

1 **Genetic ancestry and population structure in the All of Us Research Program cohort**
2 Shivam Sharma^{1,2}, Shashwat Deepali Nagar¹, Priscilla Pemu^{3,4}, Stephan Zuchner^{4,5}, SEEC Consortium⁴,
3 Leonardo Mariño-Ramírez², Robert Meller^{3,4,*}, I. King Jordan^{1,*}

4
5 ¹School of Biological Sciences, Georgia Institute of Technology, Atlanta, Georgia, United States

6 ²National Institute on Minority Health and Health Disparities, National Institutes of Health, Bethesda,
7 Maryland

8 ³Morehouse School of Medicine, Atlanta, Georgia, United States

9 ⁴SEEC Investigators, University of Miami, Coral Gables, Florida, United States *A list of investigators is
10 given at the end of the paper.

11 ⁵University of Miami, Coral Gables, Florida, United States

12

13 *Corresponding authors:

14 Rob Meller

15 RMeller@msm.edu

16

17 I. King Jordan

18 king.jordan@biology.gatech.edu

19 **Abstract**

20 The NIH All of Us Research Program (*All of Us*) aims to build one of the world’s most diverse population
21 biomedical datasets in support of equitable precision medicine. For this study, we analyzed participant
22 genomic variant data to assess the extent of population structure and to characterize patterns of genetic
23 ancestry for the *All of Us* cohort (n=297,549). Unsupervised clustering of genomic principal component
24 analysis (PCA) data revealed a non-uniform distribution of genetic diversity and substantial population
25 structure in the *All of Us* cohort, with dense clusters of closely related participants interspersed among
26 less dense regions of genomic PC space. Supervised genetic ancestry inference was performed using
27 genetic similarity between *All of Us* participants and global reference population samples. Participants
28 show diverse genetic ancestry, with major contributions from European (66.4%), African (19.5%), Asian
29 (7.6%), and American (6.3%) continental ancestry components. Participant genetic similarity clusters
30 show group-specific genetic ancestry patterns, with distinct patterns of continental and subcontinental
31 ancestry among groups. We also explored how genetic ancestry changes over space and time in the
32 United States (US). African and American ancestry are enriched in the southeast and southwest regions
33 of the country, respectively, whereas European ancestry is more evenly distributed across the US. The
34 diversity of *All of Us* participants’ genetic ancestry is negatively correlated with age; younger participants
35 show higher levels of genetic admixture compared to older participants. Our results underscore the
36 ancestral genetic diversity of the *All of Us* cohort, a crucial prerequisite for genomic health equity.

37 **Introduction**

38 The biomedical research community has become increasingly aware of the genomics research gap,
39 whereby the vast majority of participants in genetics research cohorts are of European ancestry^{1, 2, 3}. The
40 Eurocentric bias in genomics research threatens to exacerbate health disparities, since discoveries made
41 with European ancestry cohorts may not transfer to diverse ancestry groups⁴. The NIH All of Us Research
42 Program (*All of Us*) is a large cohort study of people who live in the US that combines participant genomic,
43 phenotypic, and environmental data, with health-related outcome data gleaned from surveys and
44 electronic health records^{5, 6}. *All of Us* has emphasized the recruitment of participants from population
45 groups that are underrepresented in biomedical research in an effort to close the genomics research gap
46 and to ensure that the benefits of precision medicine are shared equitably among all people^{7, 8}.

47
48 *All of Us* demonstration projects are being used to describe and validate the initial genomic data release
49 and the cloud-based Researcher Workbench, where registered users can access and analyze participant
50 data⁹. The aim of this demonstration project was to characterize the patterns of population structure and
51 genetic ancestry among *All of Us* participants. Population structure refers to differences in the
52 frequencies of genetic variants (alleles) among different groups or populations within a species, and
53 population structure can be revealed by the presence of clusters of genetically similar individuals¹⁰.
54 Genetic ancestry is closely related to the concept of population structure, and it can be defined
55 conceptually, mechanistically, and operationally. Conceptually, genetic ancestry reflects the geographic
56 origins of an individual's ancestors^{11, 12, 13, 14}. Mechanistically, genetic ancestry has been defined as the
57 subset of genealogical paths through which an individual's DNA has been inherited from their ancestors¹⁵.
58 For any individual, only a subset of their genealogical ancestors contributes DNA to their genome.
59 Operationally, genetic ancestry is typically characterized via genetic similarity between query individuals
60 (e.g. *All of Us* participants) and individuals from global reference populations, which are taken as
61 surrogates for ancestral populations^{16, 17, 18, 19}.

62
63 For this demonstration study of the *All of Us* cohort, we analyzed participant genomic variant data to (1)
64 assess the extent of population structure in the cohort, (2) characterize the patterns of participant genetic
65 ancestry at continental and subcontinental levels, and (3) explore how participants' genetic ancestry
66 changes over space and time in the US. Our results reveal substantial population structure and
67 heterogeneous patterns of genetic ancestry among *All of Us* participants, consistent with the consortium's
68 efforts to recruit a diverse participant cohort.

69 **Materials and Methods**

70 ***All of Us participant cohort, consent, and IRB review***

71 This study was performed as an *All of Us* genomic data demonstration project⁵. *All of Us* demonstration
72 projects are intended to describe and validate data and analysis tools for the participant cohort. Details
73 on the initial *All of Us* data release and Researcher Workbench used for this study were previously
74 published⁶. The genomic data demonstration project and experimental protocols were approved by the
75 *All of Us* Institutional Review Board (#2016–05-TN-Master), and informed consent was obtained from all
76 participants. *All of Us* inclusion criteria include adults 18 and older, with the legal authority and decisional
77 capacity to consent, and currently residing in the US or a territory of the US. *All of Us* exclusion criteria
78 exclude minors under the age of 18 and vulnerable populations (prisoners and individuals without the
79 capacity to give consent). Details on participant recruitment, informed consent, inclusion and exclusion
80 criteria are available online at https://allofus.nih.gov/sites/default/files/All_of_Us_operational_protocol_v1.7_mar_2018.pdf. Results reported here comply with the *All of Us* Data and
81 Statistics Dissemination Policy disallowing disclosure of group counts under 20.
82

83
84 The *All of Us* Researcher Workbench was used to build the participant cohort for this study
85 (Supplementary Figure 1). The cohort was built from the *All of Us* Controlled Tier dataset v7 (curated
86 version C2022Q4R9), which includes participants enrolled from 2018-2022, with a data cutoff date of
87 7/1/2022. Participants who self-identified as American Indian or Alaska Native were not included in the
88 analysis.
89

