity. *J. Am. Stat. Assoc.* **108**,
- 237 22–33 (2013).
- 238 27. Baker, J. L., Rotimi, C. N. & Shriner, D. Human ancestry correlates with language and reveals that race is not an objective
- 239 genomic classifier. *Sci. Rep.* **7**, 1572 (2017).
- 240 28. Barbujani, G. & Sokal, R. R. Zones of sharp genetic change in Europe are also linguistic boundaries. *Proc. Natl. Acad. Sci. U.*
- 241 *S. A.* **87**, 1816–1819 (1990).
- 242 29. Barbujani, G. & Belle, E. M. S. Genomic boundaries between human populations. *Hum. Hered.* **61**, 15–21 (2006).
- 243 30. Pagani, L. *et al.* Genomic analyses inform on migration events during the peopling of Eurasia. *Nature* **538**, 238–242 (2016).
- 244 31. Rosenberg, N. A. *et al.* Clines, clusters, and the effect of study design on the inference of human population structure. *PLoS*
- 245 *Genet.* **1**, e70 (2005).
- 246 32. Gibson, G. Rare and common variants: twenty arguments. *Nat. Rev. Genet.* **13**, 135–145 (2012).
- 247 33. Mathieson, I. & McVean, G. Differential confounding of rare and common variants in spatially structured populations. *Nat.*
- 248 *Genet.* **44**, 243–246 (2012).

249 Material and Methods

250 Merging pipeline

251 We obtained SNP genotype data from 27 different studies (Extended Data Table 1). Processing

252 was done using a reproducible snakemake pipeline³⁴ available under

253 <http://github.com/NovembreLab/eems-merge>, heavily relying on plink 1.9³⁵ for handling

254 genotypes. The sources differ in the input format and pre-processing, however in general we

255 performed the following steps:

256

- 257 1. Remove all non-autosomal, non-SNP variants
 - 258 2. Map SNP to forward strand of human reference genome b37 coordinates using chip
 - 259 manufacturer metadata files or SNP identifiers
 - 260 3. Remove strand-ambiguous A/T and G/C variants
- 261

262 The remaining SNPs were then merged using successive plink --bmerge commands into a

263 single master dataset with 9,003 individuals and 1.9M SNPs but a total genotyping rate of only

264 20.6%. 46 SNPs were removed because different studies reported different alternative alleles.

265 We used a relationship filter of 0.6 using the "--rel-cutoff 0.6" flag in plink to remove 667 closely

266 related individuals or duplicates. After merging, each analysis panel had missingness rates
267 <0.5% (AEA=0.2%, WEA=0.3%, CEA=0.2%, SEA=0.5%, AFR=0.2%, SAHG=0.1%). In all
268 panels, all SNPs passed a one-sided HWE-test (p -value < 10^{-5}), with the exception of SEA,
269 where nine (out of 7553 SNPs) failed and were excluded.

270 Data Retrieval and Filtering

271 Human Origins data set²⁵

272 Sampling location information was obtained from table S9.4 of ref. ²⁵, and the data were shared
273 by David Reich. We used the population information in the `vdata` subset of all ascertainment
274 panels, except for the analysis where we assess ascertainment bias. The utility `convert` from
275 `admixtools`²² was used to convert the data into plink format.

276 Estonian Biocentre data

277 The data generated by the Estonian Biocentre³⁶ were provided in plink format by Mait Metspalu
278 on 10/30/15, along with location information where it was available. This data set contained
279 1,282,568 SNPs. Of those, 6770 SNPs had non-unique ids and were removed.

280 HUGO Pan-Asian SNP consortium³⁷

281 The data were downloaded on 6/24/15 from www.biotec.or.th/PASNP. Location-metadata were
282 obtained on the same day from the map on the same website, and individuals were matched to
283 populations using the individual identifiers. All individuals with the same tag were assigned the
284 median of all locations from that tag. The data were first lifted onto hg19 (with 5 out of 54794
285 SNPs being removed), and then re-formatted into binary plink format. Due to the small size of
286 the chip used and the low overlap with the human origins array in particular, we only consider
287 this data in the South-East Asian panel.

288 Uniform global sample³⁸

289 This data were downloaded on 6/20/15 from [http://jorde-](http://jorde-lab.genetics.utah.edu/pub/affy6_xing2010/)
290 [lab.genetics.utah.edu/pub/affy6_xing2010/](http://jorde-lab.genetics.utah.edu/pub/affy6_xing2010/). Sampling locations were provided by Jinchuan Xing.
291 We used version 32 of the annotation file obtained on 6/19/15 from affymetrix.com to map SNPs
292 onto hg19, remove strand-ambiguous SNPs and to flip SNPs that were on the minus-strand.

293 POPRES data³⁹

294 POPRES data were obtained under dbGAP study accession phs000145 to John Novembre,
295 and we used the data as processed in ref ¹⁷, and only retain individuals for which all
296 grandparents were from the same country, and labelled the Swiss sample according to self-
297 reported language. We used version 32 of the annotation file obtained on 6/19/15 from
298 www.affymetrix.com ("Mapping250K_sp.na32.annot.csv" and
299 "Mapping250K_Sty.na32.annot.csv") to filter SNPs that did not map onto hg19 and we removed
300 strand-ambiguous AT and GC polymorphisms.

301 African data

302 Data from refs ^{40,41} were obtained on 04/19/17 from David Comas' website under
303 <http://www.biologiaevolutiva.org/dcomas/?p=607>. We used version 32 of the annotation file
304 "GenomeWideSNP_6.na32.annot.csv" obtained on 6/19/15 from affymetrix.com to map SNPs
305 onto hg19, remove strand-ambiguous SNPs and to flip SNPs that were on the minus-strand.

306 South-East Asian data⁴²

307 The data were obtained on 7/14/15 from Mark Stoneking in three different source files. After
308 merging the three different source files, SNPs not mapping to hg19 using the annotation file
309 "GenomeWideSNP_6.na32.annot.csv" were removed, as were AT and GC SNPs. Sampling
310 locations were extracted from Figure 1 of ref ⁴²

311 Mediterranean Panel⁴³

312 Data were obtained on 8/13/15 in binary plink format from
313 http://drineas.org/Maritime_Route/RAW_DATA/PLINK_FILES/MARITIME_ROUTE.zip.
314 Sampling location information was obtained from Supplementary Table 3 in ref. ⁴³. SNPs not
315 mapping to hg19 using the annotation file "GenomeWideSNP_6.na32.annot.csv" were removed,
316 as were AT and GC SNPs.

317 Tibetan and Himalayan data

318 Data from refs ⁴⁴⁻⁴⁶ were obtained from Choongwon Jeong and Anna Di Rienzo. We used the
319 same filtering as in the ⁴⁴ study, but only added the samples originating from these three studies
320 with permission from the respective authors.

321 Combining Meta-information

322 All sources with the exception of the Estonian Biocentre data provided (approximate) sampling
323 coordinates. However, the level of accuracy varied between sources, with some providing
324 specific ethnicities, some (such as POPRES) only providing country information and others just
325 providing city- or state-level information. For POPRES-derived data, and most countries, we
326 assigned individuals to the country's centerpoint, with the exception of Sweden and Finland,
327 which were assigned their capital. For the Estonian Biocentre data, sampling location data were
328 highly heterogeneous. Samples that could not be confidently assigned to a region with an
329 accuracy of 100km were excluded. For populations with samples from multiple studies, the
330 most accurate source location was used. For locations covered with different accuracy, only the
331 most accurate samples were retained. For example, we dropped all Spanish individuals from
332 POPRES (only country level data), as the Human Origins data provided higher resolution, with
333 samples from eleven different regions in Spain. The resulting table is given as Table S1.

