

## Genetics of Latin American Diversity (GLAD) Project: insights into population genetics and association studies in recently admixed groups in the Americas

Victor Borda<sup>1,\*</sup>, Douglas P. Loesch<sup>1</sup>, Bing Guo<sup>1</sup>, Roland Laboulaye<sup>1</sup>, Diego Veliz-Otani<sup>1</sup>, Jennifer N. French-Kwawu<sup>1</sup>, Thiago Peixoto Leal<sup>2</sup>, Stephanie M. Gogarten<sup>3</sup>, Sunday Ikpe<sup>1</sup>, Mateus H. Gouveia<sup>4</sup>, Marla Mendes<sup>5</sup>, Gonçalo R. Abecasis<sup>6</sup>, Isabela Alvim<sup>5</sup>, Carlos E. Arboleda-Bustos<sup>7</sup>, Gonzalo Arboleda<sup>7</sup>, Humberto Arboleda<sup>7</sup>, Mauricio L. Barreto<sup>8</sup>, Lucas Barwick<sup>9</sup>, Marcos A. Bezzerá<sup>10</sup>, John Blangero<sup>11</sup>, Vanderci Borges<sup>12</sup>, Omar Caceres<sup>13,14</sup>, Jianwen Cai<sup>15</sup>, Pedro Chana-Cuevas<sup>16</sup>, Zhanghua Chen<sup>17</sup>, Brian Custer<sup>18</sup>, Michael Dean<sup>19</sup>, Carla Dinardo<sup>20</sup>, Igor Domingos<sup>10</sup>, Ravindranath Duggirala<sup>11</sup>, Elena Dieguez<sup>21</sup>, Willian Fernandez<sup>7</sup>, Henrique B. Ferraz<sup>12</sup>, Frank Gilliland<sup>17</sup>, Heinner Guio<sup>13,22,23</sup>, Bernardo Horta<sup>24</sup>, Joanne E. Curran<sup>11</sup>, Jill M. Johnsen<sup>25</sup>, Robert C. Kaplan<sup>26,27</sup>, Shannon Kelly<sup>18,28</sup>, Eimear E. Kenny<sup>29</sup>, Stephen Kittner<sup>30</sup>, Barbara A. Konkle<sup>31</sup>, Charles Kooperberg<sup>27</sup>, Andres Lescano<sup>21</sup>, M. Fernanda Lima-Costa<sup>32</sup>, Ruth J. F. Loos<sup>33</sup>, Ani Manichaikul<sup>34</sup>, Deborah A. Meyers<sup>35</sup>, Braxton D. Mitchell<sup>36,37</sup>, Michel S. Naslavsky<sup>38</sup>, Deborah A. Nickerson<sup>39,§</sup>, Kari E. North<sup>40</sup>, Carlos Padilla<sup>13</sup>, Michael Preuss<sup>29</sup>, Victor Raggio<sup>41</sup>, Alexander P. Reiner<sup>27,42</sup>, Stephen S. Rich<sup>34</sup>, Carlos R. Rieder<sup>43</sup>, Michiel Rienstra<sup>44</sup>, Jerome I. Rotter<sup>45</sup>, Tatjana Rundek<sup>46</sup>, Ralph L. Sacco<sup>46</sup>, Cesar Sanchez<sup>13</sup>, Bruno Lopes Santos-Lobato<sup>47</sup>, Artur Francisco Schumacher-Schuh<sup>48,49</sup>, Marilia O. Scliar<sup>38</sup>, Edwin K. Silverman<sup>50</sup>, Tamar Sofer<sup>51</sup>, Jessica Lasky-Su<sup>50</sup>, Vitor Tumas<sup>52</sup>, Scott T. Weiss<sup>50</sup>, Latin American Research Consortium on the Genetics of Parkinson's Disease (LARGE-PD), NINDS Stroke Genetics Network (SiGN) Consortium, TOPMed Population Genetics Working Group, Ignacio F. Mata<sup>2</sup>, Ryan D. Hernandez<sup>53,54,55,56</sup>, Eduardo Tarazona-Santos<sup>55,57</sup>, Timothy D. O'Connor<sup>1,36,58,59,\*</sup>

<sup>1</sup>Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD 21201, USA.

<sup>2</sup>Lerner Research Institute, Genomic Medicine, Cleveland Clinic, Cleveland, OH, USA.

<sup>3</sup>Department of Biostatistics, University of Washington, Seattle, WA, USA.

<sup>4</sup>Center for Research on Genomics and Global Health, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, 20892, USA.

<sup>5</sup>Department of Genetics, Ecology, and Evolution. Instituto de Ciências Biológicas, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil.

<sup>6</sup>Department of Biostatistics, University of Michigan School of Public Health, Ann Arbor, MI, USA.

<sup>7</sup>Neuroscience and Cell Death Research Groups, Medical School and Genetic Institute, Universidad Nacional de Colombia, Bogota, Colombia.

<sup>8</sup>Instituto de Saúde Coletiva, Universidade Federal da Bahia, Salvador, BA, 40110-040, Brazil.

<sup>9</sup>LTRC Data Coordinating Center, The Emmes Company, LLC, Rockville, MD, USA.

<sup>10</sup>Department of Genetics, Federal University of Pernambuco, Av. Prof. Moraes Rego, 1235, Recife, PE, 50670-901, Brazil.

<sup>11</sup>Department of Human Genetics and South Texas Diabetes and Obesity Institute, University of Texas Rio Grande Valley School of Medicine, Brownsville TX, USA.

<sup>12</sup>Movement Disorders Unit, Department of Neurology and Neurosurgery, Universidade Federal de São Paulo, São Paulo, Brazil.

<sup>13</sup>Instituto Nacional de Salud, Lima, Perú.

<sup>14</sup>Facultad de Ciencias de la Salud, Universidad Científica del Sur, Lima, Perú.

<sup>15</sup>Department of Biostatistics, University of North Carolina at Chapel Hill.

<sup>16</sup>CETRAM, Facultad de Ciencias Médicas, Universidad de Santiago de Chile, Santiago, Chile.

<sup>17</sup>Keck School of Medicine, University of Southern California, Los Angeles, CA, USA.

<sup>18</sup>Vitalant Research Institute, San Francisco, CA, USA.

<sup>19</sup>Laboratory of Genomic Diversity, Division of Cancer Epidemiology and Genetics, NCI, NIH, Rockville, MD, USA.

<sup>20</sup>Instituto de Medicina Tropical, University of São Paulo, São Paulo, Brazil.

<sup>21</sup>Neurology Institute, Universidad de la República, Montevideo, Uruguay.

<sup>22</sup>INBIOMEDIC Medicina Preventiva y Personalizada, Lima, Perú.

<sup>23</sup>Centro de Investigación en Medicina Traslacional. Universidad Privada Norbert Wiener. Lima, Perú.

<sup>24</sup>Faculdade de Medicina, Departamento de Medicina Social, Universidade Federal de Pelotas, Pelotas, RS, Brazil.

<sup>25</sup>Bloodworks Northwest Research Institute, Seattle, WA, USA.

<sup>26</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx NY USA 10461.

<sup>27</sup>Division of Public Health Sciences, Fred Hutchinson Cancer Center, Seattle WA USA, 98109.

<sup>28</sup>UCSF Benioff Children's Hospital Oakland.

<sup>29</sup>Icahn School of Medicine at Mount Sinai, New York, NY, USA.

<sup>30</sup>Department of Neurology, University of Maryland School of Medicine, Baltimore, MD 21201, USA.

<sup>31</sup>Department of Medicine, University of Washington, Seattle, WA.

- <sup>32</sup>Instituto de Pesquisas René Rachou, Fundação Oswaldo Cruz, Belo Horizonte, MG, Brazil.
- <sup>33</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
- <sup>34</sup>Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA.
- <sup>35</sup>Division of Genetics, Genomics, and Precision Medicine, University of Arizona, Tucson, AZ, USA.
- <sup>36</sup>Department of Medicine, University of Maryland School of Medicine, Baltimore, MD, USA.
- <sup>37</sup>Geriatrics Research and Education Clinical Center, Baltimore Veterans Administration Medical Center, Baltimore, MD, USA.
- <sup>38</sup>Human Genome and Stem Cell Research Center, University of São Paulo, São Paulo, SP, Brazil.
- <sup>39</sup>Department of Genome Sciences, University of Washington, Seattle, WA, USA.
- <sup>40</sup>Department of Epidemiology, Gillings School of Global Public Health, The University of North Carolina at Chapel Hill, Chapel Hill, NC, USA.
- <sup>41</sup>Genetics Department, Facultad de Medicina, Universidad de la República, Montevideo, Uruguay.
- <sup>42</sup>Department of Epidemiology, University of Washington, Seattle, WA, USA.
- <sup>43</sup>Departamento de Neurologia, Universidade Federal de Ciências da Saúde de Porto Alegre, Porto Alegre, Brazil.
- <sup>44</sup>Department of Cardiology, University of Groningen, University Medical Center Groningen, Groningen, The Netherlands.
- <sup>45</sup>The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA USA.
- <sup>46</sup>Department of Neurology, Miller School of Medicine, University of Miami, and The Evelyn F. McKnight Brain Institute, FL, USA.
- <sup>47</sup>Instituto de Ciências da Saúde, Universidade Federal do Pará, Belém, Brazil.
- <sup>48</sup>Departamento de Farmacologia, Universidade Federal do Rio Grande do Sul, Porto Alegre, Brazil.
- <sup>49</sup>Serviço de Neurologia, Hospital de Clínicas de Porto Alegre, Porto Alegre, Brazil.
- <sup>50</sup>Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital, Boston, MA, USA.
- <sup>51</sup>Division of Sleep and Circadian Disorders, Brigham and Women's Hospital, Harvard Medical School, Boston, MA USA.
- <sup>52</sup>Ribeirão Preto Medical School, Universidade de São Paulo, Ribeirão Preto, Brazil.
- <sup>53</sup>Institute for Human Genetics, University of California, San Francisco, San Francisco, CA, USA.
- <sup>54</sup>Bakar Computational Health Sciences Institute, University of California, San Francisco, San Francisco, CA, USA.
- <sup>55</sup>Quantitative Biosciences Institute, University of California, San Francisco, San Francisco, CA, USA.
- <sup>56</sup>Department of Bioengineering and Therapeutic Sciences, University of California San Francisco, San Francisco, CA 94143.
- <sup>57</sup>Facultad de Salud Pública y Administración. Universidad Peruana Cayetano Heredia, Lima, Peru.
- <sup>58</sup>Program in Health Equity and Population Health, University of Maryland School of Medicine, Baltimore, MD 21201, USA.
- <sup>59</sup>Program in Personalized Genomic Medicine, University of Maryland School of Medicine, Baltimore, MD 21201, USA.
- <sup>§</sup>Deceased December 24, 2021.

\*Corresponding authors: E-mail: vicbp1@gmail.com and timothydoconnor@gmail.com

# 1 Abstract

2 Latin America is underrepresented in genetic studies, which can exacerbate disparities in personalized genomic  
3 medicine. However, genetic data of thousands of Latin Americans are already publicly available, but require a  
4 bureaucratic maze to navigate all the data access and consenting issues. We present the Genetics of Latin  
5 American Diversity (GLAD) Project, a platform that compiles genome-wide information of 54,077 Latin  
6 Americans from 39 studies representing 45 geographical regions. Through GLAD, we identified heterogeneous  
7 ancestry composition and recent gene-flow across the Americas. Also, we developed a simulated-annealing-  
8 based algorithm to match the genetic background of external samples to our database and share summary  
9 statistics without transferring individual-level data. Finally, we demonstrate the potential of GLAD as a critical  
10 resource for evaluating statistical genetic softwares in the presence of admixture. By making this resource  
11 available, we promote genomic research in Latin Americans and contribute to the promises of personalized  
12 medicine to more people.

## 13 Introduction

14

15 Latin Americans/Latinos/Latinx/Latine, or Hispanics, as an ethnic label, represent a set of populations across the  
16 Americas characterized by admixture between populations from many parts of the world with distinct ancestry  
17 compositions<sup>1</sup>. As such, treating Latin Americans as a single group is an over-simplification that may limit  
18 opportunities to improve health and clinical treatment. Latin Americans comprise 656 million people (8.5% of the  
19 world's population)<sup>2</sup>. In the United States, Latin Americans represent 18% of the population and are the fastest-  
20 growing demographic<sup>3</sup>. Unfortunately, these populations remain understudied and underserved in biomedical  
21 research and are at risk of being left behind by the precision medicine revolution. For example, Latin Americans  
22 only represent about 0.23% of participants in genome-wide association studies (GWAS) performed<sup>4</sup>. Several  
23 important efforts have been made to understand Latin American (**LAm**) genetic history and to identify genetic  
24 variants associated with complex traits<sup>5-26</sup>. However, most of these samples are thinly spread across many  
25 projects with few initiatives (e.g., the Mexico City Prospective Study<sup>27</sup>) to obtain the 100K+ individuals necessary  
26 to have statistical power comparable to other population groups (e.g., Europeans<sup>28</sup> and East Asians<sup>29,30</sup>).

27 To remedy the under-representation of Latin Americans in genomic studies, we have created the Genetics of  
28 Latin American Diversity database (GLADdb), a resource to infer fine-scale patterns of population structure  
29 across the Americas and boost statistical power for the discovery of genetic factors contributing to LAm health  
30 and disease. By gleaning LAm individuals through dbGaP and whole genome sequencing projects across the  
31 Americas, we gathered over 54,000 unrelated individuals, either genotyped and imputed, or sequenced, from  
32 ten countries (**Figure 1A**) spanning 45 geographical groups (**Table S1** and **Table S2**). These group labels  
33 reflect administrative division level (e.g., country, state, or city level information) when available. Using GLADdb,  
34 we addressed two major goals regarding LAm genomics: (i) in population genetics: to identify recent fine-scale  
35 patterns of distant relatedness and differentiation along the Americas, providing insights into regions with  
36 genetic underrepresentation, and (ii) in genetic epidemiology at two levels: (a) by developing a web tool for  
37 matching the genetic background of GLADdb individuals with external pools of samples providing additional  
38 power to discover genotype-phenotype associations and (b) by demonstrating how GLADdb can be utilized for  
39 testing statistical genetic software in diverse LAm cohorts.

