

1 **Population histories of the United States revealed through fine-scale migration and**
2 **haplotype analysis**

3

4 Chengzhen L. Dai¹, Mohammad M. Vazifeh², Chen-Hsiang Yeang³, Remi Tachet², R. Spencer
5 Wells⁴, Miguel G. Vilar⁵, Mark J. Daly^{6,7,8,9}, Carlo Ratti^{2*}, Alicia R. Martin^{7,8,9*†}

6

7 ¹ Department of Electrical Engineering and Computer Science, Massachusetts Institute of
8 Technology, Cambridge, MA 02139, USA

9 ² Senseable City Lab, Massachusetts Institute of Technology, Cambridge, MA 02139, USA

10 ³ Institute of Statistical Science, Academia Sinica, Nankang, Taipei, Taiwan

11 ⁴ Insitome, Inc, Austin, TX 78701, USA

12 ⁵ Genographic Project, National Geographic Society, Washington, DC 20036, USA

13 ⁶ Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland

14 ⁷ Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA 02114,
15 USA

16 ⁸ Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge,
17 MA 02142, USA

18 ⁹ Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA
19 02142, USA

20

21 *These authors jointly supervised this work

22 †Corresponding author: armartin@broadinstitute.org

23 **Abstract**

24

25 The population of the United States is shaped by centuries of migration, isolation, growth, and
26 admixture between ancestors of global origins. Here, we assemble a comprehensive view of
27 recent population history by studying the ancestry and population structure of over 32,000
28 individuals in the US using genetic, ancestral birth origin, and geographic data from the National
29 Geographic Genographic Project. We identify migration routes and barriers that reflect historical
30 demographic events. We also uncover the spatial patterns of relatedness in subpopulations
31 through the combination of haplotype clustering, ancestral birth origin analysis, and local
32 ancestry inference. These patterns include substantial substructure and heterogeneity in
33 Hispanics/Latinos, isolation-by-distance in African Americans, elevated levels of relatedness
34 and homozygosity in Asian immigrants, and fine-scale structure in European descents. Taken
35 together, our results provide detailed insights into the genetic structure and demographic history
36 of the diverse US population.

37

38 **Keywords:** population genetics, human history, human genomics, USA

39 **Main Text:** 3,765 words (excluding Methods, References, and Figures), 6 figures, 1 table

40 Introduction

41
42 The United States population is a diverse collection of global ancestries shaped by migration
43 from distant continents and admixture of migrants and Native Americans. Throughout the past
44 few centuries, continuous migration and gene flow have played major roles in shaping the
45 diversity of the US. Mixing between groups that have historically been genetically and spatially
46 distinct have resulted in individuals with complex ancestries while within-country migration have
47 led to genetic differentiation.¹⁻⁶

48
49 Previous genetics studies of the US population have sought to disentangle the relationship
50 between the genetic ancestry and population history of African Americans, European
51 Americans, and Hispanics/Latinos. In African Americans, proportions of African, European, and
52 Native American ancestry vary across the country and reflect migration routes, slavery, and
53 patterns of segregation between states.^{2,3,7} European American ancestry is characterized by
54 both mixing between different European populations as well as admixture with non-European
55 populations.^{6,8,9} Isolation and expansions in certain European population have also resulted in
56 founder effects.^{10,11} The mixing of European settlers with Native Americans have contributed to
57 large variations in the admixture proportions of different Hispanic/Latino populations.^{1,4,5} Among
58 Hispanics/Latinos, Mexicans and Central Americans carry more Native American ancestry;
59 Puerto Ricans and Dominicans have higher African ancestry; and Cubans have strong
60 European ancestry.^{1,4} Although much effort has been made to understand the genetic diversity
61 in the US, fine-scale patterns of demography, migration, isolation, and founder effects are still
62 being uncovered with the growing scale of genetic data, particularly for Latin American and
63 African descendants with complex admixture history.^{12,13} At the same time, there has been little
64 research on the population structure of individuals with East Asian, South Asian, and Middle
65 Eastern ancestry in the US.

66
67 In addition to being of anthropological interest, understanding fine-scale human history and its
68 role in shaping genetic variation is also important for interpreting the genetic basis of biomedical
69 traits. Currently, these roles are best understood in European populations due to Eurocentric
70 biases in studies.^{14,15} Consequently, translational interpretability gaps are evident in non-
71 European populations: more variants of unknown significance are identified via genetic testing;¹⁶
72 polygenic risk scores for complex disease risks are much less accurate;^{15,17} and false positive
73 genetic misdiagnoses are more common.¹⁸ Thus, studies of diverse, heterogeneous populations

74 offer substantial value to both our understanding of population history and biomedical
75 outcomes.¹⁹

76
77 In this study, we comprehensively explore the population structure and migration history of over
78 32,000 genotyped individuals in the US who partook in the National Geographic Genographic
79 Project, a not-for-profit public participation research initiative to study human migration history.²⁰
80 Here, we identify patterns of genetic ancestry and haplotype sharing among the project
81 participants. We combine these patterns with ancestral birth origin records and geographic
82 information to uncover recent demographic and migration trends. Taken together, we provide
83 insights into the ancestral origins and complex population histories in the US.

84

85

86 **Results**

87

88 **Genetic ancestry and diversity across the United States**

89 To assess the diversity of ancestries among individuals in the Genographic Project, we first
90 performed PCA and ADMIXTURE analysis (**Figure 1A-C; Figure S1-S2**).^{21,22} Since self-
91 reported ancestry does not always reflect genetic ancestry, we objectively assigned continental
92 ancestry to each Genographic sample using the 1000 Genomes Project data as reference
93 populations (**Methods and Materials**). We first trained a Random Forest classifier on the first
94 10 principal components (PCs) of the 1000 Genome Project samples with super population
95 classifications as ancestry labels (EUR = European, AMR = Admixed American, AFR = African,
96 EAS = East Asian, SAS = South Asian). We then used the trained model to assigned continent
97 ancestry to each individual in the Genographic cohort at 90% confidence. A total of 3,028
98 individuals (9.3% of total) did not meet the classification threshold, although many have
99 ancestry patterns similar to other European individuals (**Figure 1C; Table S1**). The inability to
100 classify these individuals may be due to the complex and variable admixture profiles of certain
101 populations such as Hispanics/Latinos.

102

103 Regional differences in genetic ancestry proportions correspond to historical demographic
104 trends. We evaluated the admixture proportions of classified individuals across the four
105 designated US Census regions: South, Northeast, Midwest, and West (**Figure 1C; Figure S2**).
106 Individuals of European descent make up the majority (78.5%) of the Genographic cohort and
107 are the most prevalent in the Midwest (82.8% of individuals in the Midwest; $P < 0.01$, Fisher's

108 exact test; **Table S1**). Individuals classified as having African ancestry are most common in the
109 South (3.2%), followed by the Northeast (3.0%). Individuals of Native American ancestry are
110 most prominent in the West and South (9.7% and 7.8% of total individuals in the West and
111 South, respectively; $P < 0.05$, Fisher's exact test). East Asians mostly reside in the West (2.1%),
112 while South Asians are most abundant in the Northeast (1.0%).

113

114 To uncover population substructure, we performed dimensionality reduction with Uniform
115 Manifold Approximation and Projection (UMAP) on the first 20 PCs of a combined Genographic
116 and 1000 Genomes Project dataset.^{23,24} By leveraging multiple PCs at once, UMAP can
117 disentangle subcontinental structure (**Figure 1D-E; Figure S3-S4**). Similar to previous
118 analysis,²⁴ populations in the 1000 Genomes Project form distinct clusters corresponding to
119 ancestry and geography. The Genographic individuals project into several clusters, overlapping
120 with the 1000 Genomes Project clusters. Consistent with the PCA and ADMIXTURE analysis,
121 the largest clusters correspond to European ancestry and cluster closely with the 1000
122 Genomes CEU and GBR populations (CEU=Utah Residents with Northern and Western
123 European Ancestry, GBR=British in England and Scotland).

124

125 While UMAP is a visualization tool with no direct interpretation on genetic distance, the
126 continuum of points connecting UMAP clusters reflects the varying degrees of estimated
127 admixture between different continental ancestries. In particular, the complex population
128 structure of Hispanics/Latinos is shown by the points spanning between the clusters of
129 European, Native American, and African ancestry. Coloring of these points based on ancestry
130 proportions affirms the relationship between the degree of admixture and their relative position
131 between reference clusters. Interestingly, African American individuals from both datasets form
132 a single continuum from the European cluster to the Yoruba (YRI) and Esan (ESN) populations
133 of Nigeria in the 1000 Genomes Project, indicative of the West African origins of most African
134 Americans. This observation is consistent with and further expands the previous finding that the
135 African tracts in the admixed 1000 Genomes populations of ACB and ASW were previously
136 found to be similar to the Nigerian YRI and ESN populations.^{2,17}

137

138 **Population differentiation and migration rate inference across the United States**

139 To better understand the relationship between genetics and geography, we investigated
140 migration rates for genetically inferred Europeans, African Americans, and Hispanic/Latinos
141 across the United States. We excluded East Asians and South Asians due to small sample size

142 and limited our analysis to the contiguous 48 states. We inferred effective migration rates with
143 the estimating effective migration surfaces (EEMS) method,²⁵ which statistically characterizes
144 genetic differentiation via resistance distance across non-homogenous landscapes. By
145 overlaying a dense regular grid of demes and measuring genetic dissimilarities between
146 neighboring demes, EEMS quantifies and visualizes areas with high relative rates of effective
147 migration (colored in blue) and areas with low relative rates of effective migration (also called
148 migration barriers and colored in dark orange).