334

335 Language data

336 To validate troughs correlating with presumed language barriers, we cross-referenced the
337 genetic data with linguistic data from the Glottolog 3.2 database ⁴⁷. To do so, we compared the

338 correlation of pairwise genetic distance and geographic distances within and between pairs of
339 language groups. As there was frequently no primary data recording the language of speakers,
340 we proceeded as follows: For population identifiers that correspond to languages / or ethnic
341 groups with a clear majority language, we used that language. For samples with country-level
342 information where the country has a clear majority language (e.g. Germany, Slovenia), that
343 language was assigned (Table S1). Otherwise, if a sample was from a region with a clear
344 majority language that is not obviously due to recent colonization, that language was assigned.
345 All other samples were not assigned a language. For simplicity, we group Nilotic, Central
346 Sudanic and Mande languages into “Nilo-Saharan”, Khoe, Kxa and Tuu speakers into
347 “KhoeSan” and Armenic, Circassian, Kartvelian and Nakh-Daghesanian into “Caucasus”. For all
348 troughs we hypothesize that they align with boundaries between linguistic groups, we now
349 perform a partial mantel test comparing genetic distances and language groups as a categorical
350 variable using the implementation in the R-package “vegan”⁴⁸. We note that results need to be
351 interpreted cautiously, as the mantel test is generally poorly calibrated for spatially
352 autocorrelated data⁴⁹.

353

354 Samples omitted from model fitting

355 Besides samples whose geographic origin we could not unambiguously assign (n=74), we
356 removed a small number of samples that would violate some assumptions of the EEMS model.
357 In particular, we excluded all Jewish samples (n=379), due to complexity of the diaspora and
358 subsequent local admixture⁵⁰) and Han-Chinese in Taiwan and Singapore(n=170), who both are
359 recent migrant population to those locales. To avoid any possible distortion due to uneven
360 sampling, we downsampled all single locales to at most 50 individuals, drawn independently for
361 different panels. This resulted in a total of 6066 individuals used in at least one panel (Table
362 S1).

363 Visualization pipeline

364 We developed a second pipeline using snakemake³⁴ to perform all subsetting and demographic
365 analyses, available under github.com/NovembreLab/eems-around-the-world. The pipeline
366 allows for defining panels using a flexible set of features, latitudinal and longitudinal boundaries,
367 continent or country of samples, source study, as well as the addition and exclusion of particular
368 samples or populations. Based on these subsets, different modules allow performing EEMS and
369 PCA analyses, as well as generating all the figures, that were then annotated using inkscape.
370 All configuration variables are stored in json and yaml config files. We perform EEMS and PCA
371 for each panel independently. Structural variants are a potential confounding factor for genome-
372 wide SNP based analysis. In PCA, these variants may result in a number of neighboring SNP in
373 high LD to have very high loadings, thus overemphasizing the effect of these variants. For this
374 reason, it is advisable to remove regions containing SNP that have extremely high loadings on
375 some Principal component. Thus, for each panel, we perform a preliminary PCA analysis using
376 flashpca⁵¹. The loading-scores for each PC were normalized by dividing them by the standard
377 deviations on each PC [$\text{outlier_score} = L[i]/\text{sd}(L[i])$], and then we removed a 200kb window

378 around any SNP for which $|\text{outlier_score}| > 5$. We also dropped individuals with more than 5%
379 missingness, and SNPs with more than 1% missing data from each panel.

380 EEMS

381 To generate the map surfaces, we must choose a grid size and boundaries. Choosing a
382 coarse grid results in faster computation, but only produces a map with broad-scale patterns. A
383 finer grid, on the other hand, is able to reveal more details, but at a steep increase in
384 computational cost and with an increased danger of introducing patterns that are harder to
385 interpret. Grid density and sizes are given in Extended Data Table 1, along with population
386 level F_{ST} calculated using plink, and F_{ST} based on the mean migration rate inferred by eems and
387 equilibrium stepping stone model theory⁵².

388
389 We evaluated the impact of SNP ascertainment bias by running EEMS on the multiple,
390 documented SNP ascertainment panels of the Human Origins data²⁵. We found that while
391 ascertainment bias has an effect on the heterozygosity surfaces that EEMS estimates, the
392 migration surfaces remain relatively unaffected (Extended Data Fig. 1). Therefore, we restrict
393 our presentation to the migration surfaces.

394
395 For each panel, we performed four pilot runs of 2-8 million iterations each. The run with the
396 highest likelihood was then used for a second set of four runs of 4-10 million iteration each, with
397 the first 500,000 million discarded as burn-in. Number of iteration were chosen such that total
398 computation time was around 10 days. Every 20,000th iteration was sampled. EEMS
399 approximates a continuous region with a triangular grid, which has to be specified. We
400 generated global geodesic graphs at three resolutions (approximate distance between demes of
401 120, 240 and 500km, respectively) using dggrid v6.1⁵³ and intersected these graphs with the
402 area representing each panel (Extended Figures 2,3). All other (hyper-)parameters were kept at
403 their default values¹⁶. We compared EEMS to an isolation-by-distance model with a constant
404 migration rate by re-fitting EEMS allowing only a single migration rate tile, but arbitrary diversity
405 rate tiles using the otherwise same settings. The resulting log Bayes Factors are given in
406 Extended Data Table 2.

407 Evaluating fit of EEMS and PCA to genetic distances

408 For EEMS, the posterior samples imply an expected distance matrix between populations. For
409 PCA, the components and their loadings provide an approximation to the genetic distance
410 matrix between individuals. We use the median PCA values of individuals across ten PC
411 components to produce an expected genetic distance matrix between populations. We use ten
412 PC components as most investigators evaluate population structure based on only the first
413 several PCs. For each method the expected genetic distance matrices are compared to the
414 observed matrices using a simple linear correlation computed between all pairwise distances.

415 Acknowledgements

416 We are grateful for helpful comments from Choongwon Jeong, Matthew Stephens, Anna Di
417 Rienzo, Melinda A. Yang, Joshua G. Schraiber and members of the Novembre lab. This
418 research was supported by research grants NIH/NCI U01 CA198933 and NIH/NIGMS R01
419 GM108805 (J. N. and B. M. P) and by a Swiss National Science Foundation early postdoc
420 mobility fellowship (B. M. P.). We acknowledge the University of Chicago Research Computing
421 Center for support of this work. We dedicate the paper in memoriam of Brad McRae (1966-
422 2017) whose work on resistance distances underlies the EEMS methodology.

423

424 Author contributions.

425 B.M.P. analyzed data. B.M.P., D.P., and J.N. interpreted results. B.M.P and J.N conceived the
426 study and wrote the manuscript.

427 Competing financial interests

428 The authors declare no competing financial interests.

429 Correspondence to:

430 Benjamin Peter (benjamin_peter@eva.mpg.de), or John Novembre (jnovembre@uchicago.edu)

431 Supplementary Text on Regional Scale Analyses

432 Here we provide a more expanded discussion of the regional-scale results. To help identify
433 features that we discuss, we have added labels to discussed features in the figures, and refer to
434 them in the text here in parentheses. The labels are typically capitalized abbreviations and in
435 some cases are full words.

436

437 **Western Eurasia.** Europe appears largely homogeneous in the Afro-Eurasia panel, but a finer-
438 scale analysis (Western Eurasia panel, Fig. 2a; $n=2,049$; 122 locales, $F_{ST}=0.010$) reveals
439 abundant fine-scale structure: bodies of waters are consistently covered by lower effective
440 migration regions, with migration being lower in southern seas (Mediterranean, Adriatic, Black
441 Sea) relative to those in northern Europe (North Sea, Irish Sea, English Channel). Terrestrial
442 barriers are observed in: the Alps (and an adjacent region extending into Southern France),
443 surrounding the Mozabites in Tunisia and the Saami in Scandinavia, the western and northern
444 edges of the Arabian Desert (though we note the region has few samples). Troughs reflecting
445 historical domains are observed: between Germanic and Northern Slavic-speakers (CE,
446 Extended Data Figure 7) and between domains of Slavic-speakers and the Caucasus (NS).
447 Remaining regions are generally inferred to have average or above average migration, with one
448 obvious corridor being that between Iceland and Scandinavia (IS), presumably due to the recent
449 colonization of Iceland. One interesting feature is an area of East-West low migration between
450 the Italian peninsula and Greece (GI). A corridor between Crete and Sicily is inferred south of it,
451 and between mainland Greece and southern Italy north of it. This likely reflects a pattern of

452 close genetic similarity among coastal Mediterranean populations observed previously⁴³ but
453 suggests it may have north-south structure. Ancient DNA results suggest that the patterns we
454 observe are recent^{54,55} and have been shaped in the last 3,000-5,000 years with contributions
455 from multiple sources. Strikingly, proposed expansion routes through the Eurasian Steppe and
456 Levant into Europe partially align with corridors of high effective migration.