40 We start by exploring distant genetic relatedness among LAm countries. Several studies have focused on  
41 determining the sources and timing for admixture events that led to the current genetic composition in some  
42 LAm countries<sup>6,11-13,16,31-33</sup>. However, understanding LAm genetic diversity goes beyond the initial continental  
43 admixture and involves bottlenecks, founder effects, and migration into and along the Americas, especially as it  
44 relates to fine-scale population structure within the continental sources (i.e., Indigenous American, European  
45 and African groups). We explored population structure and recent migration among LAm regions by analyzing  
46 allelic frequencies and identity-by-descent (IBD) sharing.

47 We then address issues about data availability when performing large-scale analyses in LAm populations. Many  
48 association analyses in LAm populations have smaller sample sizes than similar studies in Europeans and other  
49 populations. Data, even when publicly available, is often prohibitively restrictive for investigators to access  
50 because of quality control efforts and data curation, in addition to the bureaucratic maze typically required to  
51 obtain the data<sup>34</sup>. Artomov et al.<sup>35</sup> showed that with a large control cohort, a matching procedure, which is the  
52 identification of individuals with similar genetic backgrounds with external data, and sharing of their summary  
53 statistics (e.g., allele counts), is possible without the transfer of individual-level data. The matching procedure is  
54 designed to guard against genetic control inflation and reduce spurious associations due to population structure.  
55 Given Artomov's approach was developed on European ancestry individuals, we adapt this idea to the complex  
56 ancestral composition of LAm individuals. We devised an enhanced matching algorithm to explore the principal  
57 component space derived from our diverse GLADdb cohorts, into which we project external samples and match  
58 them to GLADdb individuals using ancestral background summary statistics. From the selected GLADdb  
59 individuals, we will generate and return summary statistics of genome-wide genotype frequencies and aggregate  
60 local ancestry composition to increase the sample size and power of the end-user study. Since GLADdb  
61 consists of both cases and controls for different phenotypes, we will also use phenotype filters to select

62 individuals useful as controls. We implemented all these features through an interactive web portal  
63 ([glad.igs.umaryland.edu](http://glad.igs.umaryland.edu)).

64 Finally, we demonstrate the potential of GLADdb as a critical resource for evaluating the performance of  
65 statistical genetic software in the presence of admixture. We do so by comparing three polygenic risk score  
66 (PRS) algorithms for estimating PRS in admixed individuals in a scenario where the ancestries corresponding to  
67 the GWAS summary statistics do not match the target cohort. PRS, the linear summation of risk variants  
68 weighted by their GWAS effect size, are highly impacted by the European-ancestry bias underlying much of the  
69 available GWAS data, and their transferability across populations remains a critical limitation of the approach  
70 <sup>36,37</sup>. GLADdb is uniquely situated to support methods development efforts that help ensure cross-population  
71 transferability of statistical genetic applications.

## 72 Results

### 73 Data Description and QC

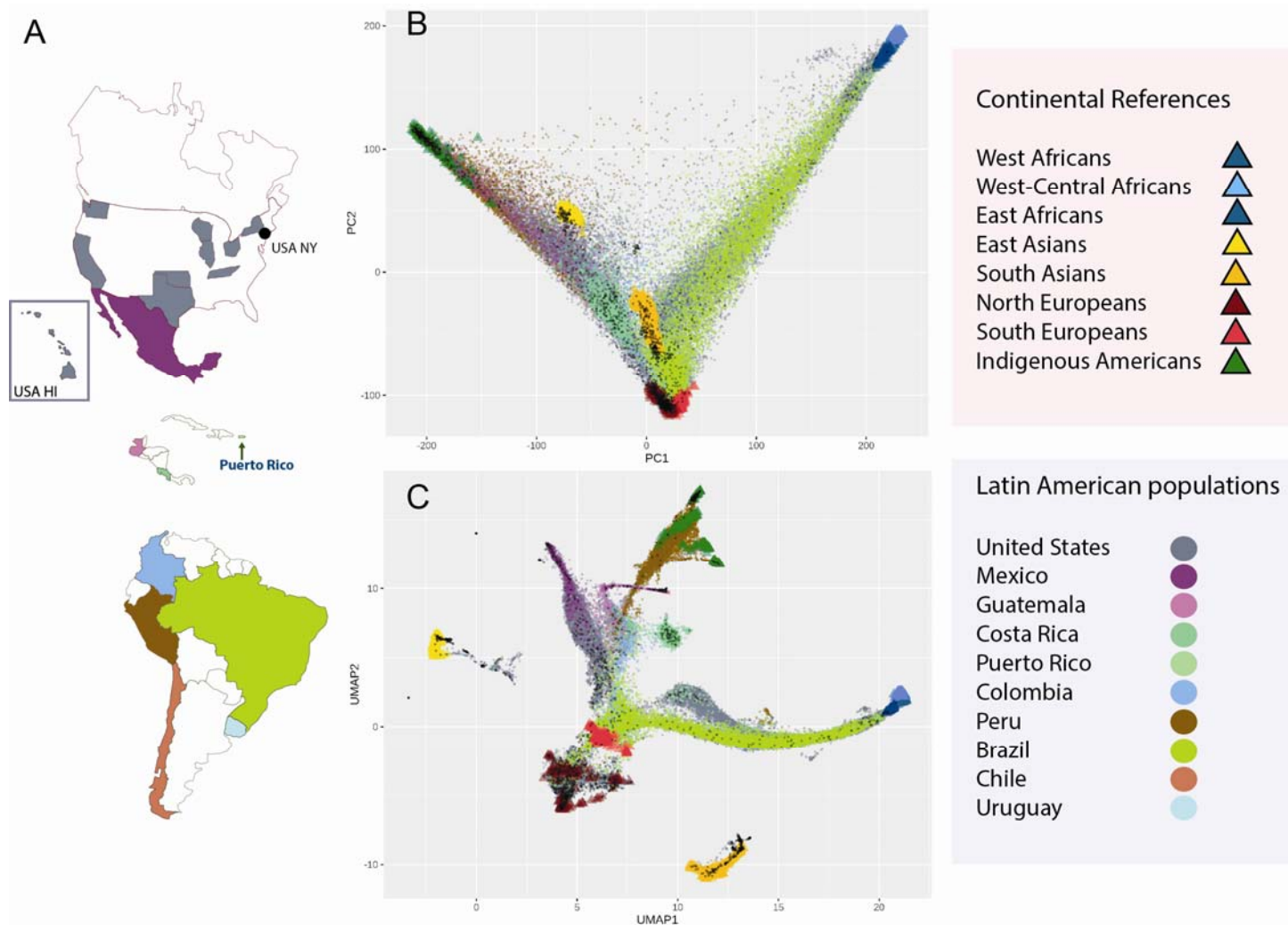
74 Our main workflow is described in [Figure S1](#) and [Supplementary methods](#). Briefly, we have explored over  
75 268K samples by gathering data from 39 dbGaP cohorts and other WGS projects that include US Hispanics /  
76 LAm individuals <sup>5-8,13,38</sup> ([Table S1](#)). As inclusion criteria, we gathered individuals self-described as “Latino” or  
77 “Hispanic” and ADMIXTURE-defined individuals. This latter criterion was applied to identify possible LAm  
78 individuals using ADMIXTURE analysis<sup>39</sup>, keeping any individuals with more than 2% Indigenous American (IA)  
79 ancestry (See Methods). For genotyped cohorts ([Table S1](#)), we imputed all self-described (GLAD-SD,  
80 n=25,627) and ADMIXTURE-defined (GLAD-AD, n=17,642) individuals within each cohort using the TOPMed  
81 Imputation server<sup>40</sup>. After imputation QC, we kept 42,539 individuals that were combined into a single dataset  
82 with sequencing data from TOPMed Project<sup>5</sup> (27,088 individuals) and 1000 Genomes Project<sup>38</sup> (345 individuals)  
83 with 9,121,629 overlapping variants with an imputation  $r^2 > 0.3$  across all datasets. For all analyses here, we  
84 kept overlapping variants with imputation  $r^2 > 0.9$  in each dataset before merging. The final merged dataset with  
85  $r^2 > 0.9$  for analysis contains 3,248,494 biallelic variants. Finally, to remove the family structure in GLADdb, we  
86 inferred kinship coefficients using IBD segments on the complete dataset, keeping 54,077 unrelated individuals  
87 (See Methods).

### 88 Continental Population Structure of GLADdb

89 Using 54K unrelated samples and ancestry-reference groups ([Table S3](#)), we explored the patterns of diversity  
90 and differentiation throughout the Americas using principal component analysis (PCA), uniform manifold  
91 approximation and projection (UMAP), and ADMIXTURE analyses ([Figure 1B and C](#), [Figure S2-S5](#)). Both  
92 results highlighted some important points. First, the samples cluster according to ancestry and not technology or  
93 other batch effects ([Figures S3 and S4](#)). Notably, GLAD-AD individuals cluster well with other GLAD-SD  
94 individuals validating our inclusion criteria ([Figure S4A and B](#)). By coupling UMAP and ADMIXTURE results, we  
95 reaffirm the heterogeneous ancestry distribution of LAm individuals, with some groups showing predominantly IA  
96 ancestry (Peru, Mexico, and Guatemala) and others showing majority admixture between European and African  
97 ancestries (USA and Brazil) ([Figure 1C](#) and [Figure S5](#)). Regarding sample sizes, the best-represented regions  
98 in GLADdb included Brazil, Central America, Mexico, Peru, and the United States.

99





**Figure 1. Dimensionality reduction of genetic data for more than 52K unrelated Latin Americans from the GLAD database.** A) Geographical distribution of GLADdb cohorts. B) Principal Component Analysis of the entire dataset based on high-quality imputed SNPs ( $r^2 > 0.9$ ) showing the sampling spread of Latin Americans. C) Uniform Manifold Approximation and Projection (UMAP) of the first 10 PCs showing clusters of different population groups.

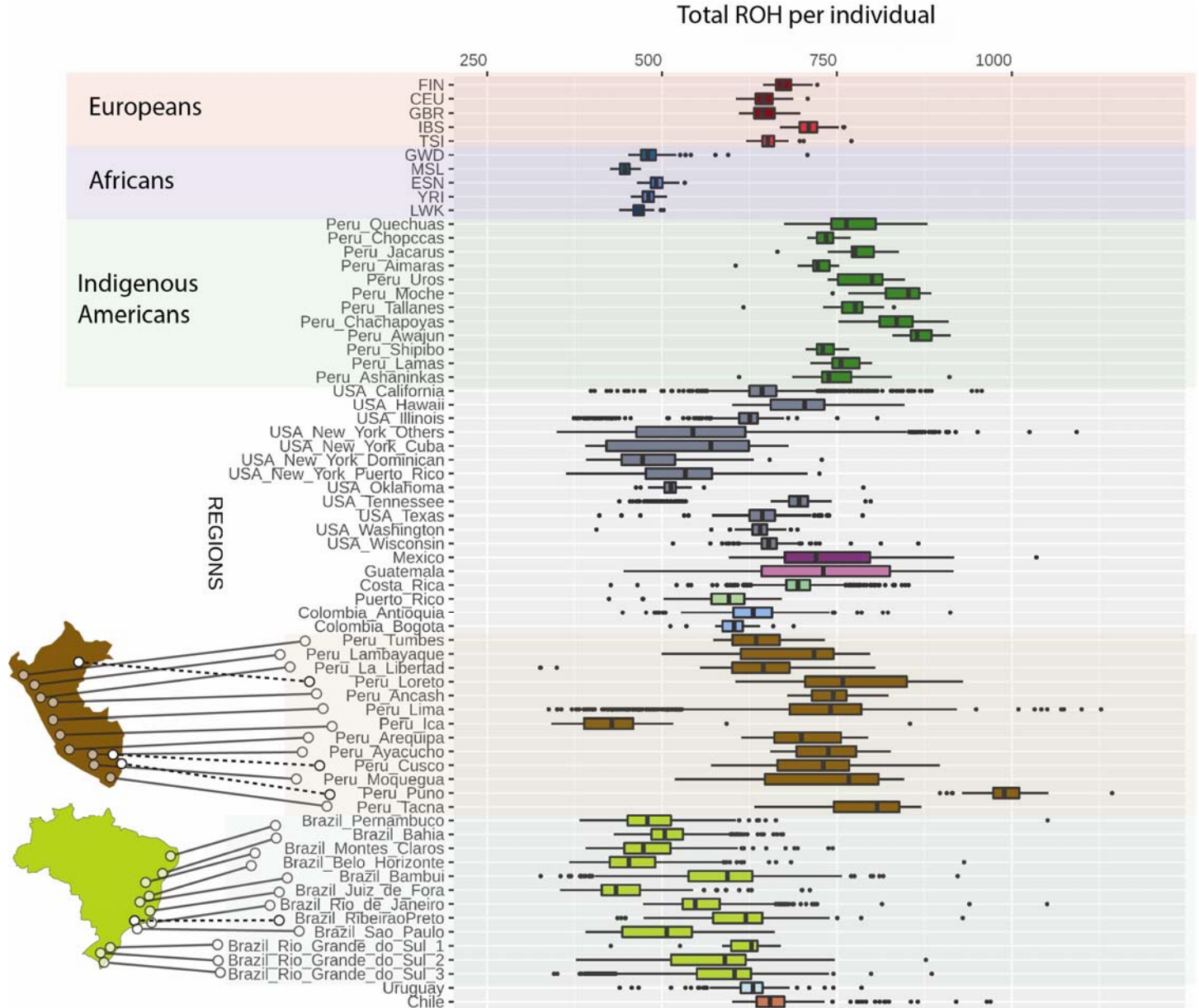
## Levels of genetic diversity within Latin American groups

Although our population structure analyses identified a wide diversity of LAm groups, these groups originated from continental progenitors that suffered a significant drop in effective population size during the colonial period of the Americas<sup>41–43</sup>. This resulted in a higher level of consanguinity and enrichment of long runs of homozygosity observed in some LAm groups (e.g., CLM and PEL from 1000 Genomes Project) compared to Finnish<sup>42</sup>, a population notably shaped by a strong founder effect. Based on demographic information available for the cohorts, we organized GLAD-SD individuals into 45 self-described LAm groups, consistent with geographic labels based on administrative division level (e.g., country, state, or city level information) (Table S2). In addition, we included 12 IA populations from the Peruvian Genome project as well as 5 European (EUR) and 5 African (AFR) populations from the 1000 Genomes Project (See Methods).