90 ***Unsupervised genetic clustering analysis***

91 Participant genomic data were accessed from the Controlled Tier dataset. Genome-wide genotypes for
92 *All of Us* participants were characterized using the Illumina Global Diversity Array with variants called for
93 1,824,517 genomic positions on the GRCh38/hg38 reference genome build. *All of Us* participant variants
94 were merged and harmonized with whole genome sequence variant data from 3,433 global reference
95 samples characterized as part of the 1000 Genomes Project (1KGP; phase 3) and the Human Genome
96 Diversity Project (HGDP; Supplementary Table 1)^{20, 21}. Biallelic variants common to the *All of Us* and
97 reference data sets were merged, with strand flips and variant identifier inconsistencies harmonized as
98 needed. Variants with >1% missingness and <1% minor allele frequency were removed from the merged
99 and harmonized dataset. Linkage disequilibrium (LD) pruning was done using window size=50, step
100 size=10, and pairwise threshold $r^2 < 0.1$, yielding a final *All of Us* and global reference sample dataset of
101 187,795 variants. Variant merging, harmonization, and LD pruning were performed using PLINK version
102 1.9²² and custom scripts as previously described^{23, 24, 25}. The final dataset of *All of Us* participant genomic
103 variants was used for unsupervised clustering analysis. Principal Component Analysis (PCA) was run on
104 the variant dataset using the FastPCA program implemented in PLINK version 2.0. The clustering tendency
105 of the resulting genomic PCA data was analyzed using the Hopkins statistic with the hopkins R package²⁶
106 and nearest neighbor search with the FNN R package version 1.1.4²⁷. Kernel density estimation was
107 performed with the MASS R package using PCs 1-3 and contour lines were extracted from the estimated
108 density distribution²⁸. Density-based clustering was performed using the HDBSCAN algorithm²⁹. HDBSCAN
109 was run on first 5 PCs for the PCA data with parameters min_samples=2,000 and min_cluster_size=2,500.
110 Cluster boundaries were visualized using the ggforce R package.

111 **Supervised genetic ancestry inference**

112 Genomic variants from *All of Us* participants and a set of four global reference populations were merged
113 and harmonized as described in the previous section to perform continental and subcontinental genetic
114 ancestry inference. Kinship analysis was performed with the KING program to eliminate related (or
115 duplicated) reference samples from the global reference populations³⁰. Continental genetic ancestry
116 inference was performed using a subset of 1,572 global reference samples from the 1KGP and the HGDP,
117 which were selected as non-admixed representatives of seven ancestry groups: African, American, East
118 Asian, South Asian, West Asian, European, and Oceanian (Supplementary Table 1). K-nearest neighbor
119 clustering of genomic PCA data was used to identify *All of Us* participants that cluster together with
120 African, East Asian, South Asian, and European reference populations, and these participants were used
121 for subcontinental ancestry inference³¹. West Asian and Oceanian reference populations were not used
122 for this purpose owing to the relatively low number of participants that clustered with these groups. Asian
123 and European reference populations for subcontinental ancestry inference were taken from the 1KGP and
124 HGDP (Supplementary Table 2). 1KGP and HGDP reference populations were used together with
125 additional reference populations to provide broader geographic coverage for African and American
126 subcontinental ancestry inference (Supplementary Table 2). African reference samples were taken from
127 a study of Bantu-speaking populations in Africa that included samples from 53 populations from east,
128 central, south, and west Africa³². The merged and harmonized African subcontinental ancestry inference
129 panel included 1,659 reference samples and 228,033 variants.

130
131 Continental and subcontinental ancestry inference was performed via analysis of merged *All of Us*
132 participant and global reference population genomic variant sets with the program Rye (Rapid ancestry
133 Estimation)³³. Rye performs rapid and accurate genetic ancestry inference based on principal component
134 analysis (PCA) of genomic variant data. PCA was run on the merged variant datasets using the FastPCA
135 program implemented in PLINK version 2.0, and Rye was then run on the first 25 PCs, using the defined
136 reference ancestry groups to assign ancestry group fractions to individual *All of Us* participant samples.
137 The continuous ancestry fractions that we report here were calculated independently of the categorical
138 ancestry predictions currently provided by the *All of Us* Researcher Workbench³⁴.

139
140 *All of Us* participant continental ancestry fractions were visualized as admixture-style plots at the state (or
141 territory) level using the geofacets R package^{35, 36}. Admixture entropy (AE) was used to quantify the
142 amount of genetic admixture for *All of Us* participants as previously described^{25, 37}: $AE_i =$
143 $-\sum_{j=1}^7 p_j \log(p_j)$, where p_j is the fraction of ancestry group j for individual i .

144
145 **Note on genetic ancestry inference**

146 As discussed in the introduction, genetic ancestry can be defined conceptually, mechanistically, and
147 operationally. We use an operational definition of genetic ancestry for *All of Us* participants in this study,
148 as measured by their levels of genetic similarity with global reference population samples^{16, 17}.
149 Accordingly, the phrase ‘African ancestry’ is used here as shorthand for similarity to African reference
150 population samples, ‘European ancestry’ is used for similarity to European reference population samples
151 and so on. ‘American ancestry’ refers to genetic similarity in Indigenous American reference population
152 samples. The relative levels of similarity to different reference population groups allows us to infer
153 percent ancestry components for *All of Us* participants³³. The genetic ancestry results reported here are
154 contingent upon the choice of reference populations, how these reference populations are delineated,
155 and the method used to infer genetic similarity between *All of Us* participants and the reference
156 population samples.

157 **Results**

158 ***Unsupervised: population structure***

159 A cohort of 297,549 *All of Us* participants, for whom genomic data are available, was created using the *All*
160 *of Us* Researcher Workbench (Supplementary Figure 1). *All of Us* participant genetic diversity was
161 analyzed using PCA of genomic variant data followed by unsupervised clustering to assess the extent of
162 population structure in the cohort. The clustering tendency of participant genomic PCA data was
163 evaluated using the Hopkins statistic, nearest neighbors, and kernel density estimation. The PCA data
164 yield a Hopkins statistic value of ~ 1 , indicating highly clustered, non-uniformly, and non-randomly
165 distributed genomic PCA data. The numbers of close neighbors per participant are highly variable across
166 PC space, and kernel density estimation shows a multimodal distribution with distinct peaks separated in
167 PC space (Figure 1A and 1B). All three of these metrics reveal highly clustered participant genomic data,
168 with dense groups of genetically similar individuals interspersed among less dense regions, indicative of
169 substantial population structure in the *All of Us* cohort.