457
458 **Central/Eastern Eurasia.** The Central/Eastern Eurasia surface (Fig. 2e; $n=2,578$; 181 locales,
459 $F_{ST}=0.042$) is overall similar to the patterns seen in the AEA panel, with a trough through the
460 Himalayas/Tien-Shan but we observe a corridor from Mongolia to the Caspian Sea, possibly
461 reflecting the Mongol empire, as the Kalmyk, Kazhaks, and Uygurs all have well documented
462 Mongolian-like genetic ancestry. The presence of this corridor depends on a small number of
463 samples; if Uygurs and Kalmyks are removed, we find a pattern similar to that in the AEA panel
464 in that region.

465
466 Where the global analysis did not reveal any strong patterns in South Asia, at the higher
467 resolution we observe troughs in the Indian subcontinent between a central Indian region (CI) of
468 mainly Indo-Aryan languages and an Eastern and Southern region with two Austroasiatic
469 speaking (Kharia, Ho), and Dravidian speaking populations, but this trough is not significantly
470 correlated with linguistic group (Extended Data Figure 7d) We also find that the Himalayas
471 generate a trough between India and Tibet, but the Kusunda population adds a corridor there,
472 which is explained by the fact that they have both Tibetan and Indian ancestry⁵⁶.

473
474 In East Asia, we observe marine troughs in the East China Sea, strait of Tartary and the
475 Andaman Sea (Onge). Terrestrially, we observe troughs between coastal China (CC), a central
476 region with several Tibeto-Burman samples (TB, along with the Tu who speak a Mongolic
477 language, and have been suggested to have received European admixture 1,200y ago²⁴), and a
478 western region anchored by Tibetan samples. The coastal Chinese region contains a high-
479 effective-migration area that extends into Korea and Japan.

480
481 Overall, the Central/East Asia panel is particularly complex with one of the lowest levels of r^2
482 between EEMS expected genetic distances and the observed distances ($r^2 = 0.66$, Extended
483 Data Fig. 5), and the residuals are very high (Extended Data Fig. 8). This is expected as the
484 relatively open steppe has been the site of repeated long-range population movements and
485 invasions, by e.g. Bronze Age Steppe populations, Mongols and Turkic speakers, that we
486 expect are difficult to depict using the model of steady-state gene flow model fit by EEMS.

487
488 **South-East Asia.** In the South-East Asian panel ($n=1,054$, 58 locales; $F_{ST}=0.037$; Fig. 2k)
489 troughs align with the many seas and channels in this region: the South-Chinese Sea (SCS),
490 the waterway running east of the Philippines (PH) and Sulawesi south to the Flores Sea (SEP),
491 the waterway between western New Guinea into the Banda Sea (BS), the Malacca strait
492 between Sumatra and Malaysia (SM), the Sunda Strait between Java and Sumatra (JS), the
493 Java Sea between Bali and Java (BJ), as well as the Makassar strait and Celebes Sea between
494 Borneo and Sulawesi (EB). Two corridors, one from Taiwan/Luzon through Western Mindanao
495 to Sulawesi, and one from Ternate through the Lower Sunda Islands (LSI) into Melanesia

496 possibly reflect the Austronesian expansion that started roughly 3,000 years ago⁵⁷. On the
497 mainland, we find low effective migration north of Bangkok (BAN) and near samples from
498 Northern Thailand (TH) (including the Southern Chinese Wa and Jinuo samples (SC)). These
499 two samples have low inferred effective migration with South-Eastern Chinese samples (SEC).
500 Two Malay Negrito samples in Northern Malaysia (MN) are placed in a trough, revealing their
501 genetic distance to other South-East Asian populations also apparent on PC2 (Figure 2l).
502

503 **Africa.** In Africa (AFR, $n=749$, 71 locales, $F_{ST}=0.055$; Fig. 2g) two troughs corresponding to the
504 Sahara desert and Mozambique Channel are observed. In Northern Africa (NAA), we see a
505 small trough of low effective migration separating two latitudinal corridors; one following the
506 Mediterranean coast and one inland (Fig 2g). The inland corridor disappears in our lower-
507 resolution Afro-Eurasia panel (Figure 1a) and presumably reflects Sub-Saharan ancestry in
508 some Moroccans, perhaps through trans-Saharan trade.
509

510 In the AFR-panel, we observe a trough reflecting the language group boundaries between
511 Niger-Congo and Afro-Asiatic language speakers⁵⁸ (Extended Data Figure 7), with the West-
512 African Afro-Asiatic speaking Hausa and Mada being placed in a barrier together with the
513 admixed Fulani⁵⁹. West Africa appears as a high-gene-flow region (WA), and two corridors
514 pass from Nigeria - one along the coast of Congo (CO) southwards and another further east
515 (EC) connecting to Kenya and Tanzania. In both Central and Eastern Africa Nilo-Saharan and
516 Niger-Congolese speakers overlap, resulting in low effective migration uncorrelated with langu:
517 the Nilo-Saharan Kaba, Dinka and Bulala are in a region of high gene flow, separated by a
518 trough from the Biaka and Mbuti Pygmies. Southern and Eastern Africans are separated by low
519 effective migration through Mozambique and South-Western Tanzania (SWT).
520

521 **Southern Africa** Patterns in southern Africa are complex, with a troughs separating out
522 Western Bantu speaking populations, and the KhoeSan Khwe and Xuun. Stratifying Southern
523 Africans into a KhoeSan (SAKS, $n=109$, 16 locales, $F_{ST}=0.025$, Fig. 2k) and Bantu speakers
524 (SAB, $n=30$, 11 locales, $F_{ST}=0.014$; Fig. 2l) reveals very different spatial structure. The Bantu
525 speakers are separated by a single barrier into an eastern and western location. For the South
526 African Hunter-Gatherers most samples fall into a central region with high effective migration,
527 including the Taa, Naro and Hoan (TNH). Troughs in the North separate this region from the
528 Sua and Tswa (ST) and in the south-west from the Khomani and Nama (Nama), respectively.
529 The remaining samples fall either into a Northern high migration area (Khwe and Xuun, KX) or a
530 North-Western low migration area (Damara and Haiom, DH). These results are broadly
531 consistent with existing work on African population structure⁵⁹⁻⁶², and emphasize African
532 population structure appears largely determined by the Sahara desert, the Bantu and Arabic
533 expansions, and the complex structure of hunter-gatherer groups specifically in South Africa.
534

535

536 **Extended Data**

537

Study	Abbrev.	Samples	Locations	Source
Bryc et al. 2009	B09	121	11	Ref ⁴⁰
Behar et al. 2010	Be10	295	22	Ref ⁵⁰
Behar et al. 2013	B13	131	20	Ref ⁶³
Bigham et al. 2010	Bi10	45	3	Ref ⁴⁵
Chaubey et al. 2011	C11	37	5	Ref ⁶⁴
Cardona et al. 2014	C14	192	37	Ref ⁶⁵
Di Cristofaro et al. 2013	D13	14	3	Ref ⁶⁶
Fedorova et al. 2013	F13	30	6	Ref ⁶⁷
HUGO Consortium 2009	H09	975	47	Ref ³⁷
Hunter-Zinck et al. 2010	H10	85	1	Ref ⁴¹
Jeong et al. 2017	J17	53	2	Ref ⁴⁴
Kovacevic et al. 2014	K14	70	6	Ref ⁶⁸
Lazaridis et al. 2014	L14	1590	159	Ref ⁶⁹
Metspalu et al. 2011	M11	127	11	Ref ⁷⁰
Migliano et al. 2013	M13	68	6	Ref ⁷¹
Nelson et al. 2008	N08	531	29	Ref ³⁹
Paschou et al. 2014	Pa14	626	29	Ref ⁴³
Pierron et al. 2014	Pi14	114	5	Ref ⁷²
Raghavan et al. 2014	R14	83	9	Ref ⁷³
Rasmussen et al. 2010	R10	101	9	Ref ⁷⁴
Rasmussen et al. 2011	Ra11	19	3	Ref ⁵⁶
Reich et al. 2011	Re11	106	16	Ref ⁷⁵
Skoglund et al. 2014	S14	15	1	Ref ⁷⁶
Xing et al. 2010	X10	92	4	Ref ³⁸
Xu et al. 2011	X11	28	3	Ref ⁴⁶
Yunusbayev et al. 2012	Y12	183	14	Ref ⁷⁷
Yunusbayev et al. 2015	Y15	299	42	Ref ³⁶