We explored the levels of diversity in each group by inferring runs of homozygosity<sup>44</sup> (ROH) (Figure 2, See Methods). As expected, individuals from Africa showed lower values for total ROH compared to individuals from Europe and Indigenous groups from Peru. Analogously, LAm groups with higher proportions of African ancestry (e.g., Peru-Ica and Northeast Brazilian regions) tend to have the lowest total ROH. Furthermore, taking advantage of the detailed sample representation for 13 Peruvian and 12 Brazilian regions, we determined the correlation between average genome-wide ancestry proportions (Table S2) and the median total ROH for each

122  
123  
124  
125

population. We observed a positive correlation between the average Indigenous American ( $r=0.81$ ,  $p\text{-value} = 0.00246$ ) and European ( $r=0.88$ ,  $p\text{-value} = 1.12 \times 10^{-4}$ ) ancestries with a higher density of ROH in Peruvians and Brazilians, respectively. Interestingly, both correlations follow a North to South line.



126  
127  
128  
129  
130  
131

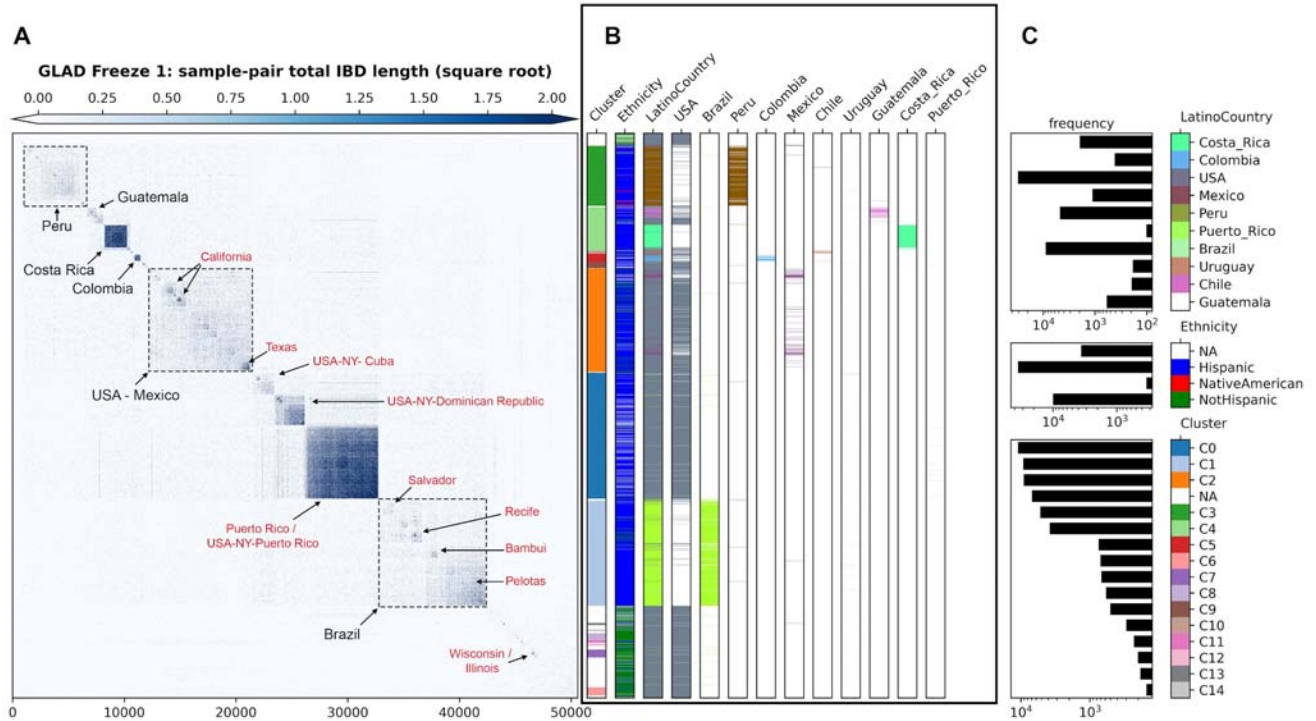
**Figure 2. Distribution of Genome-wide amount of Runs of Homozygosity for Latin American groups and Reference populations included in GLADdb.** The upper part of the plot shows continental reference populations—the lower part details the distribution in Peru and Brazil. Populations are sorted in a North-to-South pattern.

132  
133  
134  
135  
136  
137  
138  
139  
140

## Fine-scale population structure revealed by IBD network

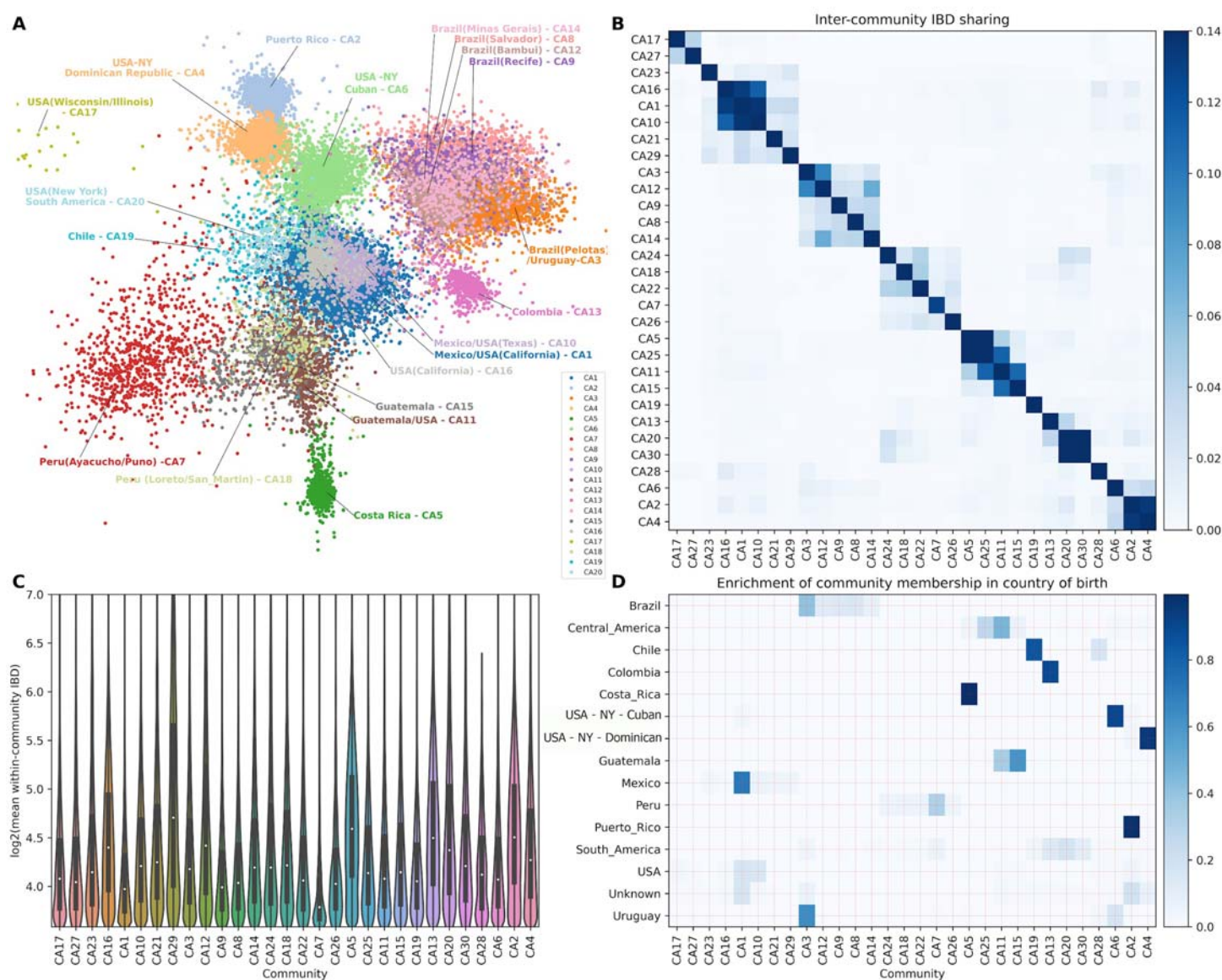
To obtain a fine-scale picture of population structure among LAm groups, we built a sample-pair genome-wide total IBD matrix using all IBD segments  $> 5\text{cM}$  shared in our 54K dataset. Clusters in this matrix are mainly consistent with geographic labels, with strong intra-cluster sharing among individuals from Puerto Rico, Dominican Republic, and Costa Rica (Figure 3). Given the sample size and genetic diversity, finer-scale population structure is observable in clusters representing the USA/Mexico, Peru, and Brazil. To reveal the substructure, we employed an IBD network-based community detection algorithm to further analyze relatedness patterns. We selected the top 20 IBD-network-based communities that accumulated 69% of GLADdb individuals (other communities each have less than 270 individuals). Each of these communities (labeled as CA1 to 20 and

141 ordered from largest to smallest) showed enrichment of individuals from a particular country, such as Costa Rica  
 142 (99.6%, IBD community CA5), Puerto Rico (98%, IBD community CA2), Dominican Republic (95.0%, CA4),  
 143 Cuba (89.8%, CA6), Colombia (89.4%, CA13), and Chile (84%, CA19) (Figure 4). In contrast, individuals from  
 144 Mexico, Peru, and Brazil were grouped in several communities (Mexico: 7, Brazil: 5, Peru: 13 communities).  
 145 These communities were represented by individuals from a particular region, reflecting the extensive sampling  
 146 performed in these countries (Figure 4).  
 147



148  
 149 **Figure 3. Clustering of total IBD matrix of unrelated individuals from GLADdb.** A) Heatmap of the square root of sample-pair total  
 150 IBD. Annotations within the heatmap represent the most enriched demographic labels in the indicated blocks. Labels with “USA-NY-  
 151 country” correspond to self-described US-Hispanic living in New York with a specific country of origin. B) Individual-level labels of  
 152 agglomerative cluster assignment (1st column), ethnicity (2nd), and sampling Country (combined:3rd and separated: 4th onward). C)  
 153 Frequency of labels (log scale) and color keys.  
 154





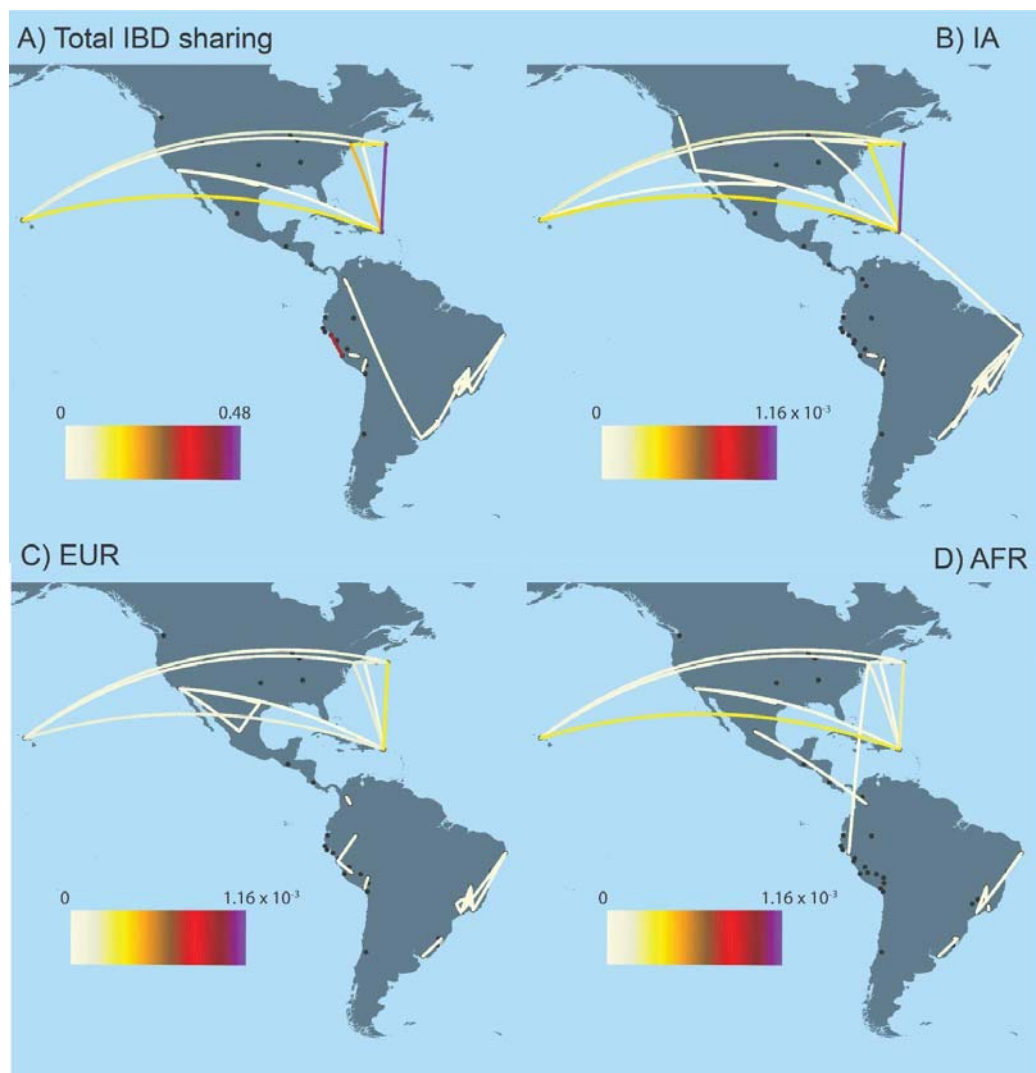
**Figure 4. IBD network community detection.** We infer the community structure using the infomap algorithm based on a matrix of IBD segments greater than 5cM. A) Top 20 IBD network communities. For visualization purposes, only individuals with connections > 30 are included in the layout calculation. The community labels, such as CA1 and CA2, are named according to the IBD version used and the rank of the community sizes, with CA1 representing the largest community when using all IBD segments. For communities inferred from short and long IBD segments, the corresponding labels are CS1 (Figure S6A) and CL1 (Figure S6B), respectively. B) IBD sharing among the top 30 inferred communities (ordered by agglomerative clustering; the same order was followed in C and D). C) Distribution of IBD shared among individuals in each community. D) Enrichment of IBD community membership in the country of origin (i.e., proportions of community labels for individuals born in a given country). To visualize the dynamics before and after the Spanish colonization of the Americas, two different IBD networks were built based on IBD segments between 5-9.3cM (Figure S6A) and those > 9.3cM (Figure S6B), respectively, which revealed distinct patterns of detected communities.