149

150 The inferred migration rates for African Americans reveal genetic signatures of historical
151 demographic events (**Figure 2A; Figure S5**). Along the Atlantic coast from the Florida
152 Panhandle to southern Maine, we find high effective migration rates, indicating the constant
153 migration and similar effective population sizes of African Americans in these states. However,
154 we also observe a strong north-south barrier to migration starting along the Appalachian
155 Mountain Range, continuing north up the Mississippi River, and extending west across the rest
156 of the country. This migration barrier, along with the migration barrier spanning Texas and New
157 Mexico, reveals a pattern of isolation-by-distance that is consistent with the Great Migration
158 from the 1910s to the 1960s in which an estimated 6 million African Americans migrated out of
159 the South to cities across the Northeast, Midwest and West.^{7,26}

160

161 A highly complex pattern of migration exists amongst Hispanics/Latinos with varying migration
162 rates across the country, capturing regional patterns of genetic similarity. Hispanics/Latinos in
163 the southwestern states including two regions bordering Mexico--one in California and another
164 extending from New Mexico to Texas--exhibit high effective migration rates and are separated
165 by a migration barrier in Arizona (**Figure 2B; Figure S5**). These two distinct regions likely reflect
166 known differences in northward migration from east versus west Mexico.^{8,27} Along the Atlantic
167 coast from Florida to New York, effective migration has also been fluid. However, barriers to
168 migration are observed west of the Atlantic coast to the Mississippi River, likely resulting from
169 varying admixture proportions.

170

171 The patterns of migration for Europeans capture subcontinental structure. Elevated migration
172 rates are observed across most of the country, except for many states in the Midwest and along
173 the Atlantic coast. We find low effective migration rates surrounding Minnesota and North
174 Dakota, potentially due to the genetic dissimilarity of Finnish and Scandinavian ancestry
175 abundant in the region (**Figure 2C; Figure S5**).⁸ We also find reduced migration rates across

176 Ohio, West Virginia, and Virginia, suggesting the existence of genetic differentiation along the
177 Appalachian Mountains. Many of the major cities, such as Chicago, Philadelphia, and Miami,
178 are also barriers to migration, perhaps due to higher admixture proportions within cities. The
179 migration barrier encompassing metropolitan New York City may be explained in part by the
180 presence of divergent European populations, such as Ashkenazi Jews (**Figure 2C**).

181

182 **Coupling fine-scale haplotype clusters and multigenerational birth records uncovers** 183 **distinct subcontinental structure**

184 To disentangle more recent and subtle population structure, we performed identity-by-descent
185 (IBD) clustering on the Genographic cohort and annotated clusters using multigenerational self-
186 reported birth origin data. We first built an IBD network from pairwise IBD sharing among 31,783
187 unrelated individuals. In this network, vertices represent individuals and edges represent the
188 cumulative IBD (in centimorgans, cM) between pairs of individuals. We employed the Louvain
189 method, a greedy heuristic algorithm, to recursively partition vertices in the graph into clusters
190 that maximize modularity at each level of hierarchy.^{8,28} The clusters of individuals resulting from
191 each iteration can be interpreted as having greater amounts of cumulative IBD shared between
192 individuals within the cluster than with those outside of the cluster. To aid in the interpretation of
193 the clusters, we merged clusters with low genetic differentiation ($F_{ST} < 0.0001$) at the lowest
194 level of hierarchy, resulting in a final set of 25 clusters (**Table 1**). We annotated each cluster
195 based on ancestral birth origin and ethnicity data and constructed a neighbor-joining tree based
196 on the F_{ST} values (**Figure 3**). 98% of the 3,028 individuals that were not classified by our
197 Random Forest model were assigned to a haplotype cluster. No single cluster was
198 overrepresented by unclassified individuals, as unclassified individuals comprised of 8-11% of
199 each cluster.

200

201 Genetic and geographic diversity is greatest amongst Hispanic/Latino haplotype clusters. We
202 identified a total of five Hispanic-related clusters. The largest of these cluster (n=810) is strongly
203 associated with south Florida (OR = 10.4; $p = 2.5e-25$; **Figure 4, Table S4**) but is also found in
204 California, and Texas (OR ≥ 2 ; $p < 0.05$). No single ancestral birthplace characterizes this
205 cluster, as the US, Mexico, and Cuba each make up more than 10% of the birth origin labels.
206 Proportions of European ancestry tracts inferred with RFMix²⁹ are higher in this cluster (mean =
207 72.7%, sd=20.4%) than in the other Hispanic/Latino clusters (mean = 48.0% - 67.4%). Puerto
208 Ricans characterize a substantial proportion of another Hispanic/Latino cluster associated with
209 Florida (OR > 4), as well as New York City (OR > 5). Unlike the other Hispanic clusters, the

210 Puerto Rican cluster shares the same branch on the F_{ST} tree as the African American clusters,
211 likely due to high proportions of African ancestry (mean = 11.2%, sd = 9.0%) among Puerto
212 Ricans.

213
214 Three distinct clusters of Hispanics were found in the Southwest (**Figure 4**): one strongly
215 associated with New Mexico (OR > 4; $p < 0.05$), another primarily in Texas (OR > 3; $p < 0.05$),
216 and the third associated with Southern California (OR > 2; $p < 0.05$). Combined with the EEMS
217 analysis, these clusters confirm our observation of parallel migration routes from east and west
218 Mexico into Southwestern United States. While the genetic differentiation of these three clusters
219 are subtle ($F_{ST}=0.001-0.003$), ancestral birth origin patterns and local ancestry proportions for
220 these clusters reveal meaningful dissimilarities. Whereas the majority of Hispanics in New
221 Mexico report US ancestral birth origins through grandparents, the recent ancestors of
222 Hispanics in Texas are predominantly from Mexico. Nonetheless, these two clusters share
223 similar local ancestry proportions with only slight genetic dissimilarity that result in a moderate
224 decrease in migration rate (from darker blue to light blue in **Figure 2B**). The reduced migration
225 rate along the Texas-Mexico border may be caused by more recent immigrants. Unlike the
226 Hispanic clusters associated with New Mexico and Texas, the Hispanics in California cluster
227 contain greater proportions of ancestors from Central and South American (e.g., Colombia and
228 El Salvador). Proportions of Native American ancestry is also highest in this cluster (**Figure 4**).
229 Taken together, these two differences further explain the presence of the migration barrier in
230 Arizona between the Hispanics in the California and the Hispanics in New Mexico.

231
232 Historical immigration of Europeans into the US occurred in successive waves, with Northern
233 and Western Europeans making up one wave from the 1840s to 1880s and another wave
234 comprising of Southern and Eastern Europeans occurring from the 1880s to 1910s.³⁰ Consistent
235 with this immigration pattern, haplotype clusters with ancestries from Northwest and Central
236 Europe have higher proportions of US ancestral birth origins than haplotype clusters from
237 Southern and Eastern Europe, suggesting earlier immigration (**Figure 5**). The two clusters with
238 the highest proportion (>75%) of US ancestral birth origin (“Northwest Europe 1” and “Northwest
239 Europe 2”) have ~4.5% of UK ancestral origins. The Central European cluster and the Irish
240 cluster both have 66.1% and 68.5% of US ancestral origins, respectively. In contrast, the US
241 makes up only 62.2% and 34.5% of ancestral birth origin for the clusters of Southern Europeans
242 and Eastern Europeans, respectively.

243

244 Unlike the larger European clusters, the smaller European clusters reflect the structure of recent
245 immigrants and genetically isolated populations, recapitulating earlier findings.⁸ The geographic
246 distributions of these subpopulations are more concentrated, and their ancestral birth origin
247 proportions are overrepresented by specific countries and ethnicities (**Figure 6**). Specifically,
248 Finns and Scandinavians are abundant in the Upper Midwest and Washington; French
249 Canadians are found in the Northeast; Acadians are present in the Northeast and Louisiana;
250 and Italians, Greeks, Ashkenazi Jews, and Admixed Jews are mostly located in the metropolitan
251 area of New York City. Of the European clusters, median cumulative IBD sharing and cROH
252 lengths are highest amongst Ashkenazi Jews (31.8cM and 11.3 Mb, respectively; **Table 1**). The
253 two Jewish-related clusters were identified using self-reported ancestral ethnicity data rather
254 than birth origin data, since Jewish ancestry is not specific to any single location. Jewish
255 ancestry, particularly Ashkenazi Jewish ancestry, was more consistently reported on both sides
256 of the family in the larger Jewish cluster (“Ashkenazi Jewish”), suggesting that individuals are
257 more admixed in the smaller cluster (“Admixed Jewish”).