538 **Extended Data Table 1:** Data Sources. Abbrev: Abbreviation; Ind: total number of individuals;
 539 Loc. Number of unique sample locations

540

541

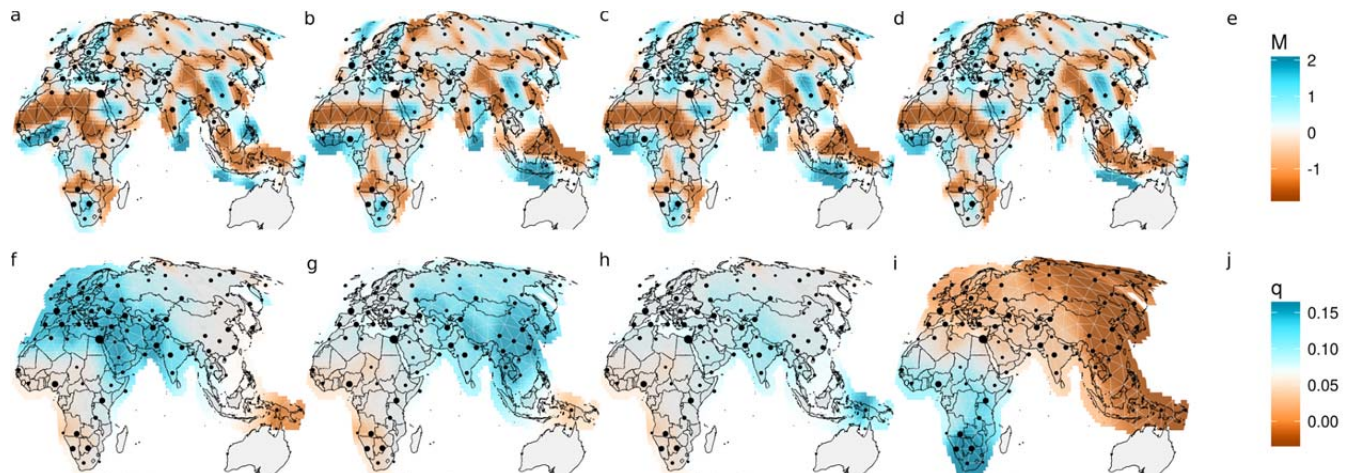
Panel	Abb.	Ind.	Locations	SNPs	Grid Size (# of demes)	Resolution (km)	F_{ST}	Model-fit F_{ST} (adjacent demes)	Model-fit F_{ST} (500km)	Support (log-BF)
Afro-Eurasia	AEA	4697	370	19972	686	500	0.071	0.99%	0.99%	254,472
Central/Eastern Eurasia	CEA	2578	181	21045	1147	240	0.042	0.22%	0.42%	129,035
Western Eurasia	WEA	2049	122	26438	1437	120	0.010	0.75%	1.08%	46,210
South-East Asia	SEA	1054	58	7553	1388	120	0.037	0.29%	0.56%	13,654
Africa	AFR	749	71	20984	694	240	0.055	0.81%	1.18%	51,771
Southern Africa	SAKS	109	16	532343	227	120	0.025	0.32%	0.62%	2298
KhoeSan										
Southern Africa	SAB	30	11	65095	227	120	0.014	0.26%	0.56%	126
Bantu										

542

543 **Extended Data Table 2:** Analysis Panels. Abb. Panel Abbreviation. Res. Avg. distance between grid
544 points (in km) ; Support: log Bayes factor in favor of complex vs constant migration model. Implied F_{ST}
545 between adjacent demes based on posterior mean migration rates. Equation 19a from ⁵² is used to
546 calculate implied F_{ST} using a torus approximation: For F_{ST} (adjacent demes): $F_{ST}=(1+32m/S(d))^{-1}$ where
547 $S(d)$ is a function of the distance between demes and given by equation A12 in ⁵². In the first column, we
548 use $S(1)$, in the second $S(4)$ for highest and $S(2)$ for medium resolution panels to get F_{ST} for demes at
549 the lowest resolution (~500km).