## Long-distance relatedness among Latin American groups

To explore recent migration among 45 LAm regions, we restricted our analyses to IBD segments greater than 21 cM, representing a recent common ancestor in the last seven generations corresponding to post-colonial times<sup>45</sup> and after the admixture process. We reasoned that sharing of larger IBD segments could be originated predominantly from gene flow among regions. At the inter-regional level, we detected higher levels of sharing between Puerto Rico with New York (Specifically with Puerto Ricans in New York) and Hawaii groups. Another tight sub-network of sharing is observed in Brazil (Figure 5A), where the South East region (São Paulo and

173 Minas Gerais states) have major connections with other Brazilian populations. Interestingly, there are IBD-  
174 sharing connections between Uruguay, South Brazil, and Colombia. On the Pacific side, two Peruvian regions  
175 (Ica and Trujillo) show high values of IBD sharing.

176 Considering the multi-way admixed origin of LAM populations, we devised a statistic (ancestry-specific IBD  
177 score) that quantifies the level of relatedness among two admixed populations for a particular ancestry (AFR,  
178 EUR, or IA) (Figure 5). We computed the ancestry-specific IBD score (asIBD score, see Methods) by coupling  
179 the IBD and local ancestry inferences. Our asIBD score explains the relationship of ancestry-specific IBD  
180 segments with respect to the global ancestry of the populations. We detected a different ancestry-sharing  
181 pattern between Puerto Rico with New York, and Hawaii (Table S4A). A three-way sharing with predominant IA  
182 ancestry characterized the sharing among Puerto Rico and New York. On the other hand, the Puerto Rico and  
183 Hawaii sharing is characterized by predominant IA and AFR-related ancestries (Figure 5). The sharing cluster in  
184 Brazil has higher values of asIBD for the IA ancestry, indicating a more homogenous composition of IA ancestry  
185 in those regions (Table S4B). For IBD sharing between Peru-Ica and Peru-La-Libertad, the EUR ancestry  
186 showed the highest value for the asIBD (Table S4A).  
187



188

189 **Figure 5. Identical-By-Descent (IBD) analyses of Latin American groups.** We explored the relationship among LAM regions by  
190 inferring the average IBD shared among regions (A) and an ancestry Specific IBD Score (asIBDScore) for Indigenous American (B),  
191 European (C), and African ancestries (D). Dots represent Latin American regions. For African and European Ancestries, we remove the  
192 sharing between Peru-Ica and Peru-La Libertad due to their higher sharing and to improve visualization. Plot A range showed the  
193 average amount of cM shared among two individuals from populations 1 and 2. On the other hand, the plot range for B-D represents the  
194 same statistic focused on segments of a specific ancestry and controlled by global ancestry proportions in each population.

## Supporting external studies through the GLADdb matching algorithm and statistical genetic software benchmarking

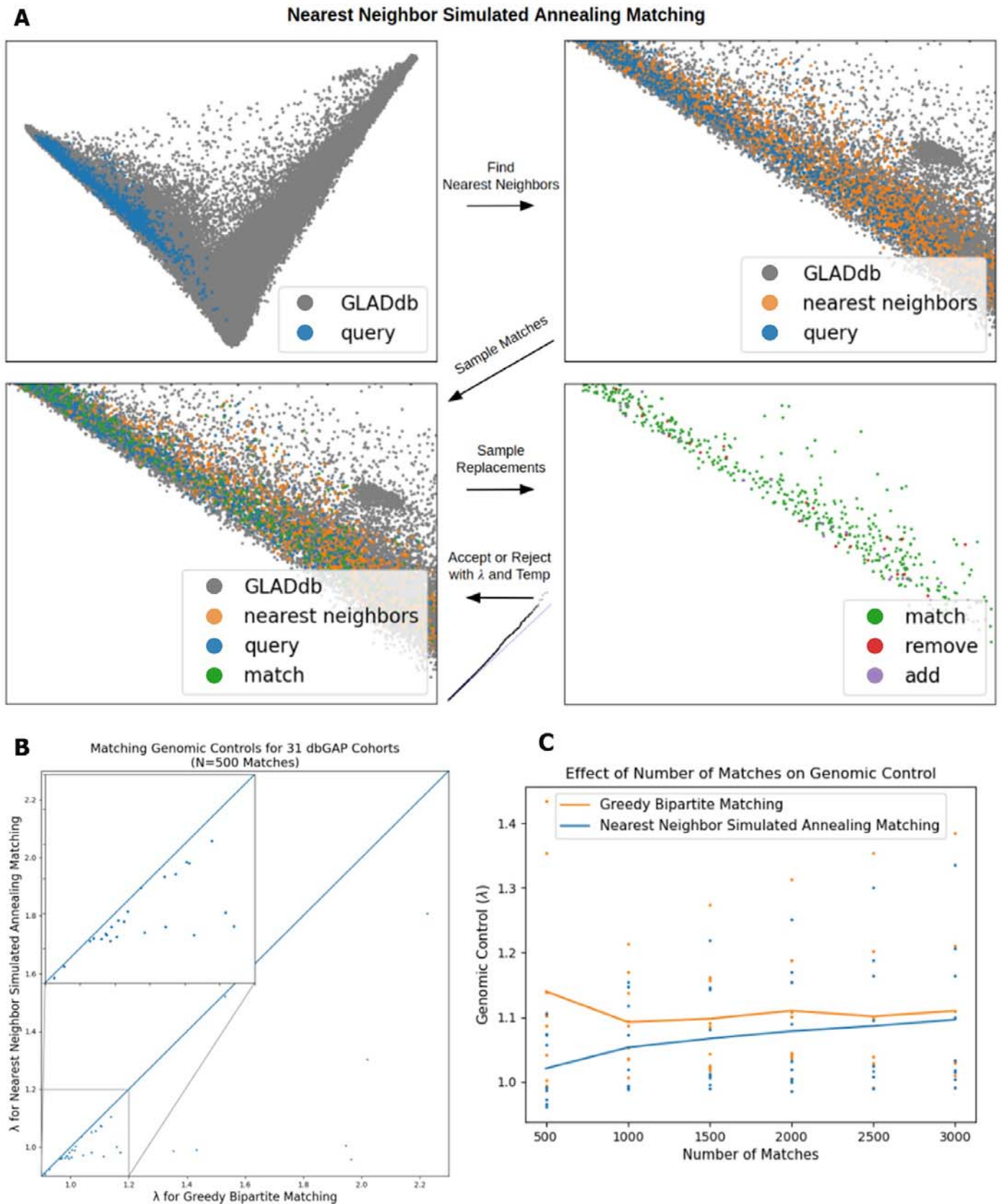
One of GLADdb's ultimate goals is to provide controls for GWAS and admixture mapping studies. We addressed this goal by developing a genetic matching algorithm. Our method, nearest neighbor simulated annealing matching, shown in **Figure 6A** and outlined in Methods, employs local search to find the optimal cohort from a set of candidates. The algorithm operates on a principal component space in which the external-user-provided query cases can be used to search for controls without needing individual genotypes. The algorithm computes variance-weighted Minkowski distance pairwise between query cases and potential controls, selects the nearest neighbors as candidate controls, samples a set of matches from the candidates, and iteratively resamples and refines the set of matches using simulated annealing, optimizing for the genomic control statistics  $\lambda$ <sup>46,47</sup>.

To evaluate both our matching algorithm and the extent to which GLAD cohorts can provide valid control sets, we performed the following experiment. Using 1000 Genomes populations and some GLAD cohorts as cases, in which the pseudo-phenotype belongs to the query cohort, we ran a greedy bipartite matching baseline<sup>48,49</sup> and our matching algorithm and returned summary statistics (i.e., alternative allele frequency, genotype counts, and haplotype ancestry counts by segment) for various control set sizes. Then, for each pair of cases and controls, we ran a GWAS for which the genomic control  $\lambda$  statistics are reported in **Table 1** and more extensively in **Figure 6B**. For the analyzed cohorts, which represent a variety of admixed groups, the matched controls yield genomic controls close to 1, suggesting that GLAD will be able to provide useful controls for a variety of cohorts, and our matching algorithm shows slight improvements for larger and more varied query cohorts. These improvements narrow progressively as the number of matches required increases (**Figure 6C**).

**Table 1.** Comparison of genomic control results ( $\lambda$  statistics) when returning 500 control individuals from GLAD using the Greedy Bipartite Matching and the Simulated Annealing Nearest Neighbor Matching algorithm.

Source		1000 Genomes populations				GLAD Cohorts			
Population or Cohort		MXL	CLM	PEL	PUR	HCHS SOL	MESA	SIGMA	LARGE-PD
Case N		64	94	85	104	6646	1046	1146	1463
Greedy Bipartite Matching	Genomic control	0.9438 ± 0.0009	1.0951 ± 0.0018	0.9555 ± 0.0012	0.9928 ± 0.0017	0.9939 ± 0.0063	1.0080 ± 0.0136	1.0721 ± 0.0074	1.0268 ± 0.0177
Nearest Neighbor Simulated Annealing Matching Algorithm		0.9400 ± 0.0001	1.0905 ± 0.0009	0.9496 ± 0.0010	0.9881 ± 0.0003	0.9612 ± 0.0045	0.9820 ± 0.0047	1.0480 ± 0.0072	1.0045 ± 0.0050





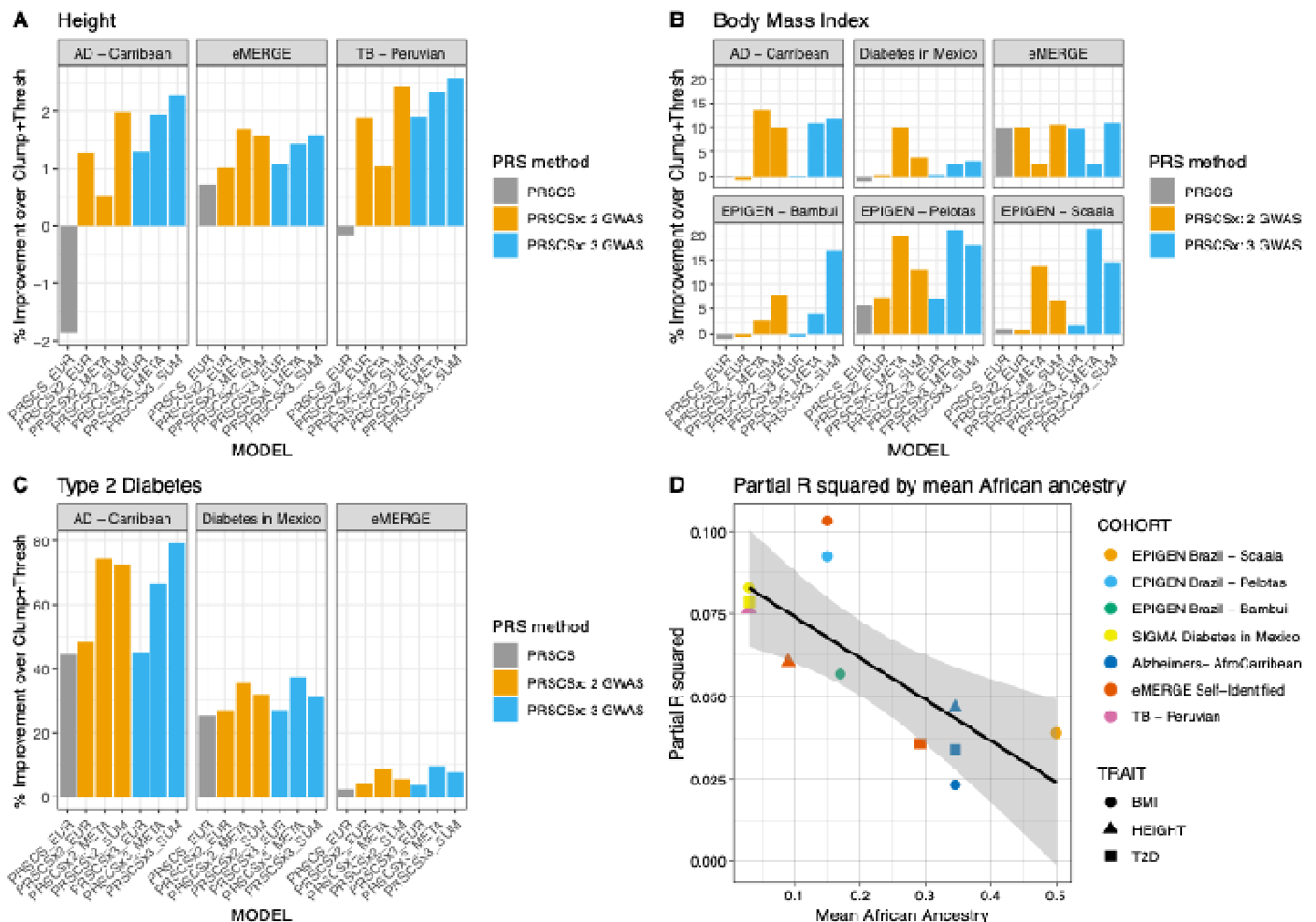
**Figure 6. Nearest Neighbor Simulated Annealing Matching Algorithm and Results.** A) Visual overview of the algorithm. B) Comparison with baseline bipartite matching algorithm (x-axis), where points below the line  $y=x$  indicate our algorithm outperforming the baseline (small box highlights high density region). C) Effect of number of matches on improvement over the baseline.