258
259 We inferred two haplotype clusters of African Americans separated along a north-south cline,
260 recapitulating the EEMS migration barrier inference. One cluster is primarily distributed amongst
261 the northern and western states (“African Americans North”) while the other is distributed
262 amongst the states southeast of the Appalachian Mountains (“African Americans South”) (**Figure S7**).
263 The proportion of US birth origin is higher in the northern cluster than the southern
264 cluster, further evidence of isolation by distance amongst African Americans in the north.⁷ These
265 two clusters share similar cROH lengths but differ in admixture proportions and median IBD
266 sharing, pointing to a cluster with consistent African American ancestors and a cluster with more
267 admixed ancestors. Median IBD sharing is higher amongst African Americans in the south
268 (median IBD = 19.6 cM, median cROH = 3.3 Mb) than in the north (median = 15.9 cM; **Table 1**)
269 while the average proportion of African ancestry is higher in the northern cluster than the
270 southern cluster.

271
272 Four of the clusters reflect recent immigrants from Asia (**Figure S8**), which grew rapidly in the
273 mid-20th Century after the elimination of national origin quotas.³¹ The recency of immigration
274 among these clusters is indicated by the less than 30% of ancestral birth origins coming from
275 the US. Geographically, individuals in these clusters primarily reside in major cities. East Asians
276 predominantly inhabit the metropolitan areas of the West and Northeast (OR > 2), Southeast
277 Asians are enriched in the West (OR > 2.5), and South Asians are strongly associated with the

278 Northeast (OR > 2.5). Despite its small size, the cluster of Middle East individuals reflects many
279 of the known demographic patterns of Arab Americans, as individuals in this cluster are
280 primarily of Lebanese origin and are distributed in the Northeast as well as metropolitan Detroit.
281 cROH lengths are particularly long for South Asians (median cROH = 10.3 cM), Southeast
282 Asians (median cROH = 7.8 cM), and Middle Easterners (median cROH = 8.2 cM), potentially
283 reflecting inbreeding patterns found in their ancestral regions.³²

284

285 **Discussion**

286

287 As the US population is becoming increasingly diverse, genomic studies are simultaneously
288 growing in scale and relevance; to increase scientific and ethical parity, these studies must
289 move beyond the current practice of evaluating genetically homogenous groups in isolation.¹⁵
290 Here, we provide an integrative framework for analyzing population structure in ancestrally
291 heterogeneous individuals. Our comprehensive approach has allowed us to capture spatial
292 patterns of gene flow within and between subpopulations that are difficult to infer from a single
293 method alone. For example, EEMS is limited in identifying unique subpopulations, while
294 haplotype clustering cannot assign partial membership for admixed individuals to multiple
295 clusters. An integrative approach can thus enable greater insights into populations with complex
296 histories.

297

298 Consistent with prior studies,^{4,9} the recent demographic history of Hispanic/Latino populations is
299 complex. Large variations in admixture proportions within and between subpopulations are
300 reflected by US Census data and can likely be explained by numerous inferred migration
301 barriers. For example, regional differences in the Southwest are highlighted by an inferred
302 migration barrier in Arizona and distinct haplotype clusters surrounding this region. These
303 differences are likely due to higher proportions of Native American ancestry as well as more
304 Central and South American origins in the California Hispanics/Latinos compared to other
305 southwestern Hispanic/Latinos. Interestingly, although the New Mexican cluster is distinct from
306 the Texan cluster, high levels of gene flow are inferred from southern New Mexico to central
307 Texas, suggesting that certain individuals in these two clusters are genetically similar and may
308 share an ancestral origin (i.e. Mexico). In contrast, those in northern New Mexico are more
309 genetically differentiated, as indicated by a migration barrier, but share the same cluster; these
310 are likely *Nuevomexicanos*, descendants of Spanish colonial settlers.

311

312 The fine-scale population structure of African Americans also reflects known historical events
313 following the transatlantic slave trade, during which millions of West Africans were forcibly
314 moved to the Americas. Subsequently, the movement of African Americans during the Great
315 Migration has been shown to correlate with current patterns of relatedness across US census
316 regions.⁷ Our results show barriers to migration and gene flow at fine-scale, particularly along
317 the Appalachian Mountains. A north-south migration barrier is also present west of the
318 Mississippi River, and is further supported by the north-south locations of two African American
319 clusters that emphasize this divide. The southern African American cluster contains more recent
320 ancestors outside the US, particularly of Caribbean origin, than the northern African American
321 cluster. These genetic signatures illustrate the impact of recent migration patterns on modern
322 population structure.

323
324 Our ability to identify population structure for certain ancestries is subject to participation among
325 individuals from those groups. In particular, individuals with Asian ancestries account for over
326 5% of US population, but they are underrepresented in US population genetics studies,
327 hindering the investigation of their ancestry in prior studies.⁸ Our analyses of East Asian,
328 Southeast Asian, South Asian, and Middle Eastern populations therefore provide initial insights
329 into their genetic structure. The ancestral origins and geographic distributions of these clusters
330 are consistent with US Census reports. Since these populations descend from more recent
331 immigrants, the observed patterns of homozygosity within several of these clusters likely reflect
332 consanguinity patterns in some of their ancestral regions. Specifically, the long cROH in South
333 Asians may reflect endogamy for example related to the caste system in India, while similar
334 patterns among the Middle Eastern and Southeast Asian clusters may be capturing
335 consanguineous marriage practices in those regions.^{33–35} Given the small size of these clusters,
336 however, further studies of more individuals are needed.

337
338 Population history in the US is best characterized among the most populous European descent
339 individuals. Genetic diversity tends to be highest in more densely populated regions, likely due
340 to the presence of multiple subpopulations in the same place. Many of the European
341 subpopulations we identified are similar to those previously found—e.g., French Canadians,
342 Acadians, Scandinavians, and Jews (Supplementary Discussion).⁸ The geographic distribution
343 of these subpopulations, particularly those that are more genetically diverged, overlap in the
344 metropolitan areas in the Northeast, Midwest, and California.

345

346 The precision of population labels assigned to clusters of individuals is a function of
347 demographic complexity and sample size. For example, Finnish ancestry is clearly European
348 but genetically distinct from several other European populations due to historical bottlenecks,
349 making this ancestry cluster relatively easily separable. By contrast, most Americans of
350 European descent have heterogeneous ancestors from several northwestern European
351 countries who have admixed over time. Additionally, while we identify and describe some
352 substantial structure among Hispanic/Latino populations, considerably more is likely to exist and
353 remains to be learned from larger and more diverse future studies. Similarly, sub-regional
354 resolution into the ancestry of recent Asian immigrants to the US has been relatively limited in
355 population genetics studies, and the structure of this immigration will be learned from larger
356 future studies. Additional considerations relating to population label precision are the accuracy
357 of self-reported birth records and variable granularity of geopolitical boundaries.

358
359 The emergence of biobank-scale genomic data is enabling more complete pedigrees,³⁶ greater
360 discoveries of fine-scale population structure, and more precise insights into health-related
361 associations. An estimated 26 million people have taken a direct-to-consumer ancestry test,³⁷
362 indicating widespread interest in ancestry and heritable factors. As participation in genetic
363 studies increase, especially in the US with the All of Us Research Program, so does the need
364 for inferring more granular demographic histories in study cohorts. Understanding such structure
365 is important to account for stratification, prevent the overgeneralization of results, and avoid
366 exacerbating existing biases.^{14,15} This study demonstrates the potential of coupling genetic data
367 with geographic and birth origin data to reconstruct such demographic histories, particularly in a
368 large and heterogeneous population.

369 **Materials and Methods**

370

371 **Human Subjects**

372 The Genographic Project and Geno 2.0 Project received full approval from the Social and
373 Behavioral Sciences Institutional Review Board (IRB) at the University of Pennsylvania Office of
374 Regulatory Affairs on April 12, 2005. The IRB operates in compliance with applicable laws,
375 regulations, and ethical standards necessary for research involving human participants. All data
376 in this study came from participants that consented to have their results be used in scientific
377 research. All data was deidentified.

378

379 Participants provided genotype data, geographic location (postal code), ancestral birth origin,
380 and self-declared ethnicity. We limited our study to those individuals who provided valid
381 geographic location. Ancestral birth origin and self-declared ethnicity data were collected up to
382 the grandparents of the participants with ~60% of individuals provided complete pedigrees.

383

384 **Genotyping and Quality Control**

385 Participants of the Genographic project were sequenced with the GenoChip array,²⁰ an Illumina
386 iSelect HD custom genotyping bead array with approximately 150,000 Ancestry Informative
387 Markers. Quality control and phasing of data is described in Supplemental Materials and
388 Methods. After QC, 32,589 individuals and 108,003 sites remained.