170
171 Density-based clustering of the genomic PCA data yield an optimal number of $K=7$ genetic diversity
172 clusters (Figure 1C). Similar clustering was performed using a Uniform Manifold Approximation and
173 Projection (UMAP) analysis of the genomic PCA data (Supplementary Methods). Density-based clustering
174 of UMAP data reveals almost twice as many clusters ($K=13$) as seen for the PCA data, but there is broad
175 concordance between the two methods with high percentages of participant overlap for each PCA cluster
176 within one or two corresponding UMAP clusters (Supplementary Figure 2). The number of *All of Us*
177 genetic diversity clusters could change with future participant data releases.

178 179 ***Supervised: genetic ancestry***

180 *All of Us* participant genetic ancestry was inferred using genomic PCA data analyzed with the Rye (Rapid
181 ancestry Estimation) program³³. Participant PCA data were compared with PCA data from global
182 reference populations, taken from the 1KGP and the HGDP, to infer individual ancestry proportions from
183 seven continental-level ancestry groups: African, American, East Asian, South Asian, West Asian,
184 European, and Oceanian (Supplementary Table 1 and Supplementary Figure 3). *All of Us* participants are
185 broadly distributed in PC space, whereas global reference samples from different ancestry groups are
186 tightly clustered in PC space (Figure 2A and 2B). Rye infers *All of Us* participant genetic ancestry
187 proportions as linear combinations of reference population ancestries. Overall, the *All of Us* participant
188 cohort shows 19.51% African, 6.33% American, 2.57% East Asian, 3.05% South Asian, 1.95% West Asian,
189 66.37% European, and 0.21% Oceanian ancestry. The *All of Us* participant genetic similarity groups
190 inferred with density-based clustering show group-specific patterns of ancestry proportions, with a
191 continuum of ancestry proportions within and between groups (Figure 2C). Groups 1, 3, 4, and 7 show
192 the most uniform patterns of ancestry within groups, whereas groups 2, 5, 6, and the remaining
193 participants that did not fall into any density-based cluster show more diverse patterns of ancestry and
194 admixture. All groups show evidence of admixture with multiple ancestry components present in
195 different proportions.

196
197 The *All of Us* Researcher Workbench predicts participant membership among six continental ancestry
198 groups, using a PCA-based machine learning method that is distinct from the continuous ancestry
199 inference approach used here³⁴. We compared the participant continental ancestry percentages inferred
200 here to the Researcher Workbench assigned categorical ancestry groups (Supplementary Figure 4). Five
201 of the six categorical ancestry groups correspond exactly with the reference population groups we use:
202 African, East Asian, South Asian, Middle Eastern (West Asian here), and European. For these five groups,
203 there is high correspondence between participants' PCA-based machine learning predicted group
204 membership and averages for the ancestry percentages that we inferred (83.02-97.71% matching

205 ancestry). The Admixed American ancestry category from the Researcher Workbench includes modern,
206 admixed reference samples from Latin America, whereas our American reference population group
207 includes Indigenous American samples only (Supplementary Table 1). The Admixed American group
208 shows 51.01% European ancestry and 35.84% American ancestry, consistent with what is expected for
209 modern Latin American populations^{38, 39}.

210
211 We also used Rye to infer subcontinental ancestry for *All of Us* participants with high levels of African
212 (n=9,291), East Asian (n=2,457), South Asian (n=2,484), and European ancestry (n=24,730; Figure 3). The
213 relationships among the reference populations used for subcontinental ancestry inference with Rye and
214 *All of Us* participants are shown in Supplementary Figures 5-7. African subcontinental ancestry is
215 characterized by a predominant West Central African component, followed by West African and Bantu
216 components. East Asian subcontinental ancestry is highly diverse with predominant Han (Chinese),
217 Japanese, and Southeast Asian components. South Asian subcontinental ancestry is mainly South Indian
218 followed by North Indian and a small Central Asian component. European subcontinental ancestry is
219 made up primarily of British ancestry followed by Italian and Iberian components.

220
221 **Genetic ancestry by geography and age**
222 *All of Us* participant continental ancestry percentages were visualized across fifty states and Puerto Rico
223 to evaluate the geographic distribution of ancestry across the US (Figure 4). African ancestry is
224 concentrated primarily in the southeast part of the country, whereas American ancestry is found primarily
225 in the southwest and California. European ancestry is more uniformly distributed across the country, with
226 the highest concentrations found in north, along the Canadian border. Relatively high levels of admixture
227 are seen in the northeast, Florida, and Hawaii.

228
229 The relationship between *All of Us* participants' age and genetic ancestry was assessed using genetic
230 admixture entropy, where higher values indicate a more diverse combination of ancestry components
231 within individual genomes and lower values indicate more homogenous ancestry (Figure 5). Genetic
232 admixture entropy is negatively correlated with participant age, indicating that younger participants have
233 more diverse ancestry combinations than older participants.

234
235 **Discussion**
236 Our analysis demonstrates the genomic and ancestral diversity of the *All of Us* cohort, consistent with the
237 project's goals to recruit participants from population groups that are underrepresented in biomedical
238 research in support of health equity. Indeed, *All of Us* is one of the most diverse population biomedical
239 datasets in the world, and this represents an important step towards making precision medicine more
240 widely available and more applicable to diverse communities in the US^{7, 8, 40}. The promise of population
241 biomedical datasets like *All of Us* rests on the integration of genetic, social, environmental, and health
242 outcome data for many thousands of diverse participants. Given that genetic ancestry is derived from the
243 genome, it should be possible to use genetic ancestry inference, together with population biomedical
244 datasets, to help elucidate genetic and socioenvironmental contributions to health outcomes and
245 disparities.

246
247 One challenge is that current methods for genetic ancestry inference, while accurate, are slow and do not
248 scale to biobank sized datasets like *All of Us*. We developed the Rye algorithm as a fast and
249 computationally efficient genetic ancestry inference method that can scale to biobank sized genomic data
250 sets³³. Application of Rye to genome-wide genetic data for 297,549 *All of Us* participants underscores its
251 utility for this purpose. Using Rye, we found the *All of Us* cohort to be ancestrally diverse with distinct
252 patterns of genetic ancestry and admixture among genetic similarity groups and geographic regions

253 (Figures 2-4). The geographic patterns of genetic ancestry seen for the *All of Us* cohort are consistent with
254 previous studies and could also reflect differences in participant recruitment across the country^{41, 42, 43}.