227

228

229 In addition to control matching, GLADdb is an optimal resource for benchmarking statistical genetic software in  
 230 complex, heterogeneous cohorts with a wide range of available traits. We demonstrated this potential by  
 231 comparing several popular PRS algorithms (Clumping + Thresholding using PRSice-2<sup>50</sup>, PRS-CS<sup>51</sup>, and PRS-  
 232 CSx<sup>52</sup>) using a subset of GLAD-SD (**Table S5**, see Methods) with type 2 diabetes (T2D) status, height, or BMI  
 233 data under a hypothetical scenario where LAm GWAS data is not available (**Table S6**). The GLAD-SD subset  
 234 includes LAm cohorts with very different population histories and ancestry proportions (e.g., Afro-Caribbeans,  
 235 Brazilians, and Peruvians). Though the use of the Bayesian PRS-CS method, in general, outperformed PRSice-  
 236 2, the inclusion of non-European GWAS data using PRS-CSx yielded the largest increase in PRS predictive  
 237 performance (Figure 7A-C, Figure S7). PRS-CSx improved single-ancestry PRS predictive performance (e.g.,  
 238 East Asian PRS from PRS-CSx versus PRS-CS or PRSice-2) in nearly every instance (**Table S7**). Combining  
 239 the posterior effect sizes estimated by PRS-CSx further improved models (Figure 7A-C, **Table S7**). Note that the  
 240 best approach for combining PRS information varied by cohort, likely reflecting cohort heterogeneity (Figure S8).  
 241 Model performance, as measured by partial  $R^2$ , was negatively associated with mean African ancestry (-0.02 per  
 242 standard deviation African ancestry, p-value 0.005, Figure 7D). While the percent improvement achieved when  
 243 leveraging non-European GWAS data can be as high as 80% over the clumping + thresholding model, the  $R^2$  of  
 244 each PRS still can be modest. For example, in the Alzheimer's cohort from the Caribbean, the T2D PRS-CSx  
 245 model improved prediction by nearly 80%, but the  $R^2$  of that model was only 0.03 on the observed scale (Figure  
 246 7D).  
 247



248

249

250 **Figure 7. PRS in select cohorts from GLAD-SD.** A) Comparison of Height model performance as percent improvement over a  
 251 European-ancestry GWAS Clumping + Thresholding PRS. Models include PRS-CS using European-ancestry GWAS, PRS-CSx using  
 252 European and East Asian-ancestry GWAS, and PRS-CSx using European, East Asian, and African-ancestry GWAS. All models were  
 compared using the correlation between the prediction and the trait. B) Comparison of BMI model performance. C) Comparison of T2D

253 model performance. D) Total R<sup>2</sup> of best PRS model by African ancestry. Cohorts are labeled by color, traits are labeled by shape. Partial  
254 R<sup>2</sup> was calculated by squaring Pearson's r followed by subtracting the full model (PRS + covariates) from the base model (covariates  
255 only, see methods). African ancestry proportions were estimated using ADMIXTURE.

## 256 Discussion

257 Latin American individuals are not well represented in genomic and epidemiological studies. This means we  
258 have a poor knowledge of their genetic diversity and environmental backgrounds, which limits the applicability of  
259 personalized medicine and our understanding of the basis of complex phenotypes<sup>53</sup>. GLADdb aims to tackle the  
260 underrepresentation of genomic data by gathering genome-wide data of LAm populations into a single resource.  
261 Through GLADdb, we have two main contributions to LAm genomics: 1) **Population genetics**: we elucidated  
262 population structure and gene flow across LAm regions. 2) **Genetic epidemiology**: we developed an algorithm  
263 and an online portal (see [Supplementary information 1](#)) to provide summary statistics from control individuals  
264 from GLADdb with a similar genetic makeup to external samples. Also, by assembling a collection of LAm  
265 cohorts with very different population histories, we have created a unique tool for evaluating the performance of  
266 statistical genetic software in the presence of admixture and other complexities.

267  
268 For population genetics, continental migrations were the initial sources of LAm diversity. However, other  
269 processes have shaped this diversity and the relationships across geographic regions. Through ROH and IBD  
270 inferences, we have explored the relationships at intra- and inter-population levels in Latin America in terms of  
271 diversity and relatedness. From both analyses, we observed that Peruvians, even with a higher level of  
272 homozygosity, have differentiated groups associated with geographical regions<sup>6,7</sup>. Moreover, IBD sharing tells  
273 us more about recent migrations when we restrict the analysis to 21 cM or greater, an interval size correlated  
274 with post-colonial events corresponding to the last seven generations before the present. We detected two main  
275 networks of sharing: Puerto Rico - New York and Hawaii, and the Brazilian internal sharing groups. In Latin  
276 America, during the 20th century, migrations have followed a rural-to-urban or outside-the-country tendency due  
277 to regional socioeconomic disparities<sup>54</sup>. Particularly, in Puerto Rico, during the early 1900s, a migration policy  
278 was enacted in response to its social and economic problems<sup>55</sup>. Hawaii, Dominican Republic, and Cuba were  
279 the primary destinations during the first stage of the Puerto Rican diaspora, followed by a strong migration to  
280 New York during the late 1940s<sup>56</sup>. It is noteworthy that there were socioeconomic differences between the  
281 groups participating in each migration stage<sup>57,58</sup>. For example, many individuals who migrated from Puerto Rico  
282 to Hawaii were recognized as *jibaros*<sup>58</sup>, which are countryside people who farm the land in a traditional way.  
283 However, Puerto Ricans who migrated to New York represented a cross-section of economic and social classes  
284<sup>57</sup>. By inferring the ancestral background of IBD segments, we found that the Puerto Rico/Hawaii sharing is  
285 characterized by predominant AFR and IA sharing compared to the IA and EUR sharing between Puerto Rico  
286 and New York. These contrasting patterns may reflect the differential composition of the two stages of migration.  
287 Brazil is another example of recent migration due to economic factors. During the 1950s, South Eastern Brazil,  
288 represented by Rio de Janeiro, São Paulo, and Minas Gerais, experienced a huge economic growth that  
289 triggered a massive migration to these regions<sup>59</sup>. We observed huge connectivity among South Eastern regions  
290 (Rio de Janeiro and São Paulo), showing higher values for EUR sharing suggesting higher mobility of European  
291 components in Brazil. Moreover, EUR sharing was detected between Southern Brazilian regions and Uruguay.  
292 This could reflect their recent shared history as Uruguay was annexed to Brazil before its independence<sup>60</sup>, and  
293 its demographic composition included a significant proportion of Brazilians at that time<sup>61</sup>.

294 For genetic epidemiology, our genotype matching algorithm and subsequent provision of control summary  
295 statistics meet a real need in the research community. Groups exploring the genetic architecture of traits in Latin  
296 American cohorts can increase their sample sizes without further straining budgets. This will help facilitate the  
297 discovery of genetic risk factors in a historically underrepresented population, which could lead to the discovery  
298 of population-specific variation and reduce bias in GWAS data. While there are initiatives that significantly  
299 increase the representation of Latin American subjects in genomics, access to that data remains a concern. In  
300 some cases, navigating the bureaucratic maze presents a real barrier, while in other cases, the data is  
301 proprietary. By constructing the first version of GLADdb, we have already acquired and aggregated Latin  
302 American data from across 39 cohorts. In addition, our matching and data transfer processes only require  
303 summary statistics (genotype counts and principal components), thus reducing the exposure of sensitive data.

304 Also, by employing our matching algorithm, we can potentially provide a better set of controls than by simply  
305 applying for individual cohorts from dbGaP or other public repositories, nor using allele frequencies from  
306 heterogeneously sampled cohorts alone.

307  
308 In addition to supporting genetic studies through control matching, GLADdb presents a valuable resource for  
309 evaluating the performance of genetic epidemiology software for methods development and benchmarking.  
310 Such software needs to be evaluated in the presence of admixture in addition to the more homogeneous  
311 cohorts. This is particularly evident for PRS estimation, where the impact of long-standing biases in GWAS data  
312 is well documented<sup>36,37,62</sup>. In our test case, we evaluated three popular PRS algorithms: clumping + thresholding  
313 implemented in PRSice-2, PRS-CS, and PRS-CSx. We found that PRS-CSx, which can model multiple GWAS  
314 populations simultaneously, significantly improved predictive performance over single ancestry methods. This  
315 was true despite not using GWAS data from any Latin American cohorts for this example. Variability in model  
316 performance likely reflected population heterogeneity across the different cohorts, and model performance was  
317 negatively associated with mean African ancestry. The sample sizes of the African-ancestry GWAS cohorts  
318 used for this study were smaller by an order of magnitude than the East Asian and European Ancestry GWAS  
319 cohorts. It is clear that well-powered, diverse GWAS is critical for equitable PRS performance. In the meantime,  
320 methodological innovation is required to improve cross-population portability for GWAS traits lacking adequate  
321 representation<sup>63</sup>. In addition to PRS-CSx, several methods such as LDPred-funct and Polypred include  
322 functional data, and TL-Multi utilizes transfer learning<sup>64-66</sup>. The robustness of existing and new PRS methods to  
323 admixture can be evaluated using the heterogeneous cohorts represented in GLADdb.

324  
325 A major challenge in our study, and for LAm genomics, is the poor representation of Indigenous American  
326 ancestries. Currently, the Indigenous American representation in public datasets is restricted to a few  
327 populations with higher levels of isolation which could lead to caveats in global and local ancestry inferences.  
328 This is important because several studies show that IA ancestry in an admixed LAm population closely relates to  
329 their local Indigenous groups<sup>6,11,16</sup>. To overcome the problems related to IA ancestry, we used a reference panel  
330 of Indigenous Peruvians and Guatemalans. These populations have higher effective population sizes compared  
331 to other native groups<sup>67</sup>, which is helpful for avoiding problems related to higher levels of genetic drift. In this  
332 way, we can get around the problem of IA inferences in Brazilians or USA individuals with some level of IA  
333 ancestry (i.e., Individuals with ancestry related to tribal nations in which genetic studies have not been allowed).  
334 Still, better ethically-aware representation in genomics is preferred. Furthermore, GLADdb allowed us to identify  
335 geographical regions better represented (e.g., Brazil, Mexico, and Peru) than others in sample size and  
336 genotyping technologies (WGS and array data). Moreover, even in these best-represented regions, there is an  
337 unbalance of ethnic diversity (e.g., European ancestry descendants are predominant in these datasets). This  
338 reality should motivate the need for urgently including regions like Bolivia or Paraguay as well as at the ethnicity  
339 level (i.e., African and Asian ancestries in the Americas).

340  
341 In conclusion, through GLADdb, we highlighted the heterogeneous ancestry composition across LAm  
342 populations and inferred ancestry differences in gene flow events relatedness among LAm regions. Also, by  
343 sharing summary statistics, we are contributing to improving global equity in genomic research, specifically in  
344 epidemiological research in which GWAS is performed routinely. This is one more step to ensuring that health  
345 disparities arising from genetic studies do not become pervasive in admixed and non-European populations.  
346

## 347 Acknowledgements

348  
349 We would like to thank Evangeline “Eevee” O’Connor for assistance in providing an acronym that is both  
350 accurate and contributes to generally uplifting our research. TDO was supported by National Human Genome  
351 Research Institute of the National Institutes of Health under Award Number R35HG010692. ETS was supported  
352 by FAPEMIG (Fundação de Amparo à Pesquisa do Estado de Minas Gerais) RED 00314-16; Programa  
353 Nacional de Genômica e Saúde de Precisão – Genomas Brasil from the Brazilian Ministry of Health (CNPq  
354 Process 403502/2020-9); Conselho Nacional de Desenvolvimento Científico e Tecnológico - CNPq. RDH is

355 supported by the National Institutes of Health under Award Number R01 GM142112. DPL was supported by the  
356 National Heart, Lung, And Blood Institute of the National Institutes of Health under Award Number T32  
357 HL007698-25. JNF was supported by Research Training in the Epidemiology of Aging funded by the National  
358 Institutes on Aging under Award Number T32 AG000262.

359  
360 Latin American Research Consortium on the Genetics of Parkinson Disease (LARGE-PD) is funded by the  
361 National Institutes of Health/National Institute of Neurologic Disorders and Stroke (NIH/NINDS) (R01  
362 NS112499).

363  
364 The NINDS-sponsored Stroke Genetics Network (SiGN) is funded by the National Institutes of Health/National  
365 Institute of Neurologic Disorders and Stroke (R01 NS105150 and R01 NS100178).

366  
367 Genome Sequencing for the Trans-Omics in Precision Medicine (TOPMed) program was supported by the  
368 National Heart, Lung and Blood Institute (NHLBI). Genome Sequencing for "NHLBI TOPMed: The Genetic  
369 Epidemiology of Asthma in Costa Rica" (phs000988.v3p1) was performed at Northwest Genomics Center  
370 (HHSN268201600032I / 3R37HL066289-13S1). Genome Sequencing for "NHLBI TOPMed: San Antonio Family  
371 Heart Study (SAFHS) " (phs001215.v4.p2) was performed at Illumina (3R01HL113323-03S1 and  
372 R01HL113322). Genome Sequencing for "NHLBI TOPMed: Women's Health Initiative (WHI)" (phs001237.v3.p1)  
373 was performed at Broad Institute Genomics Platform (HHSN268201500014C). Genome Sequencing for "NHLBI  
374 TOPMed: Hispanic Community Health Study - Study of Latinos (HCHS/SOL) " (phs001395.v2.p1) was  
375 performed at Baylor College of Medicine Human Genome Sequencing Center (HHSN268201600033I). Genome  
376 Sequencing for "NHLBI TOPMed: Multi-Ethnic Study of Atherosclerosis"(MESA) (phs001416.v2.p1) was  
377 performed at Broad Institute Genomics Platform (3U54HG003067-13S1). Genome Sequencing for "NHLBI  
378 TOPMed: Severe Asthma Research Program (SARP)" (phs001446.v2.p1) was performed at New York Genome  
379 Center Genomics (3U54HG003067-13S1). Genome Sequencing for "NHLBI TOPMed: Recipient Epidemiology  
380 and Donor Evaluation Study-III Brazil Sickle Cell Disease Cohort (REDS-BSCDC)" (phs001468.v3.p1) was  
381 performed at Baylor College of Medicine Human Genome Sequencing Center (HHSN268201600033I /  
382 HHSN268201500015C). Genome Sequencing for "NHLBI TOPMed: My Life Our Future (MLOF) Research  
383 Repository of Patients with Hemophilia A (Factor VIII Deficiency) or Hemophilia B (Factor IX Deficiency)"  
384 (phs001515.v2.p2) was performed at New York Genome Center Genomics (HHSN268201500016C). Genome  
385 Sequencing for "NHLBI TOPMed: Boston-Brazil Sickle Cell Disease (SCD) Cohort" (phs001599.v1.p1) was  
386 performed at Baylor College of Medicine Human Genome Sequencing Center (HHSN268201600033I,  
387 HHSN268201500015C, and HHSN268201600033). Genome Sequencing for "NHLBI TOPMed: Children's  
388 Health Study (CHS) Integrative Genomics and Environmental Research of Asthma (IGERA)" (phs001603.v2.p1)  
389 was performed at Northwest Genomics Center (HHSN268201600032I). Genome Sequencing for "NHLBI  
390 TOPMed: Children's Health Study (CHS) Effects of Air Pollution on the Development of Obesity in Children  
391 (Meta-AIR)" (phs001604.v2.p1) was performed at Northwest Genomics Center (HHSN268201600032I).  
392 Genome Sequencing for "NHLBI TOPMed: NHGRI CCDG: The BioMe Biobank at Mount Sinai"  
393 (phs001644.v2.p2) was performed at Baylor College of Medicine Human Genome Sequencing Center  
394 (HHSN268201600033I) and the McDonnell Genome Institute (HHSN268201600037I). Genome Sequencing for  
395 "NHLBI TOPMed: Lung Tissue Research Consortium (LTRC)" (phs001662.v2.p1) was performed at Broad  
396 Institute Genomics Platform (HHSN268201600034I). Genome Sequencing for "NHLBI TOPMed: Childhood  
397 Asthma Management Program (CAMP)" (phs0017265.v2.p1) was performed at Northwest Genomics Center  
398 (HHSN268201600032I). Core support, including centralized genomic read mapping and genotype calling, along  
399 with variant quality metrics and filtering, were provided by the TOPMed Informatics Research Center (3R01HL-  
400 117626-02S1; contract HHSN268201800002I). Core support including phenotype harmonization, data  
401 management, sample-identity QC, and general program coordination were provided by the TOPMed Data  
402 Coordinating Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I). We gratefully  
403 acknowledge the studies and participants who provided biological samples and data for TOPMed. The full study-  
404 specific acknowledgments are included in Supplementary Information 2.