389

390 **Principal Component Analysis**

391 We performed principal component analysis on the quality-controlled samples using FlashPCA
392 version 2.0.²² We included the genotypes of all 2,504 individuals from the 1000 Genomes
393 Project as reference samples. We computed PCs across 108,003 shared sites for 1000
394 Genome Project individuals and then projected the Genographic individuals on the same
395 principal component space.

396

397 **Continental Ancestry Assignment**

398 We assigned continental ancestry to each Genographic sample by using a random forest
399 classifier. Using the PCs and known super population assignment (AFR=African,
400 EUR=European, EAS=East Asian, AMR=American, and SAS=South Asian) from the 1000
401 Genome Project samples as training data, we applied the classifier to assign ancestry to each

402 Genographic sample at 90% probability. We considered unassigned ancestries as “other”
403 (OTH).

404

405 **Genetic Ancestry Proportion Estimation**

406 We estimated admixture proportions using ADMIXTURE by first analyzing the 1000 Genomes
407 Project in unsupervised mode to learn allele frequencies.²¹ Then, we projected the learned allele
408 frequencies onto the Genographic samples to obtain the admixture proportions. We ran
409 ADMIXTURE with $k=2-9$ and chose $k = 5$ as the most stable representation.

410

411 **UMAP**

412 We applied the Uniform Manifold Approximation and Projection (UMAP) method to visualize
413 subcontinental structure.^{23,24} We first combined the PCs of the Genographic samples and the
414 1000 Genome Project samples into one dataset. We then applied UMAP on the first 20 PCs
415 from the joint dataset to produce a two-dimensional plot. We tested various parameter choices
416 for UMAP and found that the default nearest neighbor value of 15 and the minimum distance
417 values of 0.5 delivered the clearest result. Coloring of UMAP plots are described in the
418 Supplemental Materials and Methods.

419

420 **Estimating Effective Migration Surfaces**

421 We estimated migration and diversity relative to geographic distance using the estimating
422 effective migration surfaces (EEMS) method for Genographic individuals that were classified
423 under African, European, and Native American ancestries.²⁵ We excluded East Asian and South
424 Asian ancestries due to low sample size and density. We used unrelated individuals with
425 available postal code data. We first computed pairwise genetic dissimilarities with the EEMS
426 *bed2diffs* tool and then ran EEMS with *runeems_snps*, setting the number of demes to 500. Per
427 the recommendation in the manual, we adjusted the variance for all proposed distributions of
428 diversity, migration, and degree-of-freedom parameters such that all were accepted 10%-40%
429 of the time. We increased the number of Markov chain Monte Carlo (MCMC) iterations until it
430 converged.

431

432 **Haplotype Calling and Network Construction**

433 We used IBDSeq version r1206 to generate shared identity-by-descent (IBD) segments from
434 genotype data for all unrelated individuals.³⁸ Unlike other IBD detection algorithms, IBDseq does
435 not rely on phased genotype data and is less susceptible to switch errors in phasing that can

436 cause erroneous haplotype breaks. We filter for IBD segments greater than 3cM. We removed
437 segments that overlapped with long chromosomal regions (1 Mb) that had no SNPs across all
438 unrelated individuals. These sites can result in false positives IBD sharing and likely correspond
439 to centromeres and telomeres. We calculate the cumulative IBD sharing between individuals by
440 summing the length of all shared IBD segments. We then constructed a haplotype network of
441 unrelated individuals by defining vertices an individuals and edge weights between vertices as
442 the cumulative IBD sharing between individuals. We filtered to keep edges with cumulative IBD
443 sharing is ≥ 12 cM and ≤ 72 cM, as previously described.⁸

444

445 **Detection of IBD Clusters**

446 To identify clusters of related individuals in the haplotype network, we used the Louvain Method
447 implemented in the igraph package for R. The Louvain Method is a greedy iterative algorithm
448 that assigns vertices of a graph into clusters to optimize modularity (a measure of the density of
449 edges within a community to edges between communities). The Louvain Method begins by first
450 assigning each node as its own community and then adds node i to a neighbor community j . It
451 then calculates the change in modularity and places i in the community with that maximizes
452 modularity. The algorithm repeats this continuously and terminates when no vertices can be
453 reassigned.

454

455 We partitioned the haplotype network into clusters by recursively applying the Louvain Method
456 within subcommunities. At the highest level, we take the full, unpartitioned haplotype graph and
457 identify a set of subcommunities. We isolate the vertices within each subcommunity, keeping
458 only the edges between those vertices to create separate new networks. We then apply the
459 Louvain Method to the new subgraphs. We repeat this process up to four levels. We combined
460 subcommunities with low genetic divergence based on F_{ST} values of < 0.0001 .

461

462 **Annotation of IBD Clusters**

463 We used a combination of ancestral birth origins and self-reported ethnicities to discern
464 demographic characteristics of each cluster. For each cluster, we quantified the proportion of
465 each birth origin (i.e. country of origin) amongst all four grandparents, treating each
466 grandparent's origin equality. We use these proportions to inform population labels. Clusters in
467 which a single non-US birth origin was in high proportions was labeled with that country. In
468 cases where multiple non-US birth locations exists in approximately equally high proportions,
469 we assigned a label representing the broader region (e.g. Eastern Europeans for Poland,

470 Lithuania, Ukraine, and Slovakia; East Asia for Japan, China). For certain clusters, annotations
471 could not be easily discerned by birth origin data. In these cases, we relied on self-reported
472 ethnicities to label the clusters as these populations were found to be less associated with a
473 non-US country (e.g. Ashkenazi Jews) or the population has resided in the US for generations
474 (African Americans, Acadians).

475

476 **Mapping IBD Clusters**

477 We mapped individuals using their present-day geographic location. We aggregated individuals
478 from the same county using the postal code to county FIPS code mapping provided by the US
479 Census. Longitude and latitude points of each county was found using the same data from the
480 US Census. We identified enriched counties for each cluster by performing a Fisher's exact test
481 on each county that had ≥ 30 individuals to obtain an odds ratio and significance value. We
482 mapped only counties with statistically significant ($p < 0.05$) enrichment and an odds ratio (OR) of
483 greater than 1. The size of the circles is scaled to the number of individuals in each location.

484

485 **Runs of Homozygosity**

486 We used PLINK v1.90b3.39 to infer runs of homozygosity with a window of 25 SNPs.³⁹ We
487 calculated the cumulative runs of homozygosity (cROH) size by summing the lengths of
488 homozygous segments.

489

490 **Local Ancestry Inference**

491 We inferred local ancestry with RFMix v1.5.4 for Genographic samples in clusters that were
492 annotated as Hispanics/Latinos and African Americans.²⁹ We used samples of African (LWK,
493 MSL, GWD, YRI, ESN, ACB, and ASW; $N = 661$), European (CEU, GBR, FIN, IBS, and TSI; N
494 $= 503$), and Native American (MXL, PUR, CLM, and PEL; $N = 347$) ancestry from the 1000
495 Genomes Project as the reference population. We ran RFMix with the default minimum window
496 size (0.2 cM) and a node size of 5 with the flags: -w 0.2, -n 5. Global ancestry proportions were
497 derived by quantifying the proportions of total local ancestry tracts for each ancestry.

498

499 **Genetic Divergence**

500 We computed weighted Weir-Cockerham F_{ST} estimates for each pair of haplotype clusters using
501 PLINK v1.90b3.39.³⁹ Using the distance matrix of F_{ST} values between clusters, we constructed
502 an unrooted phylogenetic tree using the neighbor joining method implemented in *scikit-bio*.⁴⁰ We
503 visualized the tree using Interactive Tree Of Life.⁴¹

504 **Data and Code Availability**

505 Genotype data and associated metadata are available to researchers through an application
506 process and data usage agreement. We encourage qualified researchers to email the
507 Genographic team at National Geographic Society (genographic@ngs.org) for information on
508 and access to the Genographic database.

509

510 Custom scripts generated to analyze the data in this paper are available through GitHub
511 (https://github.com/chengdai/genographic_ancestry).

512

513 **Acknowledgement**

514 We thank the National Geographic Genographic Project participants who consented to research
515 participations for making this study possible. We also thank Gregory Vilshansky for helping
516 organize and manage the data for the Genographic Project.

517

518 This work was supported by funding from the National Institutes of Health (K99MH117229 to
519 A.R.M.). C.L.D., M.M., R.T., and C.R. would also like to thank all the members of the MIT
520 Senseable City Lab Consortium for supporting this research. M.G.V. acknowledges support
521 from the National Geographic Society.

522

523 **Author Contributions**

524 C.L.D. and A.R.M. designed the study, performed research, and wrote the manuscript. R.S.W.
525 founded and formerly directed the Genographic Project. M.G.V. coordinated and supervised the
526 Genographic Project. M.M.V., C.H.Y., and R.T. contributed to the data aggregation and data
527 analysis. A.R.M., C.R. and M.J.D. supervised research. All authors reviewed the manuscript.