255
256 The extent to which human genetic diversity is characterized by clusters of closely related individuals, i.e.
257 population structure, versus clines of continuous genetic variation has long been a subject of interest^{44, 45,}
258 ^{46, 47, 48}. The *All of Us* cohort allows for an assessment of the extent of population structure in the US given
259 the large size of the cohort, the extensive sampling of participants across the country, and the
260 demographic diversity of the participants. The application of several different cluster analysis methods
261 to participants' genomic PCA data revealed evidence for substantial population structure in the cohort,
262 with dense clusters of relatively closely related participants interspersed among less dense regions in PC
263 space (Figure 1). The population structure and genetic clusters that can be gleaned from clustering
264 analysis of genomic PCA data are not readily apparent from visual inspection of these same data, owing
265 to large size of the cohort and over-plotting of participants in dense regions of PC space (Figure 2A).

266
267 Finally, we show that genetic diversity in the US is increasing over time. Younger *All of Us* participants are
268 far more ancestrally diverse than older participants, and this trend is evident across the entire age range
269 of the cohort. This finding suggests that genetic ancestry categories and group designations will become
270 increasingly obsolete over time⁴⁹.

271

272 **Acknowledgements**

273 We thank our colleagues, Kelsey Mayo, Ashley Able, Ashley Green, Andrea Ramirez, Anji Musick and Sokny
274 Lim for providing their support and input throughout the demonstration project lifecycle. We thank
275 Jennifer Zhang for providing input on the project's code review. We thank Lee Lichtenstein and Jennifer
276 Zhang for providing the data artifacts used for the project. We thank the DRC's Research Support team
277 for their help during implementation. We also thank the All of Us Science Committee and All of Us Steering
278 Committee for their efforts in evaluating and finalizing the approved demonstration projects. The All of
279 Us Research Program would not be possible without the partnership of contributions made by its
280 participants. To learn more about the All of Us Research Program's research data repository, please visit
281 <https://www.researchallofus.org/>.

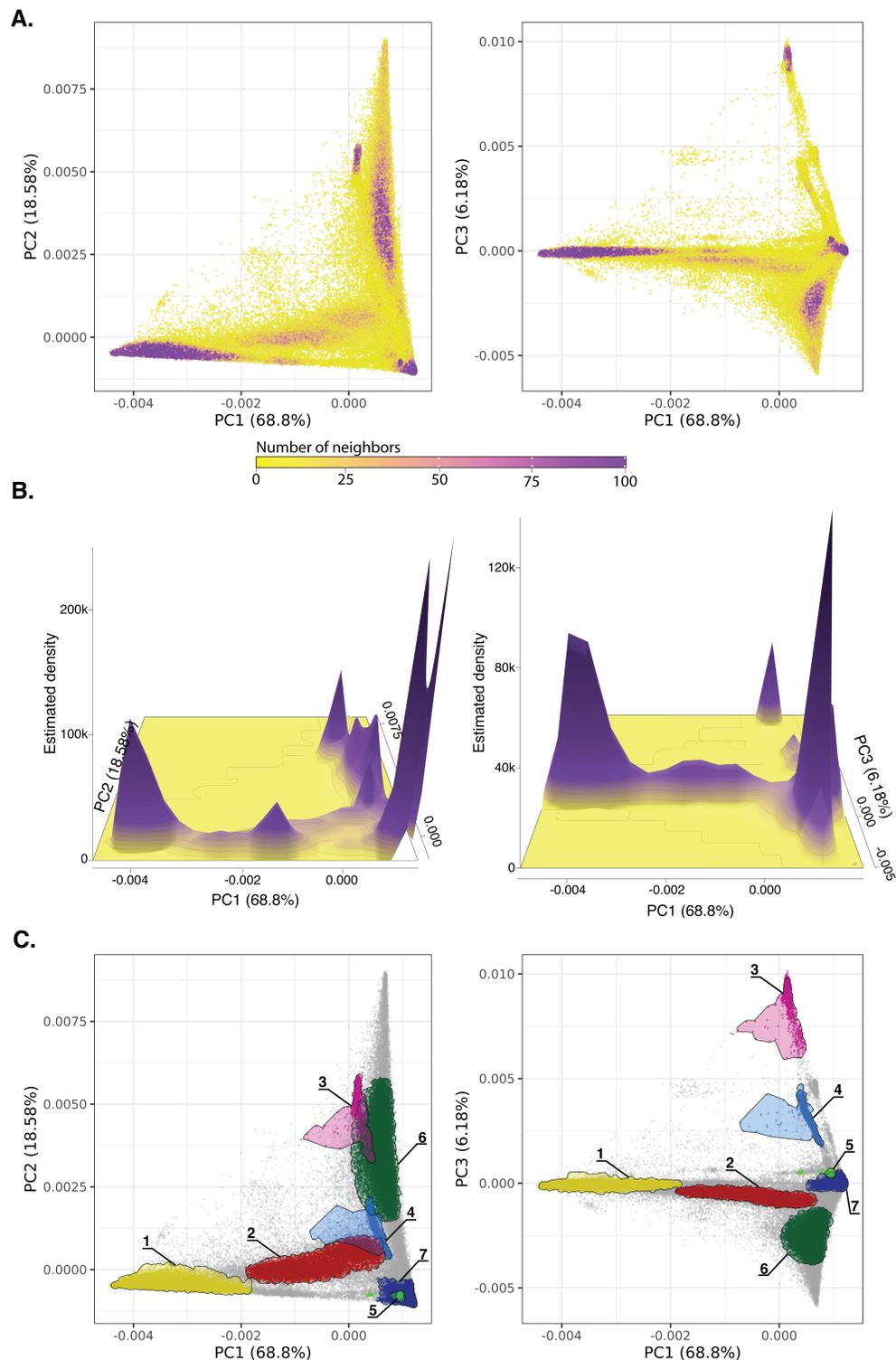
282 The All of Us Research Program is supported by the National Institutes of Health, Office of the Director:
283 Regional Medical Centers: 1 OT2 OD026549; 1 OT2 OD026554; 1 OT2 OD026557; 1 OT2 OD026556; 1 OT2
284 OD026550; 1 OT2 OD 026552; 1 OT2 OD026548; 1 OT2 OD026551; 1 OT2 OD026555; IAA#: AOD 16037;
285 Federally Qualified Health Centers: 75N98019F01202.; Data and Research Center: 1 OT2 OD35404;
286 Biobank: 1 U24 OD023121; The Participant Center: U24 OD023176; Participant Technology Systems
287 Center: 1 OT2 OD030043; Community Partners: 1 OT2 OD025277; 3 OT2 OD025315; 1 OT2 OD025337; 1
288 OT2 OD025276. In addition, the All of Us Research Program would not be possible without the partnership
289 of its participants.