## 405 Methods

### 406 Data Description, Quality control, and imputation

407 We have gathered data sets for the GLADdb by combining accessible genomic information from Whole-Genome  
408 Sequencing (WGS) and microarray genotyping chip sources. We have requested and received access to 39  
409 dbGaP cohorts. Another important source was the WGS projects in TOPMed<sup>5</sup>. In total, we have explored over  
410 268K samples in detail to find 70,702 Latin American subjects for this initial set. This search includes 172K from  
411 general dbGaP datasets including the eMERGE<sup>68</sup>, PAGE<sup>69</sup>, and SIGMA<sup>9</sup> projects (**Table S1**). **Figure S1**  
412 shows our general workflow. For each non-WGS dataset (**Table S1**), we converted their genome coordinates  
413 (liftover) from the original reference (NCBI36/hg18 or GRCh37/hg19) to the genome reference GRCh38/hg38  
414 using picard<sup>70</sup>. After a first liftover run, we used the strand flip option of PLINK<sup>71</sup> on the rejected variants and  
415 performed a second liftover run. Furthermore, variants were filtered using PLINK for 5% missingness, a p-value  
416 less than  $1 \times 10^{-6}$  on the Hardy Weinberg exact test (HWE), keeping only biallelic autosomal variants with a  
417 minimum minor allele frequency (MAF) of 1%. Samples were filtered for 5% missingness and heterozygosity  
418 exceeding three times the standard deviation from the mean. Also, a linkage disequilibrium (LD) pruned dataset  
419 was created using PLINK's indep-pairwise algorithm using the parameters 50 10 0.1.

420  
421 For each data set for which we acquired genomic information and appropriate consent, we evaluated self-  
422 described demographic variables such as an ethnic designation of Hispanic/Latino. We included the entire  
423 cohorts where the primary study design was focused on Latin American individuals, e.g. SIGMA<sup>9</sup>. For the  
424 remaining datasets, many without demographic information provided via dbGaP, we identified possible Latin  
425 American individuals using genetic clustering analysis<sup>39</sup>.

426  
427 We merged each of these remaining datasets (the LD pruned data) with a custom panel of 361 individuals to  
428 assess genome-wide ancestry proportions for European, African, East Asian, and Indigenous American  
429 ancestry. This custom panel included 100 each for European, African, and East Asian from the high coverage  
430 1000 Genomes Project data<sup>38</sup> (**Table S3**). In addition, we included 61 unrelated, previously estimated as near  
431 100%, Indigenous American high coverage genomes from the Peruvian Genome Project<sup>6</sup>. Each data set was  
432 combined with this reference sample, then we ran a supervised ADMIXTURE analysis<sup>39</sup>. These results were  
433 then evaluated for admixture proportions and any sample found to have greater than 2% Indigenous American  
434 ancestry was extracted and included for additional analyses. These samples were then designated as  
435 *admixture-defined*, which will persist in our evaluations of the database as to their utility as matches or  
436 exclusion.

437  
438 After we collected all self-described and admixture-defined individuals in each dataset (non-LD pruned data), we  
439 imputed the genotype panel against the TOPMed Imputation server<sup>40</sup>. The TOPMed imputation panel contains  
440 over 90K individuals and was shown to accurately impute Latin Americans<sup>5</sup>. To date, after we have combined  
441 across all studies analyzed, including the non-imputed TOPMed WGS data, we have 63,589 non-duplicate  
442 samples. This comprises 9,121,629 variants with an imputation  $r^2 > 0.3$  across all datasets (i.e. no missing data)  
443 and includes 8,626,916 SNPs and 494,713 INDELS.

444  
445 Importantly, GLADdb includes 30,078 individuals with non-ambiguous geographical information (**Table S2**). This  
446 means that we have country-level or, in some cases, state or city-level information like Peru, Brazil, and the  
447 USA. For the latter three groups, we did not include individuals without state-level information. A particular case  
448 is the Rio Grande do Sul state in South Brazil. Two of three cohorts that were sampled in this state correspond  
449 to specific cities (Porto Alegre and Pelotas) and were considered as independent groups. To support the  
450 clustering of individuals of different project into groups of similar geographical regions (e.g., USA-Wisconsin,  
451 Chile, Brazil-São Paulo), we performed an Fst analysis. We calculated the Fst among individuals sampled by  
452 different projects but of the same sample region. No regional cluster showed an Fst value above 0.07 (**Table**  
453 **S2**). Finally, these 30K individuals were organized into 45 different regions (**Table S2**). We used this information  
454 for ROH and IBD analyses.

455 After imputation, for each dataset, we kept only variants with  $r^2 > 0.9$ . Then, we merged all datasets and  
456 removed variants with missing information in more than 0.1% of the final dataset using bcftools:

```
457 bcftools filter -e 'F_MISSING > 0.001' ${mergedGLAD} -O b -o $QC1
```

459 For normalizing and keeping biallelic SNPs we applied the following command line:

```
460 bcftools norm -m +any -s $QC1 | bcftools view -m2 -M2 -v snps | bcftools sort -O b -o $GLAD
```

463  
464 Our initial freeze of GLADdb consists of 3,248,494 biallelic SNPs ( $r^2 > 0.9$ ) and 63,589 individuals (**R0.9 dataset**).

466  
467 To avoid any phase issues during the merging process, we infer the haplotype phase for the complete GLADdb  
468 using SHAPEIT ver4<sup>72</sup> using the TOPMed freeze9 dataset<sup>5</sup> (160K individuals) as a reference panel. We ran  
469 SHAPEIT with the following parameters

```
470 shapeit4 --input $GLAD --map $map --thread 60 --region chr${chr} --reference $TOPMEDRef --output $Phased_GLAD  
471 --log phased_chr${chr}.log --mcmc-iterations 10b,1p,1b,1p,1b,1p,1b,1p,10m
```

## 473 Identical-by-Descent and Relatedness analyses

474 Phased biallelic **R0.9 dataset** together with HapMap genetic maps (GRCh38) were used as input for inferences  
475 of IBD (Identical-By-Descent) segments using hap-ibd<sup>73</sup>. For hap-ibd, we set the parameters “min-seed=3” and  
476 “min-output=3” to reduce the rate of false positiveness; defaults were used for all the other parameters. Given  
477 IBD coverage is dramatically increased by the paucity of SNP markers, we defined low SNP density regions as  
478 1-cM windows with the number of SNPs less than 30 and processed all IBD segments overlapping with these  
479 regions by splitting them and removing the parts within the low SNP density regions. The processed IBD  
480 segments were then used as input for ancestry-specific downstream analysis. For non-ancestry specific  
481 analyses, we further merged and flattened the processed IBD segments for each sample-pair when two  
482 segments are either overlapping or close (gap no longer than 0.6cM and the number of phasing-informative  
483 discordant markers no more than 1)<sup>74</sup>. The flattened and merged IBD segments were kept if the segment length  
484  $\geq 5$ cM. Genome-wide total IBD length of all segments shared by each sample pair was then calculated and  
485 organized to an IBD matrix with each element representing the relatedness between a pair of individuals. For  
486 agglomerative clustering, we transformed the matrix into a dissimilarity matrix by the formula  $X = (\max - \min) / (X - \min + 1e-9)$ . The IBD post-processing steps including encoding, removing low SNP density regions, decoding,  
487 sorting, merging, filtering, and matrix-building were implemented in a C++ toolkit *ibdtools*  
488 (<https://github.com/umb-oconnorgroup/ibdtools>) to accelerate the computation for large IBD datasets, for  
489 instance, hundreds of billions of IBD segments.

491  
492 We estimated the kinship coefficient for each pair of individuals in GLADdb with IBDkin<sup>75</sup>. After kinship  
493 coefficient inferences, we pruned for relatedness in GLADdb using NAToRA<sup>76</sup> to exclude the minimum number  
494 of related individuals while removing the main kinship relationships in the dataset. We used 0.03125 as the  
495 kinship coefficient threshold which is the theoretical kinship coefficient expected for a 4th degree relationship.

## 496 Continental Population Structure

### 497 PCA and UMAP

498 Prior to performing dimensionality reduction, we used PLINK<sup>71</sup> to narrow our biallelic **R0.9 dataset** by applying  
499 LD pruning with a threshold of 0.5 to all 54K samples. Then, using the scikit-learn<sup>77</sup> implementation, we ran  
500 Principal Components Analysis (PCA) on the LD-pruned sites, keeping the top 50 components. To help with  
501 cluster visualization, we reduced the 50 principal components down to 2 dimensions by applying the UMAP  
502 algorithm, using the umap-learn package<sup>78</sup>, with `n_neighbors` set to 10 and `min_dist` set to 0.25.

## 504 **Runs of Homozygosity (ROH)**

505 We inferred the ROH segments for our 45 Latin American groups and 21 reference populations to explore the  
506 level of homogenization in each group. For each group, we used PLINK to filter for monomorphic variants and  
507 generate a transpose format. Then, we ran GARLIC<sup>44</sup>, software that infers ROH based on the Pemberton *et al*  
508 <sup>79</sup> pipeline detecting short (tens of kb), medium (hundreds of kb to several Mb) and long (tens of Mb) ROH  
509 segments. We set the --auto-winsize mode to allow GARLIC to estimate the best window size for ROH inference  
510 starting from a 50 SNP window. We used the following command line:

```
511  
512 garlic --tped ${pop}.tped --tfam ${pop}.tfam --build hg38 --error 0.001 --cm --winsize 50 --auto-winsize --  
513 auto-winsize-step 10 --out roh_autosize_${pop} --threads 20 --map ${geneticmap}
```

514  
515 For each individual in each group, we summed all ROH sizes to determine the genome-wide amount of ROH.  
516 Considering the good representation of Peruvian and Brazilian regions, 13 and 12 respectively, we estimated  
517 the correlation between the median for each group and the average genome-wide ancestry proportion in each  
518 country using the Pearson correlation. Processing and plotting scripts are available in: [https://github.com/umb-](https://github.com/umb-oconnorgroup/GLAD_DemographicAnalysis)  
519 [oconnorgroup/GLAD\\_DemographicAnalysis](https://github.com/umb-oconnorgroup/GLAD_DemographicAnalysis)

## 521 **Local ancestry Inferences**

522 We ran local ancestry inference using RFMix ver2<sup>80</sup> on GLADdb. We inferred local ancestry for the phased  
523 dataset considering two Expectation-Maximization runs and eight generations since admixture. For the ancestry  
524 reference panel, we selected 982 individuals including 250 Europeans, 250 East Asian, 250 Africans and 232  
525 individuals with predominant Indigenous American ancestry (Table S2). Europeans, Africans and East Asian  
526 reference populations are part of the 1000 Genomes Project. Individuals with predominant Indigenous American  
527 ancestry includes Indigenous Americans from the Peruvian Genome Project<sup>6,7</sup> and individuals with predominant  
528 Indigenous American ancestry (above 99% of Indigenous American ancestry) from Guatemala (Table S2).

## 529 **Distant genetic relatedness among Latin American groups**

### 530 **IBD-community detection**

531 For community detection, we calculated an IBD matrix by summing up all IBD segments with length within a  
532 specific range (>5cM, 5-9.3 or >9.3cM) across the genome for each pair of individuals, and set all elements with  
533 values < 12 cM to 0 in this matrix to reduce the density of non-zero elements in the matrix. The resulting  
534 symmetrical matrix was used as a weighted-adjacency matrix to build a bidirectional relatedness network. We  
535 used the infomap algorithm implemented with the python-igraph<sup>81</sup> package to infer the community structure of  
536 the relatedness network. We kept individuals within the top 20 communities and with a degree >= 30  
537 connections and used the Fruchterman Reingold layout<sup>82</sup> for visualization purposes. Community enrichment in  
538 a given birth country is defined as the largest proportion of community labels for individuals born in the country.  
539 The number of communities enriched in a birth country is determined by counting the communities that have  
540 >1% enrichment in this country.

### 542 **IBD sharing among Latin American regions**

543 To explore the recent relationship among Latin American regions, we focused on IBD segments greater than  
544 21.4 cM. We calculated the IBD sharing at intra and interregional levels. For intraregional sharing, we summed  
545 the total amount of shared IBD and divided it by the number of pairs:  $N(N-1)/2$ , where N is the total number of  
546 individuals included for that region. For interregional sharing, we summed the total amount of shared IBD among  
547 individuals of populations 1 and 2 and divided it by  $N_1 \times N_2$ , where  $N_1$  and  $N_2$  are the total number of individuals  
548 included for populations 1 and 2 involved in the sharing, respectively.

### 550 **Ancestry Specific IBD**

551 From the multi-way admixed origin of Latin American populations, IBD (segments greater than 21.4 cM) and  
552 local ancestry analyses provide an opportunity to detect ancestry-specific signatures related to bottleneck  
553 (whitin-region analysis) and recent migration (across-region analysis) along the Americas.

We implemented a python algorithm called *GAFIS* ( that stands for “Getting Ancestry For IBD Segments” ) that uses RFMIX outputs to identify local ancestry labels for an IBD segment shared by a pair of individuals under a certain probability threshold. As a probability threshold for local ancestry inferences in *GAFIS*, we set 90% for a genomic region being of the K ancestry. For this analysis, we included our processed IBD segments to reduce the proportion of false positives. Moreover, if an IBD segment contained several ancestries, we split the segment into segments corresponding to independent ancestries for each pair of individuals.

After ancestry identification of the IBD segments, we filter out ancestry specific-IBD segments based on the following criteria:

- The ancestry profile of one of the individuals for the IBD region was unknown for having a local ancestry probability lower than 90%.

- Both individuals have different ancestry labels of the IBD segment.