550

551



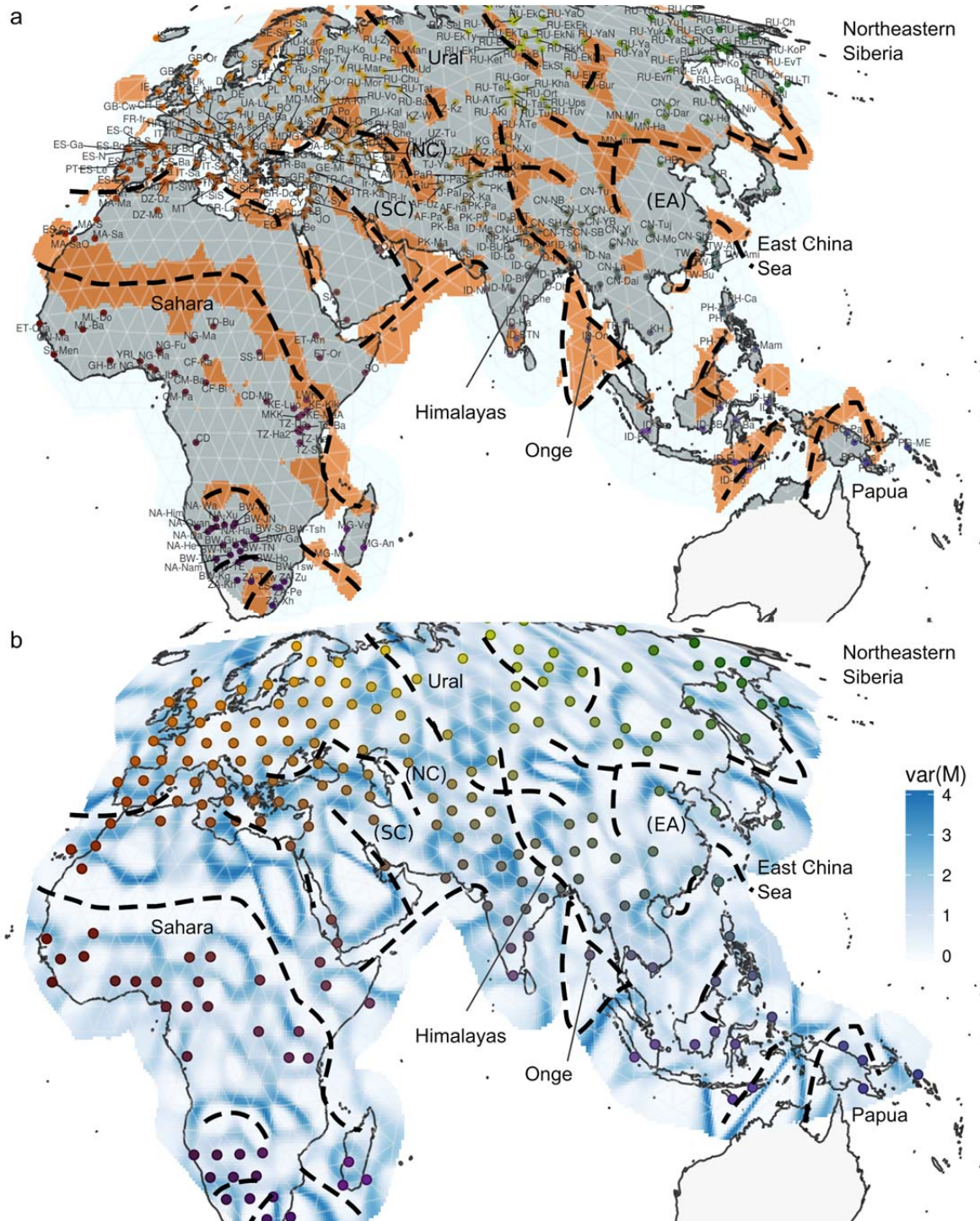
552

553

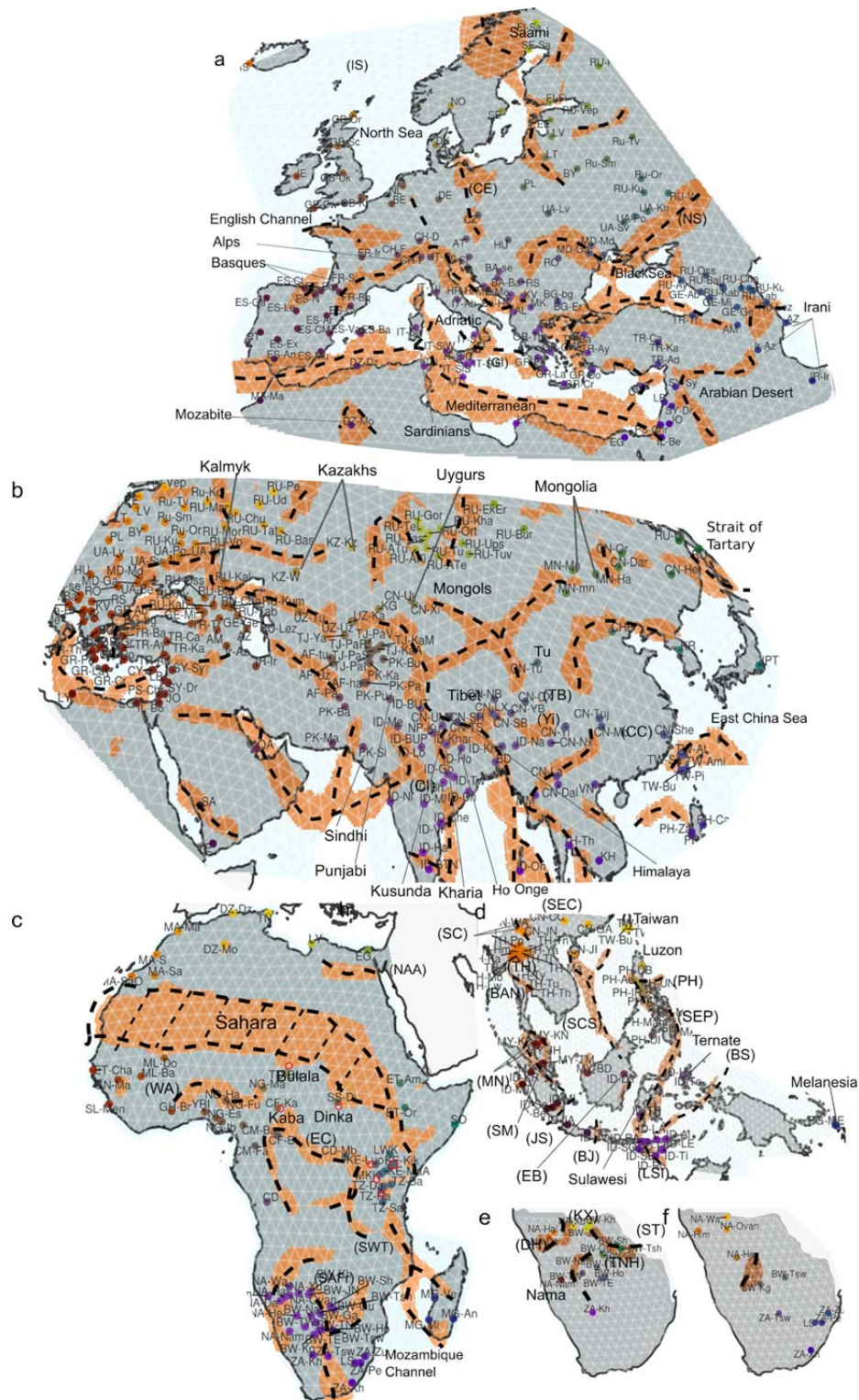
554 **Extended Data Figure 1:** Ascertainment bias. We run EEMS only using the Human Origin data
555 ²⁵, using SNPs ascertained in a French (a/f), Chinese (b/g), Papuan (c/h) and San(d/i)
556 individual. Migration rate surfaces (a-d) remain robust, whereas the within-deme diversity
557 surfaces (f-i) show highest diversity at the respective ascertainment location. e/j: scale bars for
558 migration rates and within-deme diversity rate parameters, respectively.