290 **References**

- 291 1. Bustamante CD, Burchard EG, De la Vega FM. Genomics for the world. *Nature* **475**, 163-165
292 (2011).
- 293 2. Petrovski S, Goldstein DB. Unequal representation of genetic variation across ancestry groups
294 creates healthcare inequality in the application of precision medicine. *Genome Biol* **17**, 157
295 (2016).
- 296 3. Popejoy AB, Fullerton SM. Genomics is failing on diversity. *Nature* **538**, 161-164 (2016).
- 297 4. Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ. Clinical use of current polygenic
298 risk scores may exacerbate health disparities. *Nat Genet* **51**, 584-591 (2019).
- 299 5. All of Us Research Program I, *et al.* The "All of Us" Research Program. *N Engl J Med* **381**, 668-676
300 (2019).
- 301 6. Ramirez AH, *et al.* The All of Us Research Program: Data quality, utility, and diversity. *Patterns (N*
302 *Y)* **3**, 100570 (2022).
- 303 7. Bianchi DW, *et al.* The All of Us Research Program is an opportunity to enhance the diversity of
304 US biomedical research. *Nat Med* **30**, 330-333 (2024).
- 305 8. Kathiresan N, Cho SMJ, Bhattacharya R, Truong B, Hornsby W, Natarajan P. Representation of
306 race and ethnicity in the contemporary US health cohort All of Us Research Program. *JAMA*
307 *Cardiol* **8**, 859-864 (2023).
- 308 9. All of Us Research Program Genomics I. Genomic data in the All of Us Research Program. *Nature*
309 **627**, 340-346 (2024).
- 310 10. Pritchard JK. An Owner's Guide to the Human Genome: an introduction to human population
311 genetics, variation and disease.) (2023).
- 312 11. Hellenthal G, *et al.* A genetic atlas of human admixture history. *Science* **343**, 747-751 (2014).
- 313 12. Nielsen R, Akey JM, Jakobsson M, Pritchard JK, Tishkoff S, Willerslev E. Tracing the peopling of
314 the world through genomics. *Nature* **541**, 302-310 (2017).
- 315 13. Royal CD, *et al.* Inferring genetic ancestry: opportunities, challenges, and implications. *Am J Hum*
316 *Genet* **86**, 661-673 (2010).
- 317 14. Wohns AW, *et al.* A unified genealogy of modern and ancient genomes. *Science* **375**, eabi8264
318 (2022).
- 319 15. Mathieson I, Scally A. What is ancestry? *PLoS Genet* **16**, e1008624 (2020).
- 320 16. Coop G. Genetic similarity versus genetic ancestry groups as sample descriptors in human
321 genetics}. *arXiv* **2207.11595**, (2023).
- 322 17. National Academies of Sciences Engineering and Medicine. *Using Population Descriptors in*
323 *Genetics and Genomics Research: A New Framework for an Evolving Field*. The National
324 Academies Press (2023).
- 325 18. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated
326 individuals. *Genome Res* **19**, 1655-1664 (2009).
- 327 19. Maples BK, Gravel S, Kenny EE, Bustamante CD. RFMix: a discriminative modeling approach for
328 rapid and robust local-ancestry inference. *Am J Hum Genet* **93**, 278-288 (2013).
- 329 20. Bergstrom A, *et al.* Insights into human genetic variation and population history from 929
330 diverse genomes. *Science* **367**, (2020).
- 331 21. Genomes Project C, *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74
332 (2015).
- 333 22. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to
334 the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
- 335 23. Conley AB, *et al.* A Comparative Analysis of Genetic Ancestry and Admixture in the Colombian
336 Populations of Choco and Medellin. *G3 (Bethesda)* **7**, 3435-3447 (2017).