After those filters, we kept individuals with demographic information and calculated an *ancestry-specific IBD score* (asIBD score) within and across the 45 Latin American groups. Our asIBD score is defined in the following equations:

Within regions:

$$\frac{\sum_i \sum_i IBD_{ancK}}{N_{region\ i} \times \frac{(N_{region\ i}-1)}{2} \times Proportion_{anc\ K\ region\ i}^2 \times callableIBDlength} \quad (\text{Equation 1})$$

Across regions:

$$\frac{\sum_i \sum_j IBD_{ancK}}{N_{region\ i} \times N_{region\ j} \times Proportion_{anc\ K\ region\ i} \times Proportion_{anc\ K\ region\ j} \times callableIBDlength} \quad (\text{Equation 2})$$

Where:

Anc K= African, European or Indigenous American ancestries

IBD<sub>anc K</sub>: The total amount of ancestry K IBD shared between a pair of individuals from region i and j.

N<sub>region i</sub>: Total number of individuals from region i.

N<sub>region j</sub>: Total number of individuals from region j.

Proportion<sub>anc K region i</sub>: Global ancestry proportion for Ancestry K in region i.

Proportion<sub>anc K region j</sub>: Global ancestry proportion for Ancestry K in region j.

callableIBDlength: Total size of the genome that was included for IBD analysis.

In both equations, in the numerator, for a specific ancestry, we summed the total amount of IBD per ancestry for each pair of individuals from the same region (Equation 1) or between region i and j (Equation 2). To control for sample size and ancestry proportions, for equation 1, we divide the total amount of shared IBD by the product of the total number of combinations of individuals and the square of ancestry proportion. For Equation 2, we divide by the product of sample size for each region and the product of the global ancestry proportion K for each region, respectively. Finally, to get a value relative to the total size of the genome, we included the genome size that was analyzed in the IBD inference in both equations. Codes and pipeline to estimate the asIBD score are available in: [https://github.com/umb-oconnorgroup/GLAD\\_DemographicAnalysis](https://github.com/umb-oconnorgroup/GLAD_DemographicAnalysis)

## Polygenic Risk Scores in Latin American populations

### Description of PRS cohorts

We utilized the following studies participating in GLAD: Columbia University Study of Caribbean Hispanics and Late Onset Alzheimer's disease (phs000496), Slim Initiative in Genomic Medicine for the Americas (SIGMA): Diabetes in Mexico Study (phs001388), eMERGE Network Phase III: HRC Imputed Array Data (phs001584), Early Progression to Active Tuberculosis in Peruvians (phs002025), and EPIGEN-Brasil (Bambui, Pelotas, and



SCAALA). These studies all ascertained one or more of the following traits: height, body mass index (BMI), and/or type 2 diabetes (T2D). See Table S5 for a complete description of cohorts.

### Ancestry proportions, relationship inference, principal components and imputation

Within each cohort, PCs were calculated using PC-Air<sup>83</sup> to utilize as covariates. Related individuals were resolved to the 3rd degree using a kinship matrix generated in *Identical-by-descent and relatedness analyses* section. Genotyped data from each cohort was separately merged with the 1000 Genomes Project (1KGP)<sup>38</sup>. Global ancestry proportions were estimated using ADMIXTURE<sup>39</sup>, a K of 5, and 20 replicates. For PRS estimation, imputed variants were filtered for a minimum imputation  $r^2$  of 0.9 and a MAF of 0.01. Both imputed and genotyped data were down-sampled to Hapmap Phase Three variants as required by PRS-CS<sup>51</sup>. Phenotype data was harmonized across cohorts, though all analyses were conducted on a per-cohort basis.

### GWAS summary statistics

Genome-wide association statistics were obtained from the GWAS Catalog<sup>85</sup>, Biobank Japan<sup>30</sup> (BBJ), and UK Biobank<sup>28</sup> (UKBB). African-ancestry GWAS summary statistics were combined using a random-effects meta-analysis using the GAP package in R to improve sample size. See Table S6 for a description of summary statistics used for this study.

### Heritability estimation

Per-cohort additive heritability for each trait was estimated using GCTA<sup>86</sup>, adjusting for sex, age, age<sup>2</sup>, and PCs 1-10. For each set of GWAS summary statistics, heritability was estimated using LD score regression<sup>87</sup>, using the appropriate 1KGP super-population for the calculation of LD scores.

### Polygenic risk score calculation

Pruning/Thresholding PRS: We used PRS calculated with PRSice-2<sup>50</sup> as the representative pruning and thresholding (P+T) method. For P+T, we trained the  $r^2$  parameters ( $r^2$  thresholds of 0.2, 0.4, 0.6, and 0.8), window size (+/- 250 kb, 500kb, 750kb, 1000 kb), and p-value thresholds (iterated by PRSice-2) in one cohort (eMERGE) and validated the parameters in the other cohorts.

Bayesian Mixture PRS: We used PRS estimated with PRS-CS<sup>51</sup> as the baseline Bayesian mixture method. For PRS-CS, we trained the phi ( $\phi$ ) parameter (phi=1e-06, 1e-04, 1e-02, and 1e+00) in one cohort (eMERGE, as this cohort included information for all tested traits) via a small grid search and validated it in the other cohorts. In addition, we also evaluated the fully Bayesian pseudo-validation method (phi=auto) for obtaining phi.

Multi-ancestry PRS using PRS-CSx: We leveraged PRS-CSx<sup>52</sup> to compute a multi-ancestry PRS, which simultaneously fits multiple sets of GWAS summary statistics while modeling population-specific LD, resulting in more accurate posterior effect sizes for any relatively underpowered GWAS. PRS-CSx outputs a PRS corresponding to each GWAS population and an inverse variance meta-analysis of the posterior effect sizes. We trained the best linear combination of each single-population PRS in one cohort using the mixing weights method proposed by Márquez-Luna et al.<sup>88,89</sup> (Equations 3 and 4) with validation in other cohorts. Prior to combining, each PRS is scaled (mean 0, standard deviation 1). In addition, we also evaluated weighting PRS by ancestry proportions (Equation 5), weighting by ancestry proportions after collapsing East Asian and Indigenous American ancestries (Equation 6), and regressing on ancestry proportions prior to model fitting. We compared these linear combinations to the PRS generated from the inverse-variance meta-analysis of PRS-CSx posterior effect sizes.

Equation 3:  $PRS_i = a PRS_{EAS_i} + (1 - a) PRS_{EUR_i}$ ,

Equation 4:  $PRS_i = a_1 PRS_{EAS_i} + a_2 PRS_{EUR_i} + a_3 PRS_{AFR_i}$ , where  $a_1 + a_2 + a_3 = 1$ ,

Equation 5:  $PRS_i = PRS_{EAS_i}(p_{EAS_i}) + PRS_{EUR_i}(p_{EUR_i}) + PRS_{AFR_i}(p_{AFR_i})$ ,

Equation 6:  $PRS_i = PRS_{EAS_i}(p_{EAS_i} + p_{NAT_i}) + PRS_{EUR_i}(p_{EUR_i}) + PRS_{AFR_i}(p_{AFR_i})$ ,

where  $\alpha$ ,  $\alpha_1$ ,  $\alpha_2$ , and  $\alpha_3$  represent mixing weights,  $PRS_{AFR_i}$ ,  $PRS_{EUR_i}$ , and  $PRS_{EAS_i}$  represent a PRS calculated using African, European, East Asian ancestry GWAS, respectively, for individual  $i$ .  $p_{EAS_i}$ ,  $p_{AFR_i}$ ,  $p_{EUR_i}$ , and  $p_{NAT_i}$  represent the East Asian, African, European, and Indigenous American ancestry proportions for individual  $i$ .

356 For BMI, height, and T2D, GWAS summary statistics from East Asian, European, and African populations are  
357 publicly available (see Table S6). In addition, we were able to train the full range of parameters thanks to  
358 multiple independent Latin American cohorts containing data for these traits. We first compared pruning and  
359 thresholding (P+T), PRS-CS, and PRS-CSx models. We then evaluated PRS-CSx based multi-ancestry models,  
360 comparing linear combinations (the best performing linear combination model for each cohort) and inverse-  
361 variance meta-analyses of PRS-CSx posterior effects. These multi-ancestry models were derived from East  
362 Asian and European GWAS (referred to as SUM2 and META2) or derived from East Asian, European, and  
363 African GWAS (referred to as SUM3 and META3). Finally, we compared these multi-ancestry models against  
364 the best single ancestry PRS (EUR2 and EUR3 estimated using PRS-CSx).

### 365 366 367 **PRS model evaluation**

368 All models were evaluated using the 10-fold cross validation framework outlined by Pain et al <sup>90</sup>. In this  
369 approach, the primary metric is the Pearson correlation between the predicted and true values with a standard  
370 error of  $SE_r = (1 - r^2) / \sqrt{(n - 2)}$ , where  $r$  is the Pearson correlation and  $n$  is the sample size. Correlations were  
371 compared using the two-sided William's test implemented in the psych R package that accounts for the non-  
372 independence of the model predictions.  $R^2$  was calculated as the square of Pearson's  $r$ ; partial  $R^2$  was estimated  
373 by subtracting the  $R^2$  of the base model (only covariates) from  $R^2$  of the full model (covariates and PRS). In  
374 general, the base model included age, age<sup>2</sup>, sex, and PCs 1-10 with the exception of cohorts with a categorical  
375 age variable (eMERGE for T2D). The Pelotas cohort, as a birth-year cohort, age and age<sup>2</sup> were not included as  
376 all subjects were the same age. We tested the association of mean ancestry proportions of the cohorts with  
377 model performance using linear regression, adjusting for the GWAS trait ( $R^2$ - scaled ancestry proportion +  
378 trait).

### 379 380 **Code Availability**

381 From the code utilized for this project, we developed an R package called PRSHelpDesk that supports PRS  
382 estimation and evaluation. It is available on GitHub at <https://github.com/dloesch/PRSHelpDesk>.

## 384 **Matching**

385 Both the baseline bipartite matching <sup>49</sup> algorithm and the nearest neighbor simulated annealing matching  
386 algorithm operate on a principal components space composed of the first 50 components computed using  
387 246,799 LD-pruned SNPs from GLADdb. The external-user-provided query is also embedded into the PCA  
388 space with a saved transformation matrix and pairwise distances are computed with a variance-weighted  
389 Minkowski distance metric. Once a suitable matching set has been found, we return summary statistics to the  
390 external user including alternate allele frequency, genotype counts, and haplotype ancestry counts by segment.

391  
392 The baseline algorithm is outlined in **Algorithm 1** and consists of iteratively applying scikit-learn's <sup>77</sup> bipartite  
393 matching implementation until enough controls have been found.

394  
395 Given a desired control cohort size  $m$  and hyperparameters  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $n$ , the nearest neighbor simulated  
396 annealing matching algorithm, outlined in **Algorithm 2**, proceeds as follows. The computed pairwise distances  
397 between query and GLADdb PCA embeddings are used to find the  $\alpha$  nearest neighbors of each query genome  
398 from the potential controls, which we then merge into a candidate set. We sample  $m$  controls from the candidate  
399 set and do so  $\beta$  times to generate  $\beta$  control cohorts. We use the genomic control  $\gamma$ , calculated between a  
700 control cohort and the query, to evaluate the  $\beta$  control cohorts. The  $\gamma$  values are then used to select the optimal  
701 starting control cohort and a function of their standard deviation is used to initialize our simulated annealing  
702 temperature. We perform simulated annealing for  $n$  iterations, randomly swapping  $\alpha$  genomes between our  
703 control cohort and the candidate set at each iteration, evaluating the control cohort by its genomic control  $\gamma$ .

---

**Algorithm 1 Greedy Control Match**

---

**Input:**  $m$  - number of matches,  $dist$  - distance metric,  $E$  -  $q \times e$  embedding matrix of the query,  $\mathcal{E}$  -  $d \times e$  embedding matrix of the database, ( $q$  is the number of query genotypes,  $d$  is the number of dataset genotypes, and  $e$  is the embedding dimension)

**Output:**  $M$  satisfying  $M \subseteq \mathbb{Z} \cup [1, d]$ ,  $|M| = m \vee |M| = d$

```

1:  $M \leftarrow \emptyset$ 
2: while  $|M| < m \wedge |M| < d$  do
3:    $X = \arg \min_X (\sum_{i=1}^q dist(E_{Q,i}, E_{D,x_i}) | x_i \in X \subset \mathbb{Z}, |X| = q, 1 \leq x_i \leq d)$ 
4:   if  $m - |M| \geq q$  then
5:      $M \leftarrow M \cup X$ 
6:   else
7:     while  $|M| < m \wedge |M| < d$  do
8:        $x \sim \mathcal{U}(X)$ 
9:        $M \leftarrow M \cup \{x\}$ 
10:       $X \leftarrow X \setminus \{x\}$ 
11:    end while
12:  end if
13: end while
14: return  $M$ 

```

---

**Algorithm 2 Simulated Annealing Control Match**

---

**Input:**  $m$  - number of matches,  $dist$  - distance metric,  $eval$  - evaluation metric,  $E$  -  $q \times e$  embedding matrix of the query,  $\mathcal{E}$  -  $d \times e$  embedding matrix of the database,  $G$  -  $s \times q \times 2$  genotype tensor of the query,  $\mathcal{G}$  -  $s \times d \times 2$  genotype tensor of the database,  $n$  - number of simulated annealing iterations,  $\alpha$  - number of nearest neighbors to consider,  $\beta$  - number of starting configurations to choose from,  $\gamma$  - number of individuals per iteration to swap, ( $s$  is the number of SNPs,  $q$  is the number of query genotypes,  $d$  is the number of dataset genotypes, and  $e$  is the embedding dimension)