559

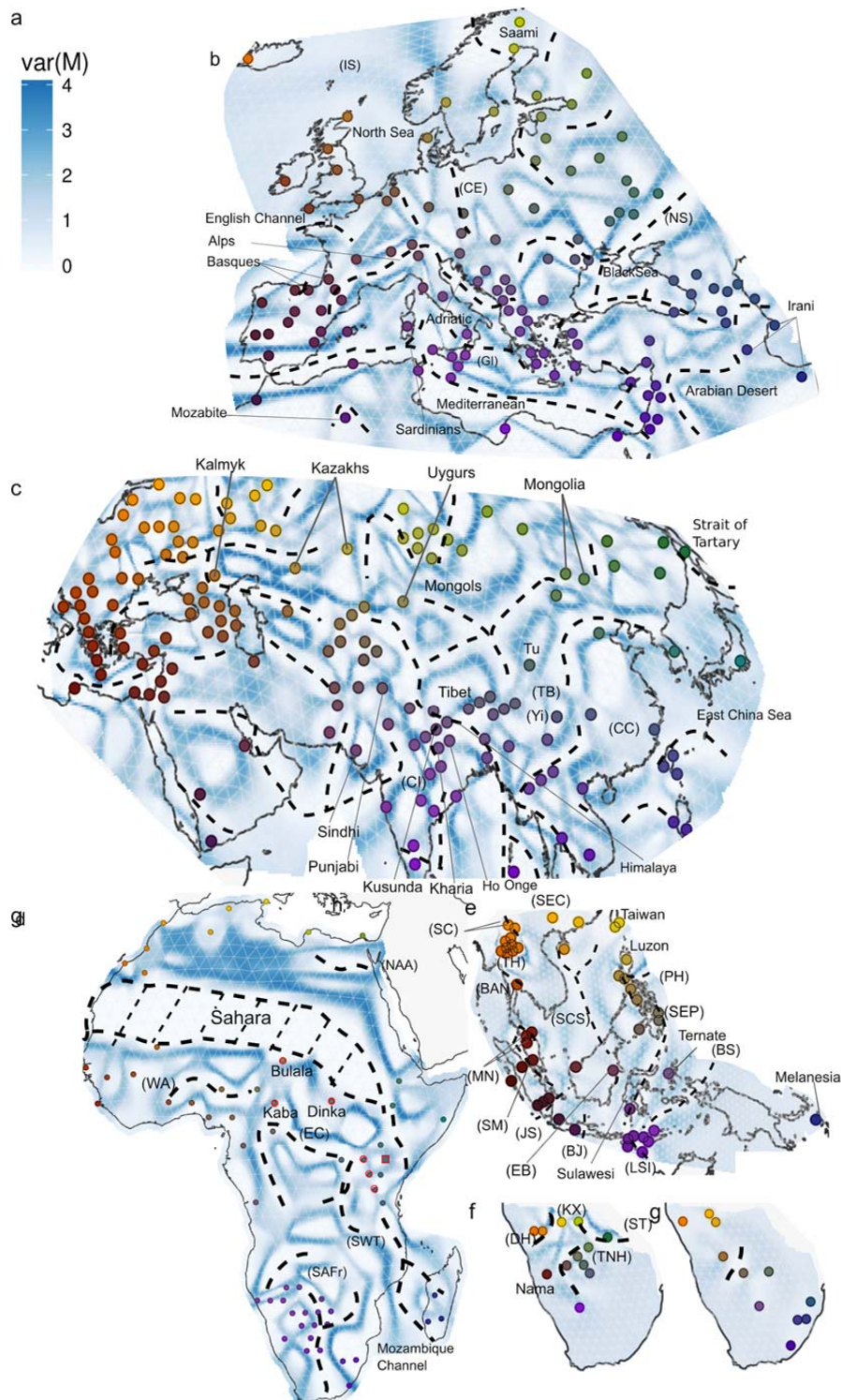
560



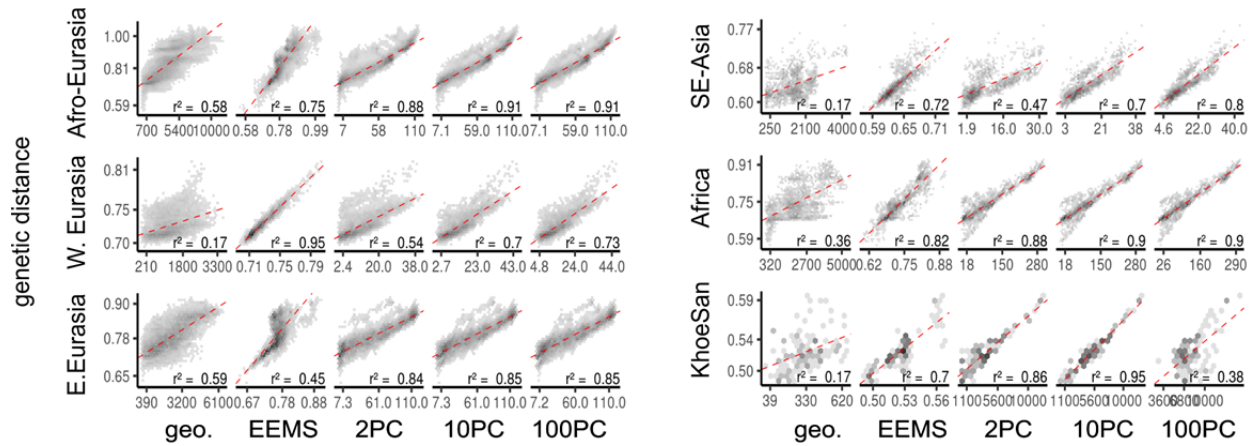
561
562 **Extended Data Figure 2:** **a:** Location of troughs (below average migration rate in more than 95% of
563 MCMC iterations) are given in brown. Sample locations and EEMS grid are displayed. **b:** Posterior
564 variance on migration rate parameters. Note that most significant features are in low variance regions, but
565 that they are often surrounded by high-variance regions, implying the exact boundary of troughs is
566 estimated with uncertainty. Grid-fitted sample locations are displayed. Annotation in both panels is
567 identical to Figure 1a.



568
 569 **Extended Data Figure 3:** Location of troughs (below average migration rate in more than 95%
 570 of MCMC iterations) are given in brown. Sample locations and EEMS grid are displayed for a:
 571 WEA b: CEA c: AFR d: SAHG and e: SEA analysis panels. Annotation in all panels is identical
 572 to Figure 2.

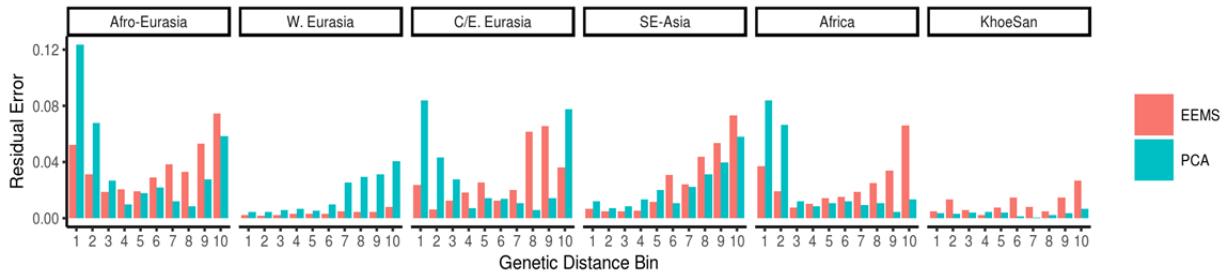


573
 574 **Extended Data Figure 4:** Posterior variances in migration rate parameters. Grid-fitted sample
 575 locations are displayed. **a:** scale bar **b:** WEA **c:** CEA **d:** AFR **e:** SAHG and **f:** SEA analysis
 576 panels. Note that most significant features are in low variance regions, but that they are often
 577 surrounded by high-variance regions, implying the exact boundary of troughs is estimated with
 578 uncertainty. Annotation of troughs and select features is identical to Figure 2.



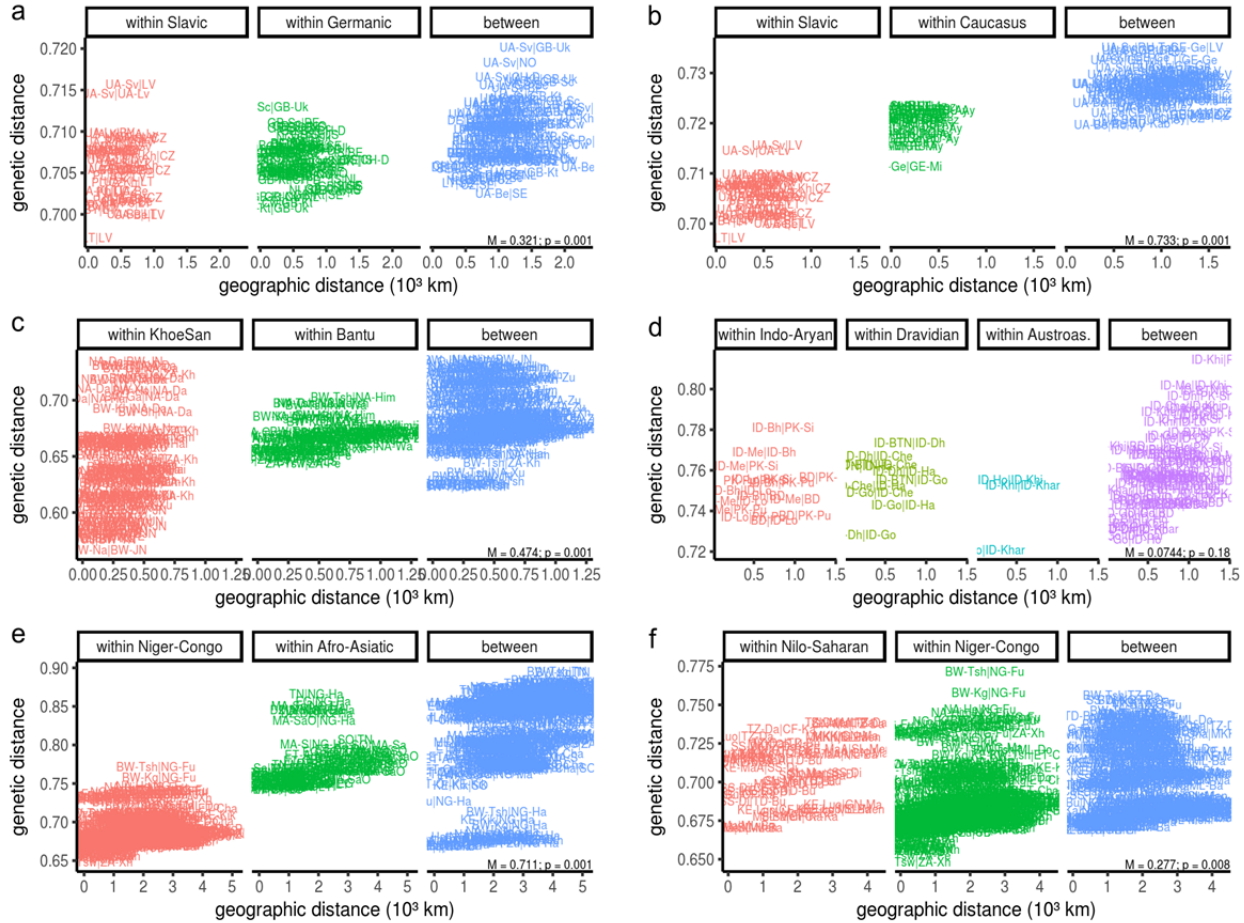
579
580
581
582
583
584

Extended Data Figure 5: Hex-binned scatterplots of genetic distance versus geographic distance (in km), predicted distance via EEMS model fit, and predicted distance via a ten-component PCA, for all panels. Darker areas correspond to bins with more points. The fit of a simple linear regression (red dashed lines) and r^2 are given.