528

529 **Conflicts of Interest**

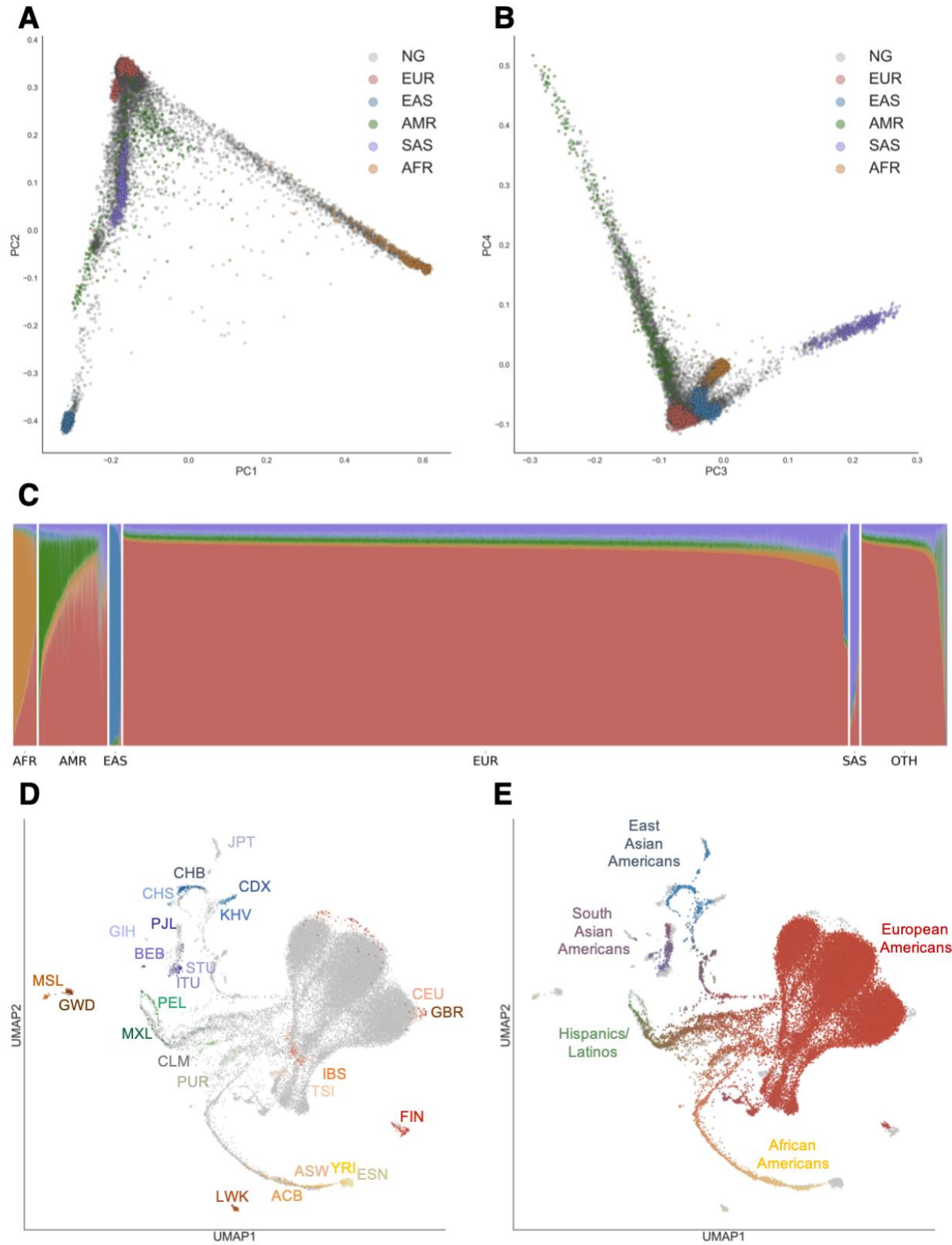
530 M.G.V. is the Senior Program Officer for the National Geographic Society and lead scientist for
531 the Genographic Project. R.S.W. was the former Director of the Genographic Project and is a
532 cofounder for Insitome. M.J.D. is a member of the Scientific Advisory Board at Ancestry.com
533 LLC.

534 **References**

- 535 1. The 1000 Genomes Project Consortium. A global reference for human genetic variation.
536 *Nature* **526**, 68–74 (2015).
- 537 2. Tishkoff, S. A. *et al.* The Genetic Structure and History of Africans and African Americans.
538 *Science* **324**, 1035–1044 (2009).
- 539 3. Bryc, K. *et al.* Genome-wide patterns of population structure and admixture in West Africans
540 and African Americans. *Proc. Natl. Acad. Sci.* **107**, 786–791 (2010).
- 541 4. Bryc, K. *et al.* Genome-wide patterns of population structure and admixture among
542 Hispanic/Latino populations. *Proc. Natl. Acad. Sci.* **107**, 8954–8961 (2010).
- 543 5. Reich, D. *et al.* Reconstructing Native American population history. *Nature* **488**, 370–374
544 (2012).
- 545 6. Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).
- 546 7. Baharian, S. *et al.* The Great Migration and African-American Genomic Diversity. *PLOS*
547 *Genet.* **12**, e1006059 (2016).
- 548 8. Han, E. *et al.* Clustering of 770,000 genomes reveals post-colonial population structure of
549 North America. *Nat. Commun.* **8**, 14238 (2017).
- 550 9. Bryc, K., Durand, E. Y., Macpherson, J. M., Reich, D. & Mountain, J. L. The Genetic
551 Ancestry of African Americans, Latinos, and European Americans across the United States.
552 *Am. J. Hum. Genet.* **96**, 37–53 (2015).
- 553 10. Wang, S. R. *et al.* Simulation of Finnish Population History, Guided by Empirical Genetic
554 Data, to Assess Power of Rare-Variant Tests in Finland. *Am. J. Hum. Genet.* **94**, 710–720
555 (2014).
- 556 11. Bray, S. M. *et al.* Signatures of founder effects, admixture, and selection in the Ashkenazi
557 Jewish population. *Proc. Natl. Acad. Sci.* **107**, 16222–16227 (2010).
- 558 12. Mooney, J. A. *et al.* Understanding the Hidden Complexity of Latin American Population
559 Isolates. *Am. J. Hum. Genet.* **103**, 707–726 (2018).
- 560 13. Belbin, G. M. *et al.* Genetic identification of a common collagen disease in Puerto Ricans via
561 identity-by-descent mapping in a health system. *eLife* **6**, e25060 (2017).
- 562 14. Popejoy, A. B. & Fullerton, S. M. Genomics is failing on diversity. *Nat. News* **538**, 161
563 (2016).
- 564 15. Martin, A. R. *et al.* Current clinical use of polygenic scores will risk exacerbating health
565 disparities. *bioRxiv* 441261 (2019). doi:10.1101/441261
- 566 16. Caswell-Jin, J. L. *et al.* Racial/ethnic differences in multiple-gene sequencing results for
567 hereditary cancer risk. *Genet. Med. Off. J. Am. Coll. Med. Genet.* **20**, 234–239 (2018).

- 568 17. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction across
569 Diverse Populations. *Am. J. Hum. Genet.* **100**, 635–649 (2017).
- 570 18. Manrai, A. K. *et al.* Genetic Misdiagnoses and the Potential for Health Disparities. *N. Engl.*
571 *J. Med.* **375**, 655–665 (2016).
- 572 19. Morales, J. *et al.* A standardized framework for representation of ancestry data in genomics
573 studies, with application to the NHGRI-EBI GWAS Catalog. *Genome Biol.* **19**, (2018).
- 574 20. Elhaik, E. *et al.* The GenoChip: A New Tool for Genetic Anthropology. *Genome Biol. Evol.* **5**,
575 1021–1031 (2013).
- 576 21. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in
577 unrelated individuals. *Genome Res.* **19**, 1655–1664 (2009).
- 578 22. Abraham, G., Qiu, Y. & Inouye, M. FlashPCA2: principal component analysis of Biobank-
579 scale genotype datasets. *Bioinformatics* **33**, 2776–2778 (2017).
- 580 23. McInnes, L. & Healy, J. UMAP: Uniform Manifold Approximation and Projection for
581 Dimension Reduction. *ArXiv180203426 Cs Stat* (2018).
- 582 24. Diaz-Papkovich, A., Anderson-Trocme, L. & Gravel, S. Revealing multi-scale population
583 structure in large cohorts. (2018). doi:10.1101/423632
- 584 25. Petkova, D., Novembre, J. & Stephens, M. Visualizing spatial population structure with
585 estimated effective migration surfaces. *Nat. Genet.* **48**, 94–100 (2016).
- 586 26. US Census Bureau. The Great Migration, 1910 to 1970. *U.S. Census* Available at:
587 <https://www.census.gov/dataviz/visualizations/020/>. (Accessed: 21st February 2019)
- 588 27. Massey, D. S., Rugh, J. S. & Pren, K. A. The Geography of Undocumented Mexican
589 Migration. *Mex. Stud. Mex.* **26**, 129–152 (2010).
- 590 28. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities
591 in large networks. *J. Stat. Mech. Theory Exp.* **2008**, P10008 (2008).
- 592 29. Maples, B. K., Gravel, S., Kenny, E. E. & Bustamante, C. D. RFMix: A Discriminative
593 Modeling Approach for Rapid and Robust Local-Ancestry Inference. *Am. J. Hum. Genet.* **93**,
594 278–288 (2013).
- 595 30. Passel, J. S. & Fix, M. U.S. Immigration in a Global Context: Past, Present, and Future.
596 *Indiana J. Glob. Leg. Stud.* **2**, 5–19 (1994).
- 597 31. Grieco, E. M., Trevelyan, E., Larsen, L., Acosta, Y. D. & Gambino, C. The Size, Place of
598 Birth, and Geographic Distribution of the Foreign-Born Population in the United States: 1960
599 to 2010. *Popul. Div. Work. Pap. No 96 US Census Bur.* 38
- 600 32. Pemberton, T. J. *et al.* Genomic Patterns of Homozygosity in Worldwide Human
601 Populations. *Am. J. Hum. Genet.* **91**, 275–292 (2012).