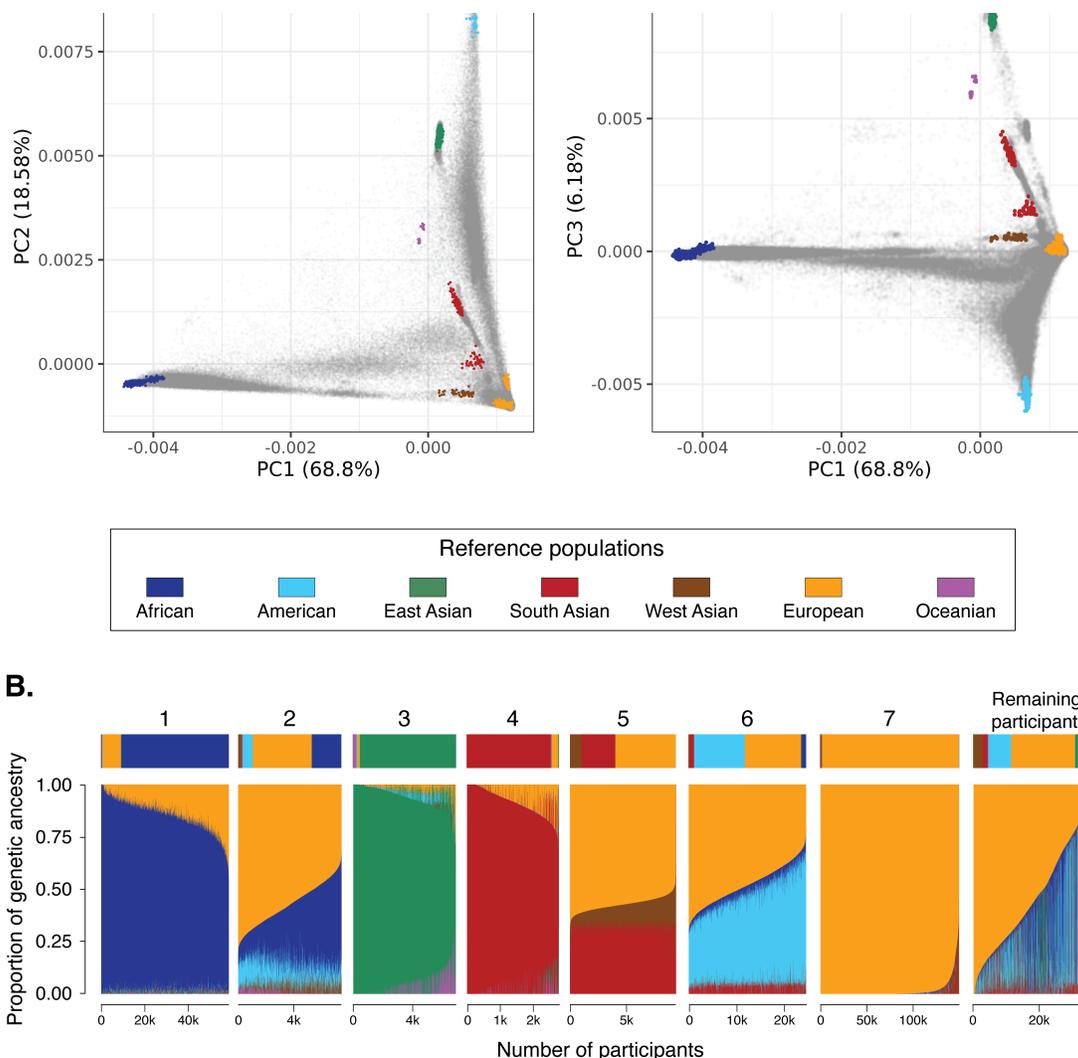
- 337 24. Jordan IK, Rishishwar L, Conley AB. Native American admixture recapitulates population-specific
338 migration and settlement of the continental United States. *PLoS Genet* **15**, e1008225 (2019).
- 339 25. Nagar SD, *et al.* Genetic ancestry and ethnic identity in Ecuador. *HGG Adv* **2**, 100050 (2021).
- 340 26. Hopkins B, Skellam JG. A new method for determining the type of distribution of plant
341 individuals. *Annals of Botany* **18**, 213-227 (1954).
- 342 27. Beygelzimer A, Kakadet S, Langford J, Arya S, Mount D, Li S. FNN: Fast nearest neighbor search
343 algorithms and applications. (2024).
- 344 28. Venables WN, Ripley BD. *Modern Applied Statistics with S*, Fourth edn. Springer (2002).
- 345 29. McInnes L, Healy J, Astels S. hdbSCAN: Hierarchical density based clustering. *J Open Source Softw*
346 **2**, 205 (2017).
- 347 30. Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust relationship inference
348 in genome-wide association studies. *Bioinformatics* **26**, 2867-2873 (2010).
- 349 31. Pedregosa F, *et al.* Scikit-learn: Machine learning in Python. *Journal of machine Learning*
350 *research* **12**, 2825-2830 (2011).
- 351 32. Patin E, *et al.* Dispersals and genetic adaptation of Bantu-speaking populations in Africa and
352 North America. *Science* **356**, 543-546 (2017).
- 353 33. Conley AB, *et al.* Rye: genetic ancestry inference at biobank scale. *Nucleic Acids Res* **51**, e44
354 (2023).
- 355 34. All of Us Research Program. Genomic Research Data Quality Report.) (2022).
- 356 35. Bivand R, Keitt T, Rowlingson B. rgdal: Bindings for the 'Geospatial' Data Abstraction Library.
357 (2023).
- 358 36. Hafen R. geofacet: 'ggplot2' Faceting Utilities for Geographical Data. (2023).
- 359 37. Medina-Rivas MA, *et al.* Choco, Colombia: a hotspot of human biodiversity. *Rev Biodivers*
360 *Neotrop* **6**, 45-54 (2016).
- 361 38. Homburger JR, *et al.* Genomic insights into the ancestry and demographic history of South
362 America. *PLoS Genet* **11**, e1005602 (2015).
- 363 39. Ruiz-Linares A, *et al.* Admixture in Latin America: geographic structure, phenotypic diversity and
364 self-perception of ancestry based on 7,342 individuals. *PLoS Genet* **10**, e1004572 (2014).
- 365 40. Abul-Husn NS, Kenny EE. Personalized medicine and the power of electronic health records. *Cell*
366 **177**, 58-69 (2019).
- 367 41. Dai CL, *et al.* Population histories of the United States revealed through fine-scale migration and
368 haplotype analysis. *Am J Hum Genet* **106**, 371-388 (2020).
- 369 42. Han E, *et al.* Clustering of 770,000 genomes reveals post-colonial population structure of North
370 America. *Nat Commun* **8**, 14238 (2017).
- 371 43. Bryc K, Durand EY, Macpherson JM, Reich D, Mountain JL. The genetic ancestry of African
372 Americans, Latinos, and European Americans across the United States. *Am J Hum Genet* **96**, 37-
373 53 (2015).
- 374 44. Serre D, Paabo S. Evidence for gradients of human genetic diversity within and among
375 continents. *Genome Res* **14**, 1679-1685 (2004).
- 376 45. Rosenberg NA, *et al.* Genetic structure of human populations. *Science* **298**, 2381-2385 (2002).
- 377 46. Rosenberg NA, Mahajan S, Ramachandran S, Zhao C, Pritchard JK, Feldman MW. Clines, clusters,
378 and the effect of study design on the inference of human population structure. *PLoS Genet* **1**,
379 e70 (2005).
- 380 47. Mountain JL, Cavalli-Sforza LL. Multilocus genotypes, a tree of individuals, and human
381 evolutionary history. *Am J Hum Genet* **61**, 705-718 (1997).
- 382 48. Bowcock AM, Ruiz-Linares A, Tomfohrde J, Minch E, Kidd JR, Cavalli-Sforza LL. High resolution of
383 human evolutionary trees with polymorphic microsatellites. *Nature* **368**, 455-457 (1994).