585
586
587
588
589
590
591
592
593

Extended Data Figure 6: Comparing Fit of PCA and EEMS. We show the relative error of EEMS (red) and PCA (blue, first 10 PCs) for all pairs, stratified by genetic distance. For each panel, all pairwise genetic distances were distributed in ten bins of equal size, for which we then computed the median absolute error of the fitted model vs the observed distances. For W. Eurasia and SE-Asia, EEMS fits uniformly better than PCA. In the Afro-Eurasian, Central/Eastern Eurasian and African panel, EEMS fits better for smaller distances, but the fit is worse for larger distances. For the KhoeSan, EEMS fits worse than PCA for all distance bins.



594

595

596 **Extended Data Figure 7: Genetic vs. geographic distance within and between language**

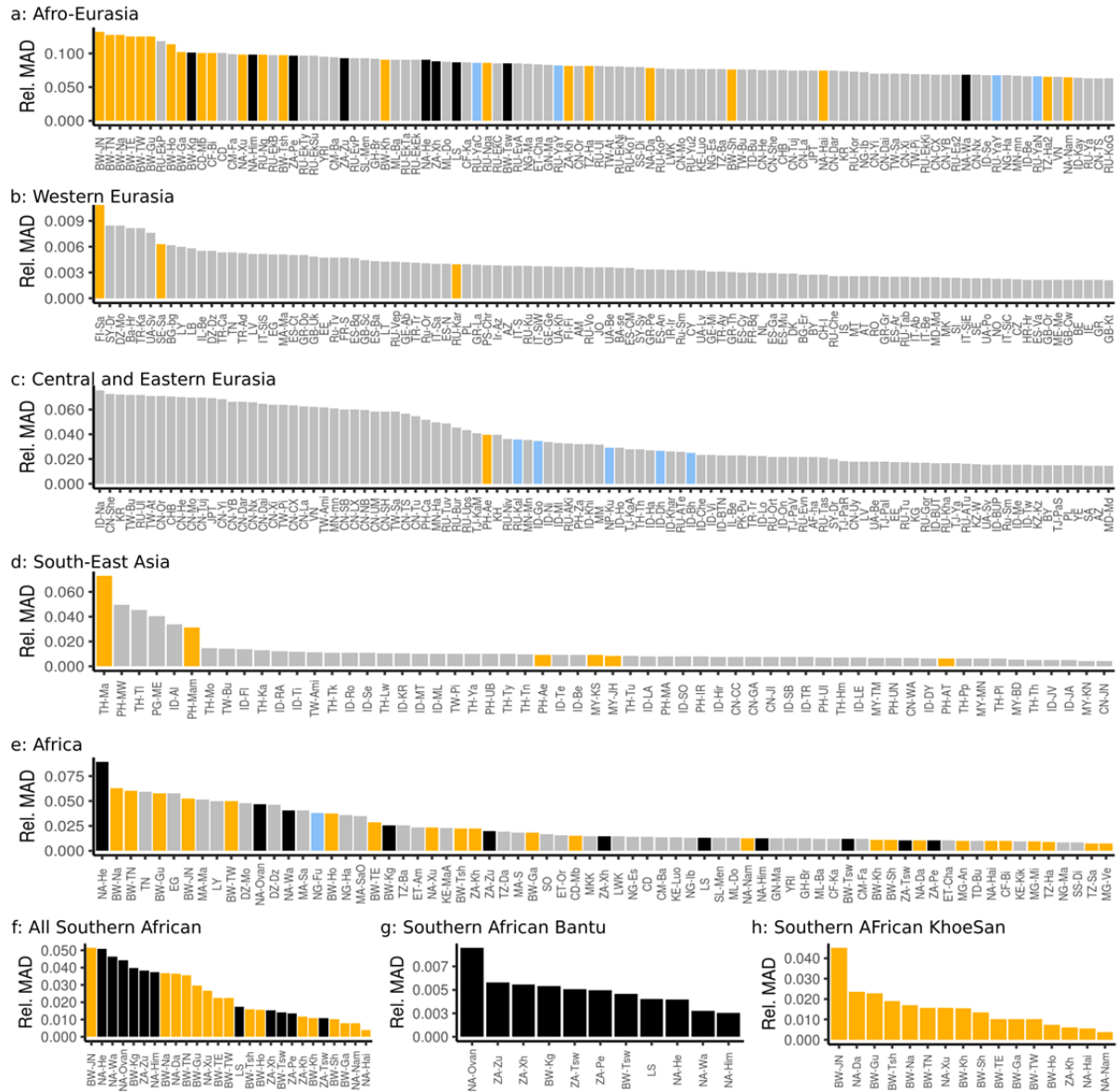
597 **groups.** The eems-plots revealed several troughs aligning with differences in linguistic groups.

598 We show the pairwise relationship of genetic and geographic differences within- and between

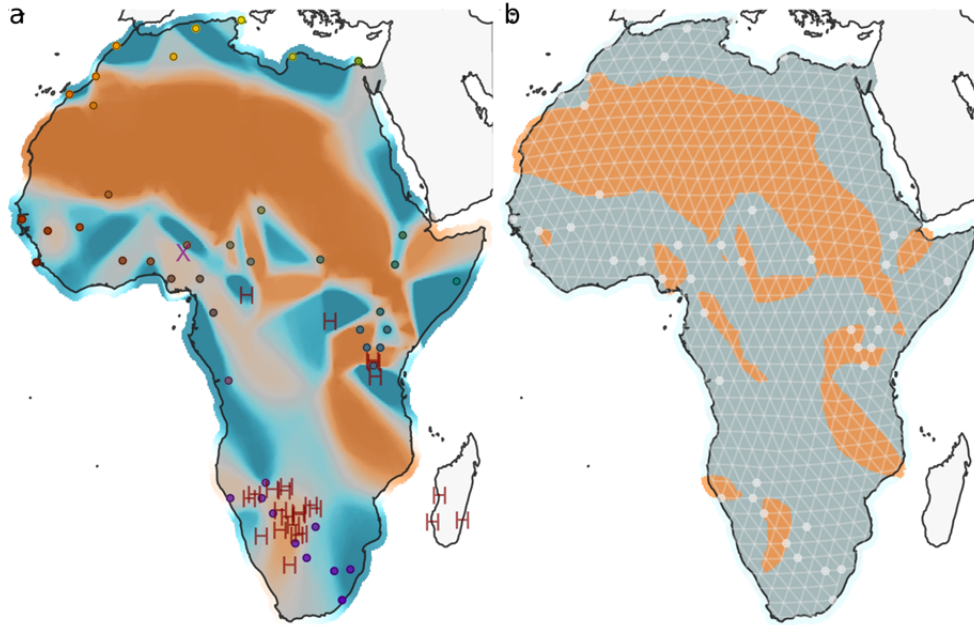
599 adjacent language groups mentioned in the main text for a. Slavic and Germanic speakers

600 (Southern Africa) d. Indo-Aryan, Dravidian and Austroasiatic (CEA) e. Niger-Congo and Afro

601 Asiatic (AFR) and f. Nilo-Saharan and Niger-Congo (AFR).



602
 603 **Extended Data Figure 8: EEMS-fit residuals.** For each population, we show the median
 604 absolute deviation (MAD) of the observed vs EEMS-fitted genetic distances, normalized by the
 605 median distance for this population. yellow: Hunter-Gatherers; Black: Southern African Bantu
 606 speakers; Blue: Populations with a recent admixture or displacement.