- 602 33. Moorjani, P. *et al.* Genetic Evidence for Recent Population Mixture in India. *Am. J. Hum.*
603 *Genet.* **93**, 422–438 (2013).
- 604 34. Tadmouri, G. O. *et al.* Consanguinity and reproductive health among Arabs. *Reprod. Health*
605 **6**, 17 (2009).
- 606 35. Hussain, R. & Bittles, A. H. Assessment of association between consanguinity and fertility in
607 Asian populations. *J. Health Popul. Nutr.* **22**, 1–12 (2004).
- 608 36. Erlich, Y., Shor, T., Pe'er, I. & Carmi, S. Identity inference of genomic data using long-range
609 familial searches. *Science* **362**, 690–694 (2018).
- 610 37. Regalado, A. More than 26 million people have taken an at-home ancestry test. *MIT*
611 *Technology Review* Available at: [https://www.technologyreview.com/s/612880/more-than-](https://www.technologyreview.com/s/612880/more-than-26-million-people-have-taken-an-at-home-ancestry-test/)
612 [26-million-people-have-taken-an-at-home-ancestry-test/](https://www.technologyreview.com/s/612880/more-than-26-million-people-have-taken-an-at-home-ancestry-test/). (Accessed: 21st February 2019)
- 613 38. Detecting identity by descent and estimating genotype error rates in sequence data. -
614 PubMed - NCBI. Available at: <https://www.ncbi.nlm.nih.gov/pubmed/24207118>. (Accessed:
615 21st February 2019)
- 616 39. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based
617 linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 618 40. Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing
619 phylogenetic trees. *Mol. Biol. Evol.* **4**, 406–425 (1987).
- 620 41. Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and
621 annotation of phylogenetic and other trees. *Nucleic Acids Res.* **44**, W242-245 (2016).
622



623

624

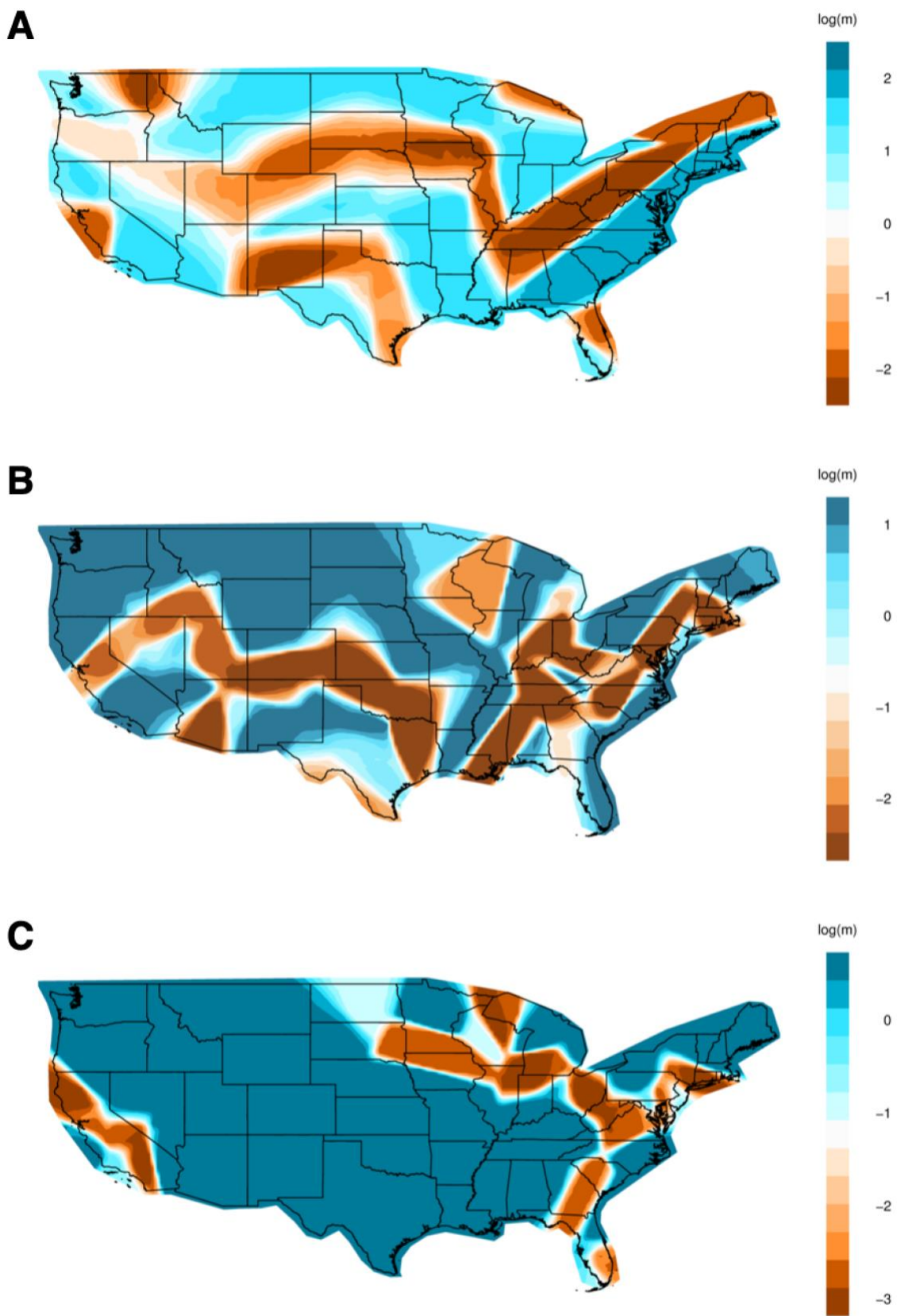
625 **Figure 1. Genetic Diversity of the US Population**

626 (A) Principal Components Analysis of individuals in the United States and in the 1000 Genome
627 Project. Each individual is represented by a single dot. Individuals in this study are colored in
628 grey while 1000 Genome Population individuals are colored by super population (EUR =
629 European, AFR = African, AMR = Admixed American, EAS = East Asian, SAS = South Asian).
630 Principal components (PC) 1 and PC2 are shown.

631 (B) Similar to (A), with PC3 and PC4 shown.

632 (C) ADMIXTURE analysis at K=5 of individuals in this study. Each individual was assigned a
633 continent-level ancestry label using a Random Forest model trained on the super population
634 labels and the first 10 PCs of the 1000 Genome Project dataset. OTH = individuals who did not
635 meet the 90% confidence threshold for classification.