- 384 49. Lewis ACF, *et al.* Getting genetic ancestry right for science and society. *Science* **376**, 250-252
385 (2022).
386



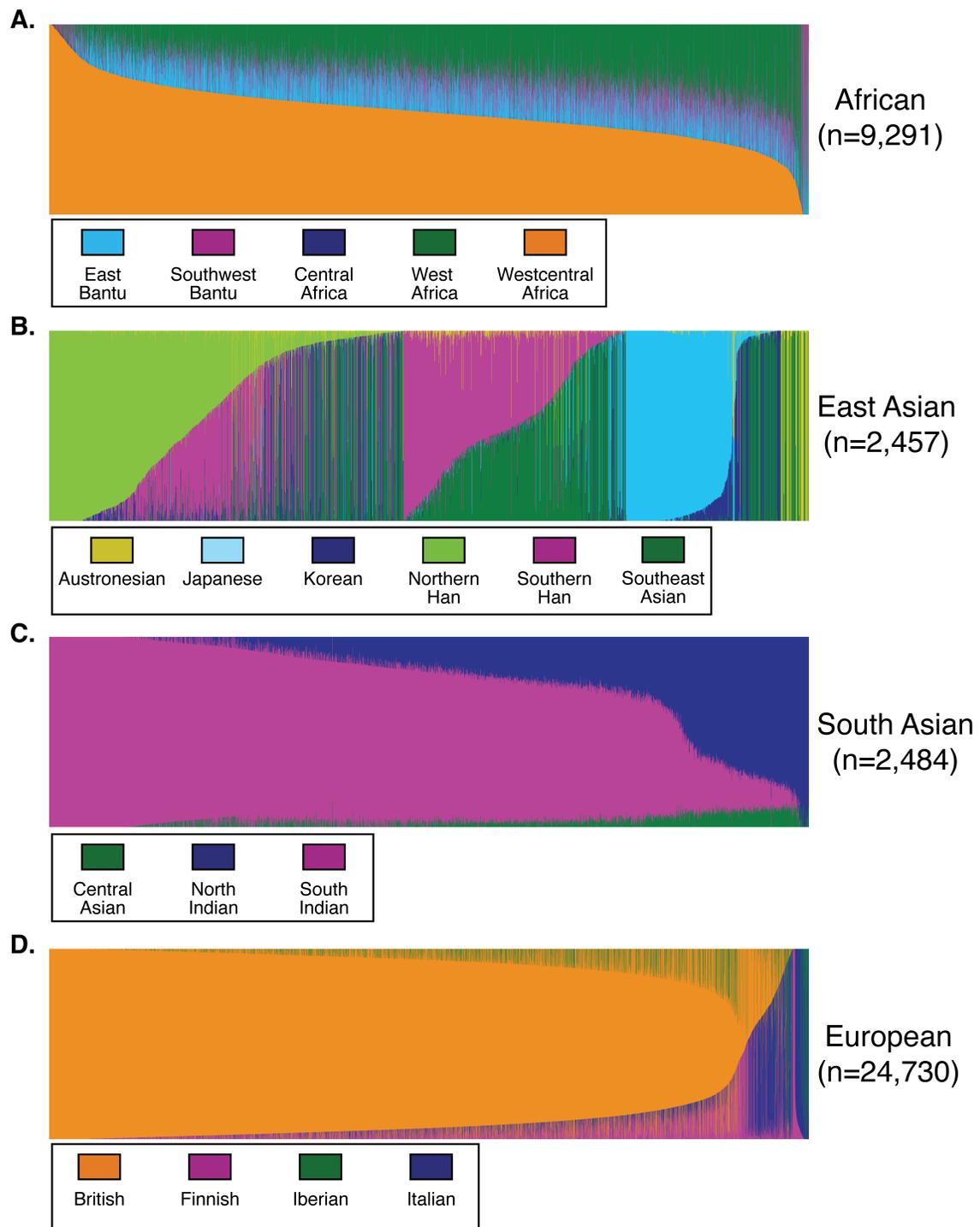
387

388 **Figure 1. Population structure.** Genomic PCA for *All of Us* participants. Left panels show PC1 versus PC2
389 comparisons, and right panels show PC1 versus PC3 comparisons, with the percent of variance explained
390 by each PC shown. (A) Participants color-coded by the number of close neighbors as defined by Euclidean
391 distance < 0.1 in PCs 1-5. (B) Kernel density estimation with peaks showing high density clusters of
392 participants in PC space. (C) High density clusters of genetically similar participants shown as groups 1-7.