607
608 **Extended Data Figure 9: Alternative Africa analysis.** To assess the effect of populations that
609 may not be modelled well by EEMS (admixed or hunter-gatherer populations), we provide
610 supplemental analyses of Africa with several populations excluded from the model fit. **a:** EEMS-
611 map and **b:** location of troughs for Africa. Excluded populations are annotated with H (Hunter-
612 gatherers) and X (admixed). With this filtering (in particular removing the Hadza and Sandawe),
613 the Eastern African trough between Afro-Asiatic speakers and Nilo-Saharan / Niger-Congo
614 speakers (seen in Figures 1 and 2g) vanishes.

615
616
617
618

619 Additional References

- 620 34. Köster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics* **28**, 2520–2522 (2012).
- 621 35. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
- 622 36. Yunusbayev, B. *et al.* The Genetic Legacy of the Expansion of Turkic-Speaking Nomads across Eurasia. *PLoS Genet.* **11**,
- 623 e1005068 (2015).
- 624 37. HUGO Pan-Asian SNP Consortium. Mapping human genetic diversity in Asia. *Science* **326**, 1541–1545 (2009).
- 625 38. Xing, J. *et al.* Toward a more uniform sampling of human genetic diversity: a survey of worldwide populations by high-density
- 626 genotyping. *Genomics* **96**, 199–210 (2010).
- 627 39. Nelson, M. R. *et al.* The Population Reference Sample, POPRES: a resource for population, disease, and pharmacological
- 628 genetics research. *Am. J. Hum. Genet.* **83**, 347–358 (2008).
- 629 40. Bryc, K. *et al.* Genome-Wide Patterns of Population Structure and Admixture in West Africans and African Americans. *Proc.*
- 630 *Natl. Acad. Sci. U. S. A.* (2009). doi:10.1073/pnas.0909559107
- 631 41. Hunter-Zinck, H. *et al.* Population genetic structure of the people of Qatar. *Am. J. Hum. Genet.* **87**, 17–25 (2010).
- 632 42. Reich, D. *et al.* Denisova admixture and the first modern human dispersals into Southeast Asia and Oceania. *Am. J. Hum.*
- 633 *Genet.* (2011).
- 634 43. Paschou, P. *et al.* Maritime route of colonization of Europe. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 9211–9216 (2014).
- 635 44. Jeong, C. *et al.* A longitudinal cline characterizes the genetic structure of human populations in the Tibetan plateau. *PLoS One*
- 636 **12**, e0175885 (2017).
- 637 45. Bigham, A. *et al.* Identifying Signatures of Natural Selection in Tibetan and Andean Populations Using Dense Genome Scan
- 638 Data. *PLoS Genet.* **6**, e1001116 (2010).
- 639 46. Xu, S. *et al.* A genome-wide search for signals of high-altitude adaptation in Tibetans. *Mol. Biol. Evol.* **28**, 1003–1011 (2011).
- 640 47. Hammarström, H., Bank, S., Forkel, R. & Haspelmath, M. Glottolog 3.2. (2018).
- 641 48. Oksanen, J. *et al.* The vegan package. *Community ecology package* **10**, 631–637 (2007).
- 642 49. Guillot, G. & Rousset, F. Dismantling the Mantel tests. *Methods Ecol. Evol.* **4**, 336–344 (2013).
- 643 50. Behar, D. M. *et al.* The genome-wide structure of the Jewish people. *Nature* **466**, 238–242 (2010).
- 644 51. Abraham, G. & Inouye, M. Fast principal component analysis of large-scale genome-wide data. *PLoS One* **9**, e93766 (2014).
- 645 52. Slatkin, M. Inbreeding coefficients and coalescence times. *Genetic Research* **58**, 167–175 (1991).
- 646 53. Sahr, K., White, D. & Kimerling, A. J. Geodesic Discrete Global Grid Systems. *Cartogr. Geogr. Inf. Sci.* **30**, 121–134 (2003).
- 647 54. Haak, W. *et al.* Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature* **522**, 207–211
- 648 (2015).
- 649 55. Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature* **522**, 167–172 (2015).
- 650 56. Rasmussen, M. *et al.* An Aboriginal Australian Genome Reveals Separate Human Dispersals into Asia. *Science* **334**, 94–98
- 651 (2011).
- 652 57. Duggan, A. T. & Stoneking, M. Recent developments in the genetic history of East Asia and Oceania. *Curr. Opin. Genet. Dev.*
- 653 **29**, 9–14 (2014).

- 654 58. Campbell, M. C. & Tishkoff, S. A. African genetic diversity: implications for human demographic history, modern human origins,
655 and complex disease mapping. *Annu. Rev. Genomics Hum. Genet.* **9**, 403–433 (2008).
- 656 59. Bryc, K. *et al.* Genome-wide patterns of population structure and admixture in West Africans and African Americans. *Proc. Natl.*
657 *Acad. Sci. U. S. A.* **107**, 786–791 (2010).
- 658 60. Tishkoff, S. A. *et al.* The genetic structure and history of Africans and African Americans. *Science* **324**, 1035–1044 (2009).
- 659 61. Pickrell, J. K. *et al.* The genetic prehistory of southern Africa. *Nat. Commun.* **3**, 1143 (2012).
- 660 62. Uren, C. *et al.* Fine-Scale Human Population Structure in Southern Africa Reflects Ecogeographic Boundaries. *Genetics* **204**,
661 303–314 (2016).
- 662 63. Behar, D. M. *et al.* No evidence from genome-wide data of a Khazar origin for the Ashkenazi Jews. *Hum. Biol.* **85**, 859–900
663 (2013).
- 664 64. Chaubey, G. *et al.* Population Genetic Structure in Indian Austroasiatic Speakers: The Role of Landscape Barriers and Sex-
665 Specific Admixture. *Mol. Biol. Evol.* **28**, 1013–1024 (2011).
- 666 65. Cardona, A. *et al.* Genome-Wide Analysis of Cold Adaptation in Indigenous Siberian Populations. *PLoS One* **9**, e98076 (2014).
- 667 66. Di Cristofaro, J. *et al.* Afghan Hindu Kush: where Eurasian sub-continent gene flows converge. *PLoS One* **8**, e76748 (2013).
- 668 67. Fedorova, S. A. *et al.* Autosomal and uniparental portraits of the native populations of Sakha (Yakutia): implications for the
669 peopling of Northeast Eurasia. *BMC Evol. Biol.* **13**, 127 (2013).
- 670 68. Kovacevic, L. *et al.* Standing at the Gateway to Europe - The Genetic Structure of Western Balkan Populations Based on
671 Autosomal and Haploid Markers. *PLoS One* **9**, e105090 (2014).
- 672 69. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature* **513**, 409–
673 413 (2014).
- 674 70. Metspalu, M. *et al.* Shared and Unique Components of Human Population Structure and Genome-Wide Signals of Positive
675 Selection in South Asia. *Am. J. Hum. Genet.* **89**, 731–744 (2011).
- 676 71. Migliano, A. *et al.* Evolution of the Pygmy Phenotype: Evidence of Positive Selection from Genome-wide Scans in African,
677 Asian, and Melanesian Pygmies. *Hum. Biol.* **85**, (2013).
- 678 72. Pierron, D. *et al.* Genome-wide evidence of Austronesian–Bantu admixture and cultural reversion in a hunter-gatherer group of
679 Madagascar. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 936–941 (2014).
- 680 73. Raghavan, M. *et al.* Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans. *Nature* **505**, 87–91 (2014).
- 681 74. Rasmussen, M. *et al.* Ancient human genome sequence of an extinct Palaeo-Eskimo. *Nature* **463**, 757–762 (2010).
- 682 75. Reich, D. *et al.* Denisova Admixture and the First Modern Human Dispersals into Southeast Asia and Oceania. *Am. J. Hum.*
683 *Genet.* **89**, 516–528 (2011).
- 684 76. Skoglund, P. *et al.* Genomic Diversity and Admixture Differs for Stone-Age Scandinavian Foragers and Farmers. *Science* **344**,
685 747–750 (2014).
- 686 77. Yunusbayev, B. *et al.* The Caucasus as an Asymmetric Semipermeable Barrier to Ancient Human Migrations. *Mol. Biol. Evol.*
687 **29**, 359–365 (2012).

688