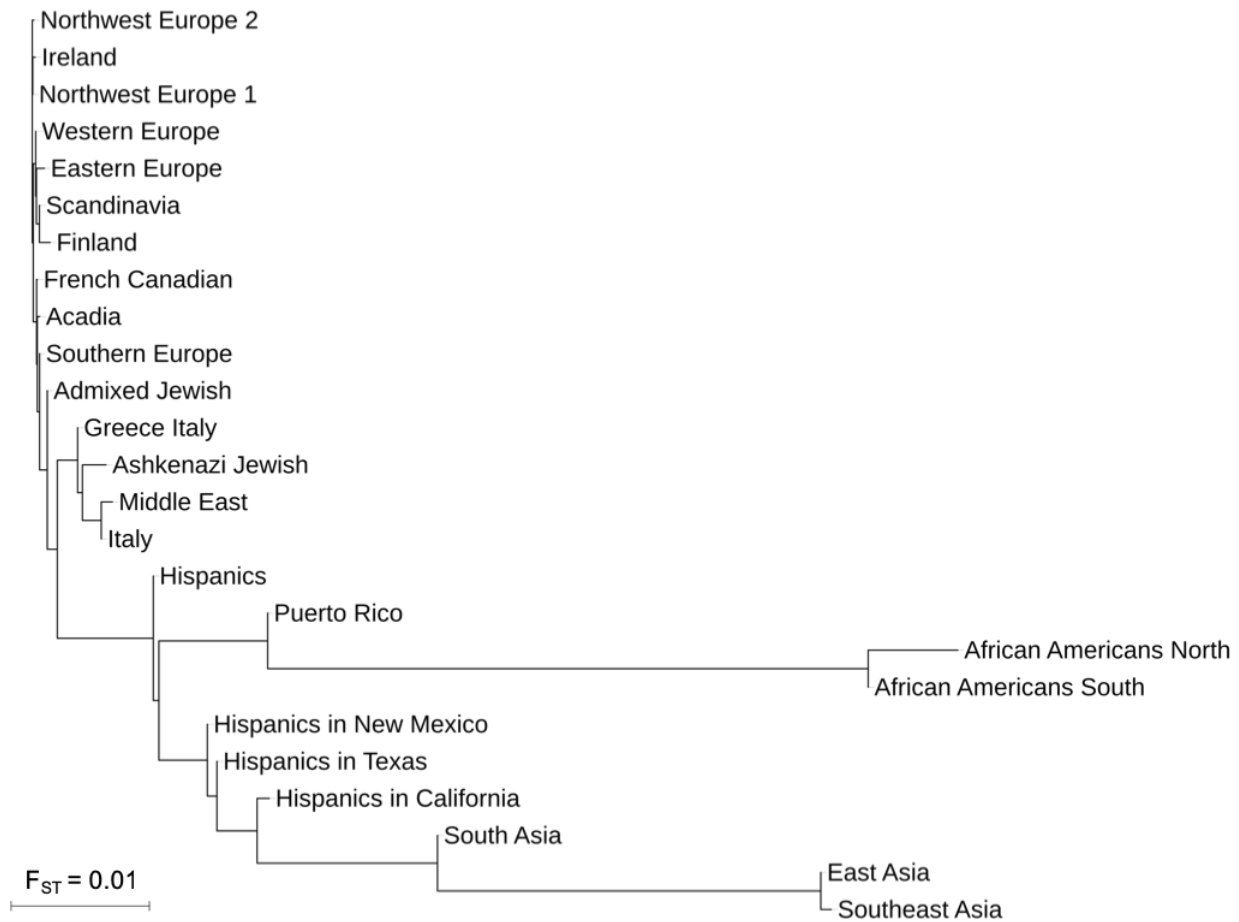
636 (D) UMAP projection of the first 20 PCs. Each dot represents one individual. In (D), individuals
637 in the 1000 Genomes Project are colored by population, while Genographic Project individuals
638 from this study are in grey. In (E), 1000 Genome Project individuals are colored in grey while
639 Genographic Project individuals are colored based on their admixture proportions from
640 ADMIXTURE. The color for each dot was calculated as a linear combination of each individual's
641 admixture proportion and the RGB values for the colors assigned to each continental ancestry
642 (EUR = red, AFR = yellow, NAT or Native American = green, EAS = blue, SAS = purple).
643 Distances in UMAP do not directly correspond to genetic distance. See Materials and Methods
644 for specific population labels.



646 **Figure 2. Migration Rates of African Americans, Hispanics/Latinos, and Europeans within**
647 **the United States.**

648 (A) - (C) Migration rates inferred with EEMS for African Americans (A), Hispanics/Latinos (B),
649 and Europeans (C). Colors and values correspond to inferred rates, m , relative to the overall
650 migration rate across the country. Shades of blue indicate logarithmically higher migration (i.e.
651 $\log(m) = 1$ represents effective migration that is ten-fold faster than the average) while shades
652 of orange indicate migration barriers.

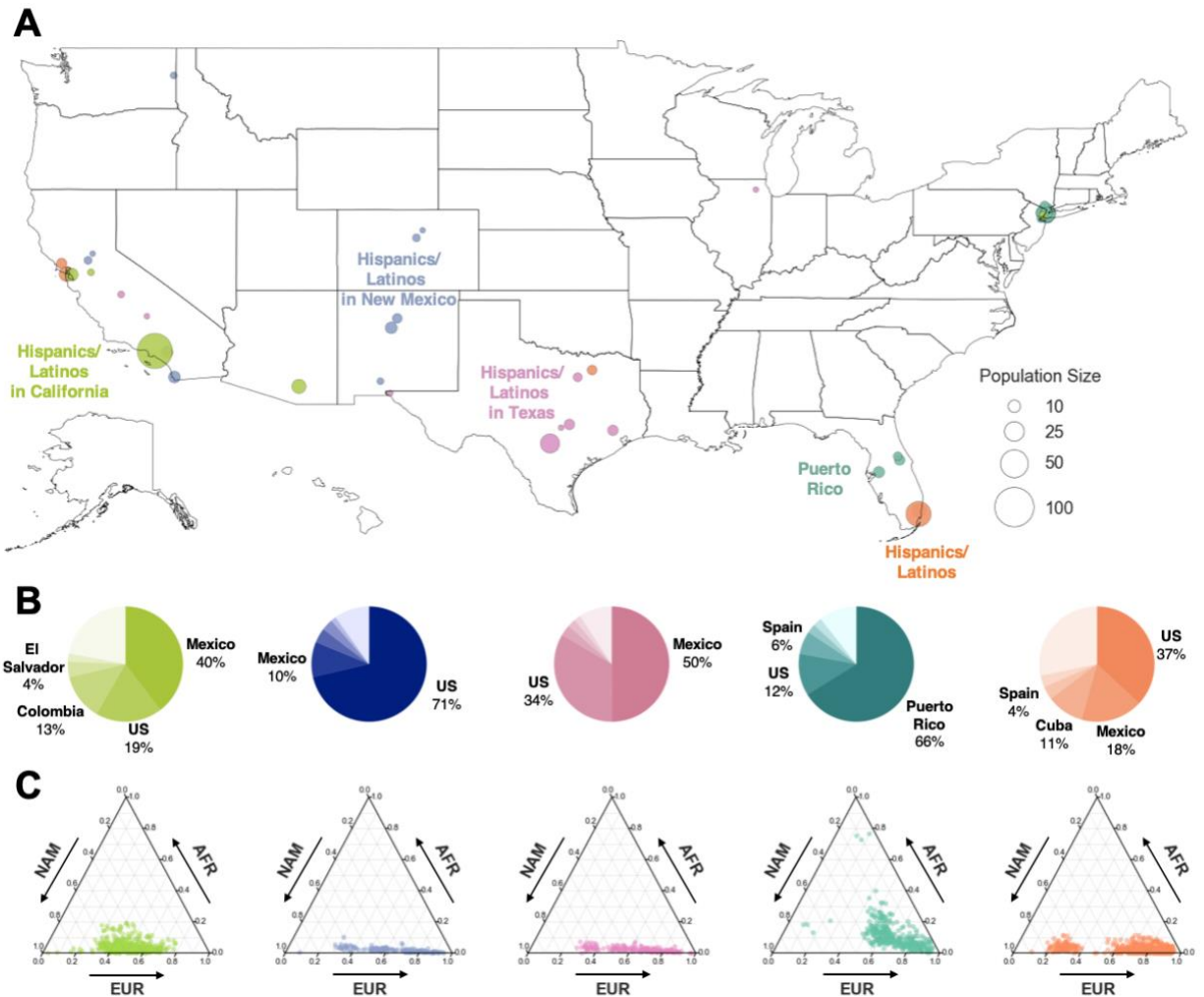
653



654

655 **Figure 3. Genetic differentiation of haplotype clusters**

656 Unrooted phylogenetic tree of haplotype clusters was constructed using the neighbor joining
657 method with F_{ST} as genetic distance. Negative branch lengths were converted to zero.