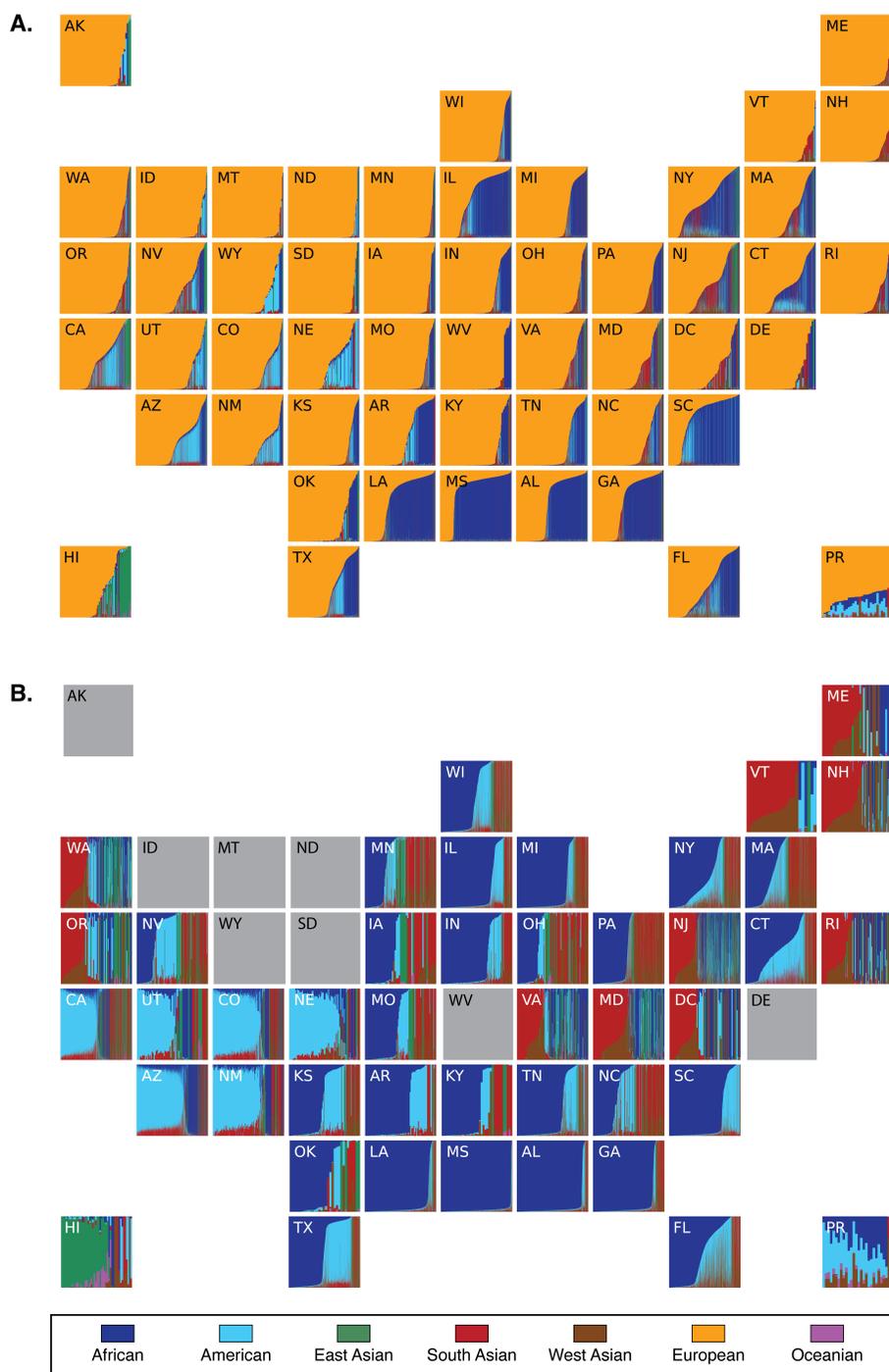
393

394 **Figure 2. Continental genetic ancestry.** (A) Genomic PCA with *All of Us* participants shown in gray and
 395 global reference population samples color-coded as shown in the key. Left panels show PC1 versus PC2
 396 comparisons, and right panels show PC1 versus PC3 comparisons, with the percent of variance explained
 397 by each PC shown. (B) Genetic ancestry proportions for *All of Us* participants stratified by the genetic
 398 similarity groups shown in Figure 1C. Average ancestry proportions are shown above each group, and
 399 numbers of participants are shown below each group. The remaining participants are individuals that did
 400 not fall into a dense PCA cluster.



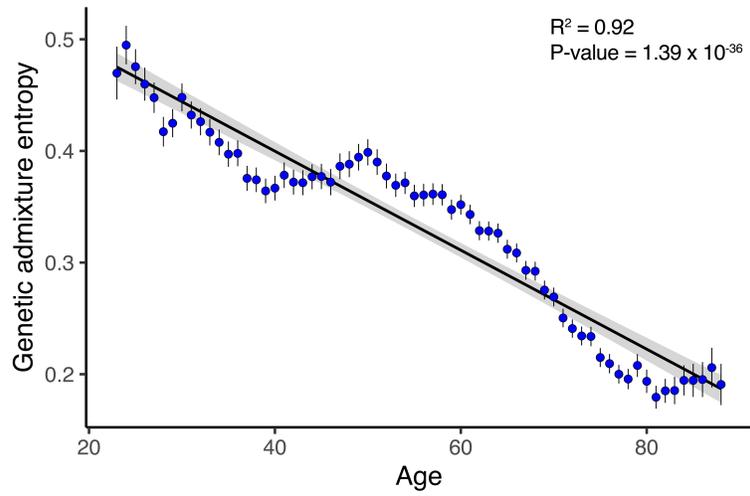
401

402 **Figure 3. Subcontinental genetic ancestry.** Subcontinental genetic ancestry proportions for *All of Us*
403 participants from African, American, East Asian, South Asian, and European continental ancestry groups.
404 Subcontinental groups (regions) for each continental ancestry group are color-coded as shown.



405

406 Figure 4. **Genetic ancestry by geography.** Genetic ancestry proportions are shown for *All of Us*
 407 participants sampled from the fifty US states and Puerto Rico. (A) All participants and ancestry
 408 components. (B) Non-European genetic ancestry proportions for all individuals with <90% European
 409 ancestry. The results for states shaded in grey are suppressed owing to <20 participants with <90%
 410 European ancestry.



411

412 Figure 5. **Genetic admixture by age.** Genetic admixture entropy (y-axis) against participant age (x-axis).

413 Ages shown in 100 bins with average and 95% CI values shown. Linear regression trend line shown with

414 95% CI shaded.

415 **SEEC Investigators**

416 Priscilla E. Pemu, MD, MS¹. Robert Meller DPhil, BSc¹. Alexander Quarshie, MD¹, MS. Kelley Carroll, MD¹.
417 Lawrence L. Sanders, MD¹. Howard Mosby, CPA, CGMA¹. Elizabeth I. Olorundare, MD, MPH¹. Atuarra
418 McCaslin, BS¹. Chadrick Anderson, MHA. Andrea Pearson¹. Kelechi C. Igwe, OD, MPH¹. Karunamuni
419 Silva¹. Gwen Daugett, PMP. Jason McCray. Michael Prude. Cheryl Franklin, MD, MPH, FACOG¹. Stephan
420 Zuchner, M.D. Ph.D². Olveen Carrasquillo, M.D. ². Rosario Isasi, JD. ², MPH. Jacob L. McCauley, PhD². Jose
421 G Melo, MSPH². Ana K Riccio, M.D. ². Patrice Whitehead, MB (ASCP) ². Patricia Guzman, MS². Christina
422 Gladfelter ². Rebecca Velez, MA. ², Mario Saporta, MD. ² Brandon Apagüño², MS. Lisa Abreu, MPH².
423 Betsy Shenkman³. Bill Hogan³. Eileen Handberg³. Jamie Hensley³. Sonya White³. Brittney Roth-Manning³.
424 Tona Mendoza³. Alex Loiacono³. Donny Weinbrenner³. Mahmoud Enani³. Ali Nouina³. Michael E. Zwick,
425 Ph.D.⁴, Tracie C. Rosser, Ph.D. ⁴, Arshed A. Quyyumi, M.D. ⁴, Theodore M. Johnson II, M.D., MPH⁴, Greg S.
426 Martin, M.D., M.Sc⁴, Alvaro Alonso, M.D., Ph.D.⁴, Tina-Ann Kerr Thompson, M.D⁴, Nita Deshpande, Ph.D.
427 ⁴, H. Richard Johnston, Ph.D. ⁴, Hina Ahmed, MPH⁴, Letheshia Husbands, MPH⁴.

428

429 **Affiliations**

- 430 1. Morehouse School of Medicine, Atlanta, Georgia, United States
431 2. University of Miami, Coral Gables, Florida, United States
432 3. University of Florida, Gainesville, FL
433 4. Emory University, Atlanta, GA.