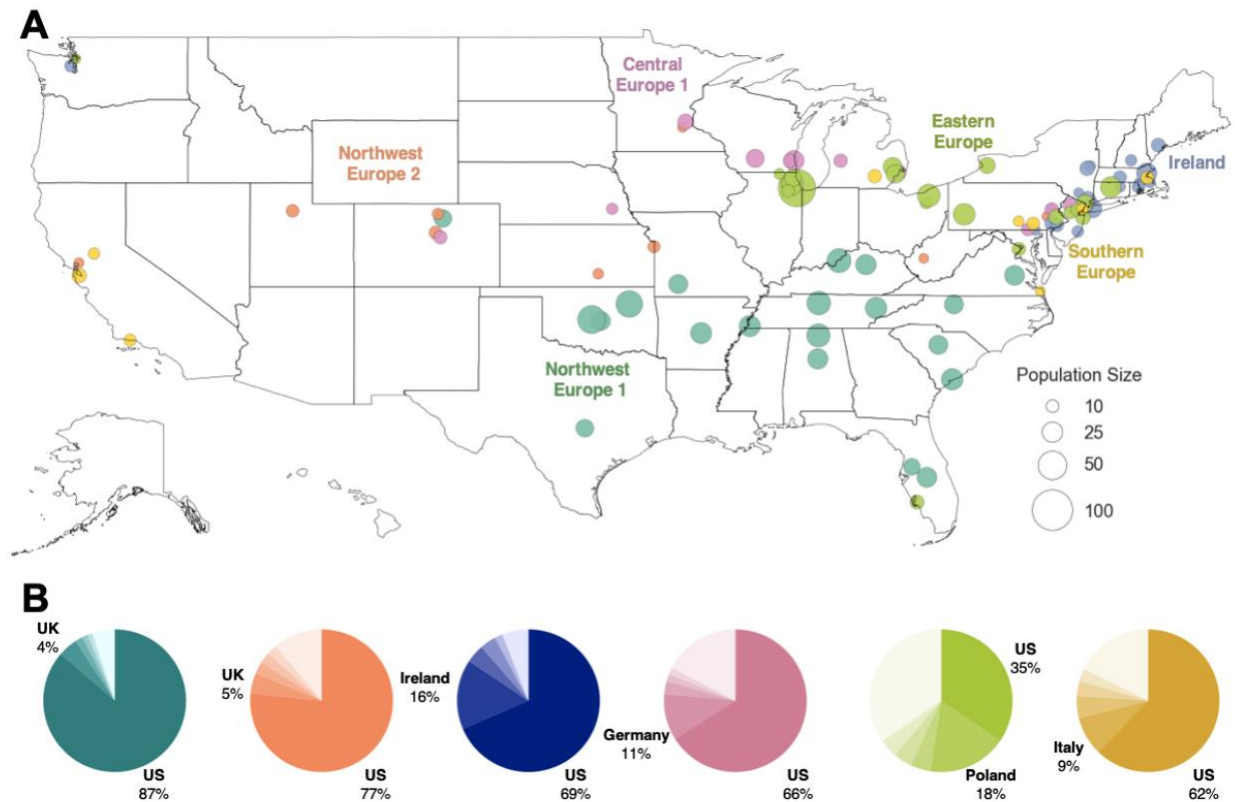
658

659 **Figure 4. Distribution of Hispanic/Latino Haplotype Clusters**

660 (A) Map of counties in which Hispanic/Latino haplotype clusters are enriched. Each dot
 661 corresponds to a county, and the size of the dot signifies the number of samples of the
 662 particular cluster in that county. Only the Hispanic/Latino cluster with the highest odds ratio is
 663 shown for each county, and only the top ten locations with the highest odds ratios are shown for
 664 each cluster. Maps showing the full distribution for each haplotype cluster can be found in the
 665 supplement (**Figure S8**).

666 (B) Ancestral birth origin proportions of each cluster for individuals with complete pedigree
 667 annotations, up to grandparent level. Proportions were calculated from aggregating the birth
 668 locations of all grandparents corresponding to members of each haplotype cluster. For each
 669 chart, only the top five birth origins are shown as individual slices; the remaining birth origins are
 670 aggregated into one slice (lightest color).

671 (C) Ternary plots of ancestry proportions based on local ancestry inference for each haplotype
672 cluster. Each dot represents one individual.

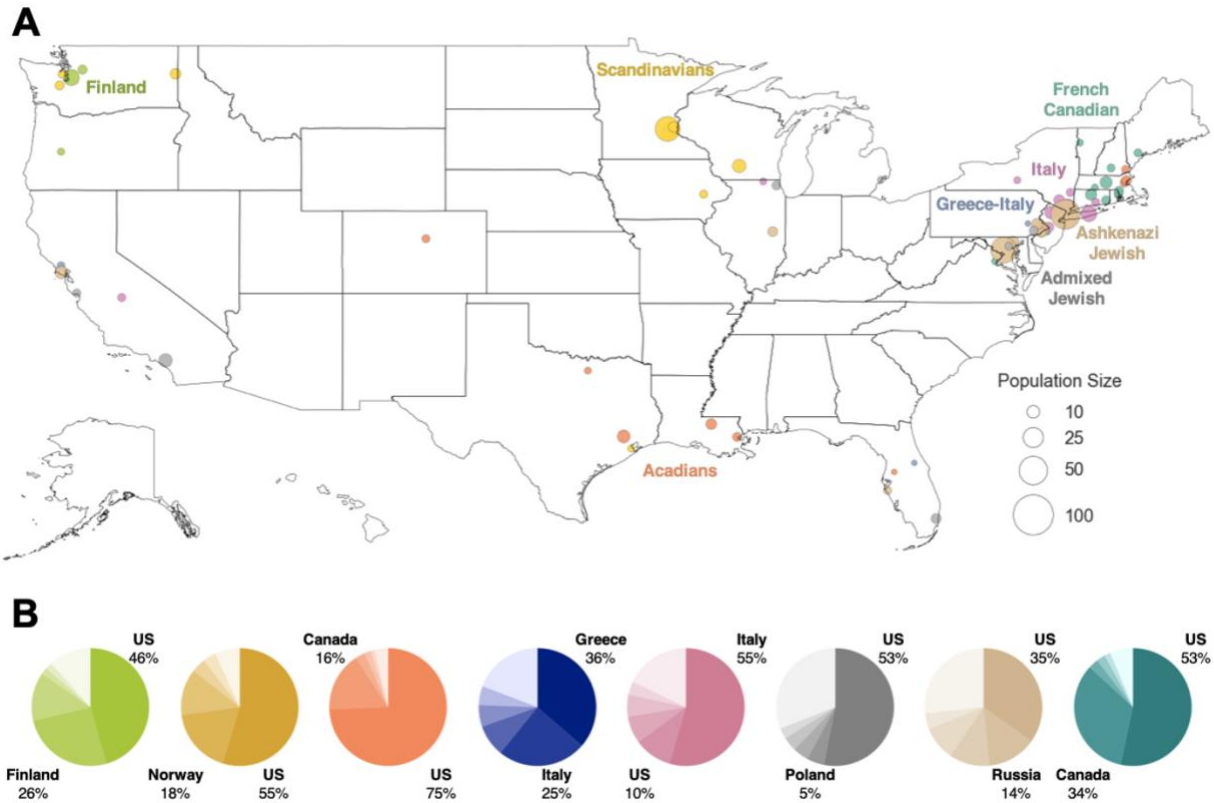


673

674 **Figure 5. Distribution of European American Haplotype Clusters**

675 (A) Geographic distributions of haplotype clusters corresponding to regional European
676 ancestries. Each county containing present-day individuals is represented by a dot. The top 20
677 locations with the highest odds ratio are shown for each cluster. Maps showing the full
678 distribution for each cluster can be found in the supplement (**Figure S8**).

679 (B) Ancestral birth origin proportions for each cluster in (A). Only individuals with complete
680 pedigree annotations, up to grandparent level, are included. For each chart, only the top five
681 birth origins are visualized as individual slices; the remaining birth origins are aggregated into
682 one slice (lightest color).



683

684

685 **Figure 6. Distribution of European American Haplotype Clusters**

686 (A) Present-day location of individuals in clusters of more genetically isolated European
687 populations, similar to Figure 5A. For clarity, the top ten locations with the highest odds ratio are
688 shown for each cluster.

689 (B) Ancestral birth origin proportions for each cluster in (A). Only individuals with complete
690 pedigree annotations, up to grandparent level, are shown. For each chart, only the top five birth
691 origins are shown as individual slices; the remaining birth origins are aggregated into one slice
692 (lightest color).

693

Cluster	Samples	Median Cumulative ROH	Median Cumulative IBD
Northwest Europe 1	11,725	2.88	15.23
Northwest Europe 2	1,571	2.80	15.15
Ireland	2,137	2.85	15.42
Central Europe	3,116	2.83	15.06
Eastern Europe	2,471	3.16	15.37
Southern Europe	1,626	2.73	14.98
Italy	697	6.91	14.64
Greece-Italy	238	7.28	15.02
Scandinavia	717	3.02	15.54
Finland	314	3.67	17.50
Acadia	249	3.89	19.48
French Canadian	314	2.89	16.60
Ashkenazi Jewish	1,475	11.26	31.75
Admixed Jewish	445	2.75	15.50
Hispanics/Latinos	810	3.53	16.38
Hispanics/Latinos in California	573	4.10	17.11
Hispanics/Latinos in New Mexico	163	5.52	21.92
Hispanics/Latinos in Texas	177	6.27	23.65
Puerto Rico	350	8.01	26.23
African Americans South	761	3.34	19.56
African Americans North	420	2.94	15.90
East Asia	561	3.65	19.63
Southeast Asia	325	8.44	17.90
South Asia	389	10.42	14.82
Greater Middle East	93	9.01	17.16

694

695 **Table 1. Summary of Haplotype Clusters**

696 Cumulative runs of homozygosity (cROH) was calculated by summing the regions of continuous
697 homozygous segments. Cumulative IBD was determined by summing IBD segments of ≥ 3 cM
698 and filtering for only pairs ≥ 12 cM and ≤ 72 cM. Statistics were determined within haplotype
699 clusters, rather than across the ancestrally heterogeneous and imbalanced full network.