**Output:**  $M^*$  satisfying  $M^* \subseteq \mathbb{Z} \cup [1, d]$ ,  $|M^*| = m \vee |M^*| = d$

```

1:  $C \leftarrow \emptyset$  ▷ Define set of candidate matches with nearest neighbors
2: for  $i \leftarrow 0$  to  $q$  do
3:    $K \leftarrow \mathbb{Z} \cup [1, d]$ 
4:   for  $j \leftarrow 0$  to  $\alpha$  do
5:      $C \leftarrow C \cup \{\arg \min_k dist(E_i, \mathcal{E}_k)\}$ 
6:      $K \leftarrow K \setminus \{k\}$ 
7:   end for
8: end for
9:  $\mathcal{M} \leftarrow \emptyset$  ▷ Select best starting match set from several random trials
10: for  $i \leftarrow 0$  to  $\beta$  do
11:    $M \leftarrow \emptyset$ 
12:    $X \leftarrow C$ 
13:   while  $|M| < m \wedge |M| < d$  do
14:      $x \sim \mathcal{U}(X)$ 
15:      $M \leftarrow M \cup \{x\}$ 
16:      $X \leftarrow X \setminus \{x\}$ 
17:   end while
18:    $\mathcal{M} \leftarrow \mathcal{M} \cup \{M\}$ 
19: end for
20:  $M^* \leftarrow \arg \min_X (eval(X, G, \mathcal{G}) | X \in \mathcal{M})$ 
21:  $M \leftarrow M^*$  ▷ Run simulated annealing
22:  $C \leftarrow C \setminus M$ 
23:  $\sigma = std(\{eval(X) | X \in \mathcal{M}\})$ 
24:  $t_0 = \frac{\sigma}{\log(\sigma)}$ 
25: for  $i \leftarrow 0$  to  $n$  do
26:    $t = \frac{t_0}{1 + \log(1+i)}$ 
27:    $X \leftarrow M$ 
28:   for  $j \leftarrow 0$  to  $\gamma$  do
29:      $x \sim \mathcal{U}(X)$ 
30:      $c \sim \mathcal{U}(C)$ 
31:      $X \leftarrow (X \setminus \{x\}) \cup \{c\}$ 
32:      $C \leftarrow (C \setminus \{c\}) \cup \{x\}$ 
33:   end for
34:   if  $eval(X) < eval(M) \vee \exp(\frac{eval(M) - eval(X)}{t}) > \mathcal{U}_{[0,1]}$  then
35:      $M \leftarrow X$ 
36:     if  $eval(M) < eval(M^*)$  then
37:        $M^* \leftarrow M$ 
38:     end if
39:   end if
40: end for

```

---

707

708

## 709 GLADdb

710 The GLADdb portal provides both visualization and control matching functionality. The visualizations are built  
711 with the Plotly library, enabling in-browser interaction, zooming, and filtering. The control matching page enables  
712 filtering by self-identified ethnicity, phs numbers, and some phenotypic traits. The external user is asked to  
713 prepare and anonymize their data using a Dockerfile provided at [github.com/umb-oconnorgroup/gladprep](https://github.com/umb-oconnorgroup/gladprep).  
714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

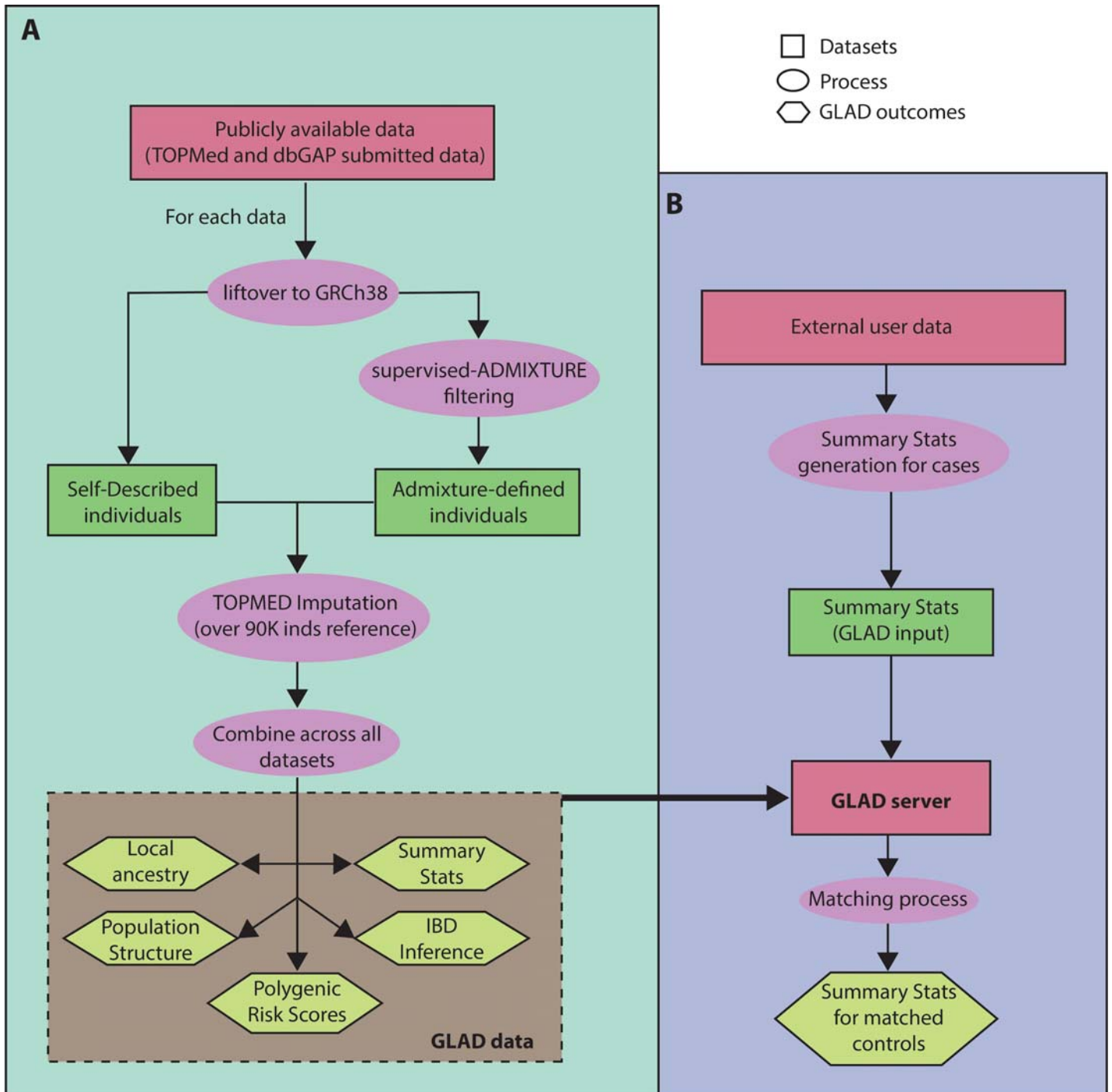
754



755

756

## SUPPLEMENTARY FIGURES



757

758

759

760

761

762

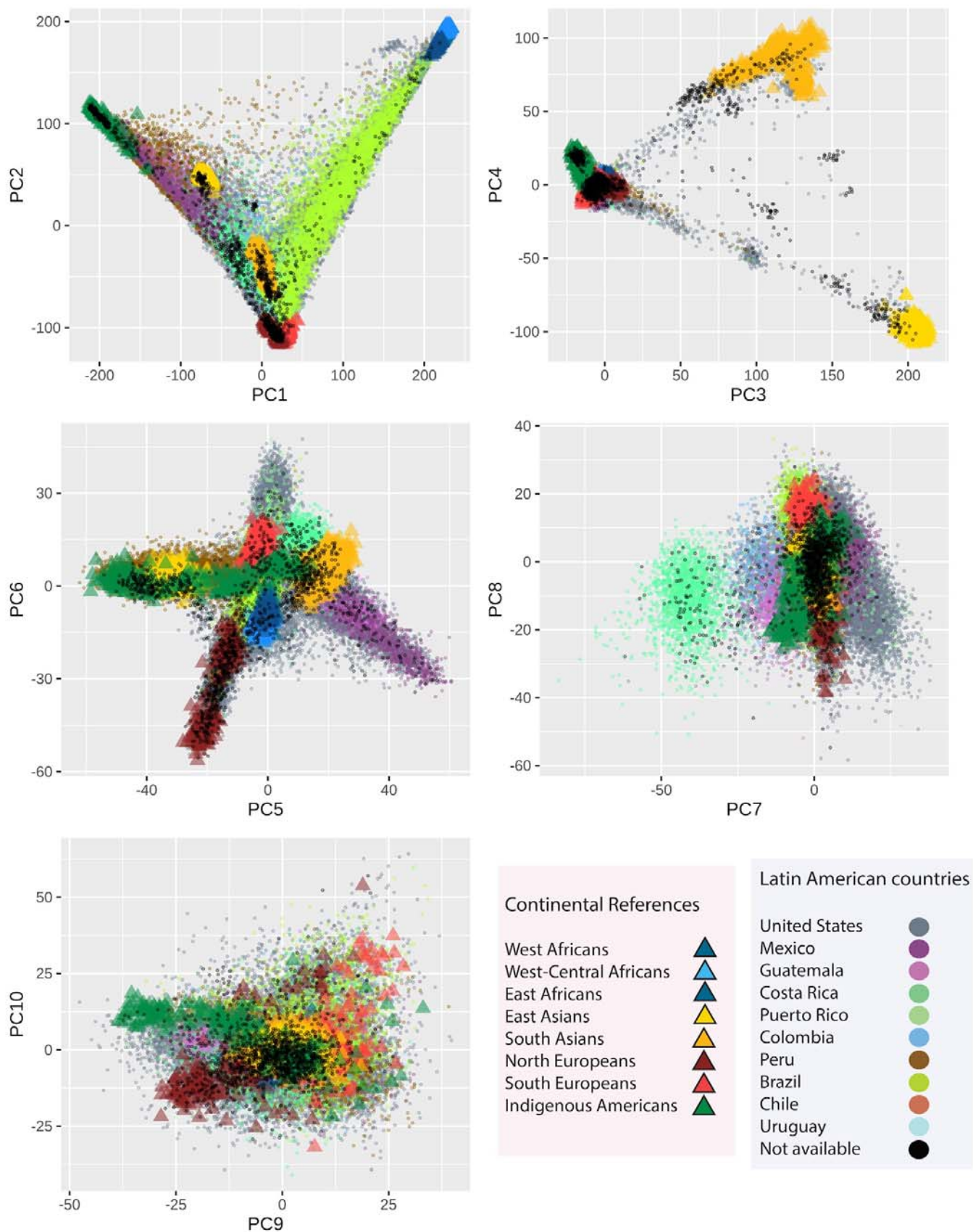
763

764

765

**Figure S1. Workflow for the building and the use of the GLAD database.** A) For each dbGaP cohort, we extracted and self-described Latino and ADMIXTURE defined subjects with at least 2% of Indigenous American ancestry. Then each cohort was imputed in the Michigan Imputation Center using the TOPMED Imputation panel. After imputation, we selected the best imputed loci ( $r^2 > 0.9$ ) and merged the data. We characterized the GLADdb using PCA, IBD and local ancestry analyses. B) By identifying the GLAD individuals that have similar genetic patterns of a query sample, we provide summary statistics of the control subjects from GLADdb.

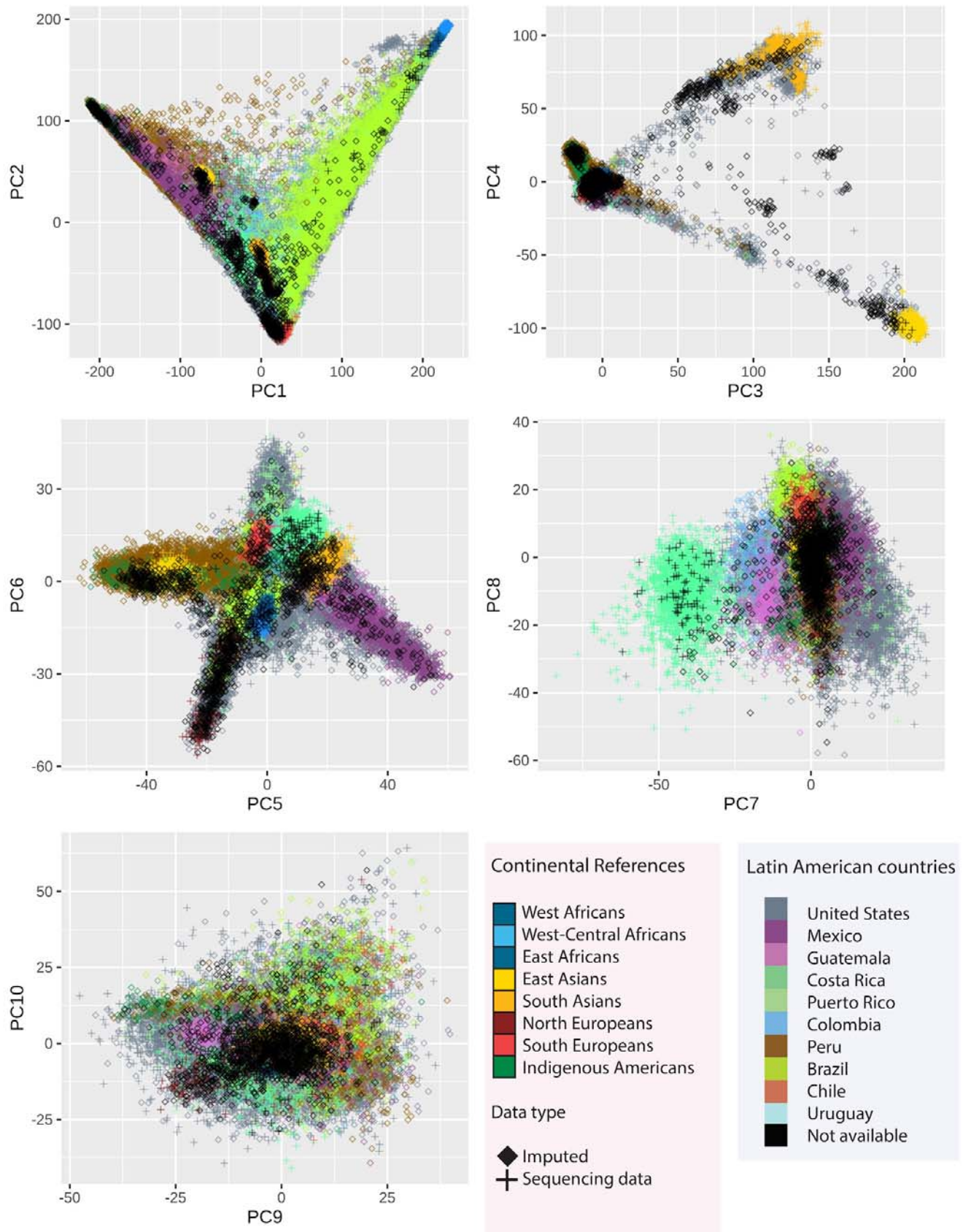
766  
767



768  
769  
770

**Figure S2. Principal component (PC) analysis GLADdb and ancestral reference groups individuals.** First ten PCs that include reference groups (triangles) and GLAD individuals (circles).

771

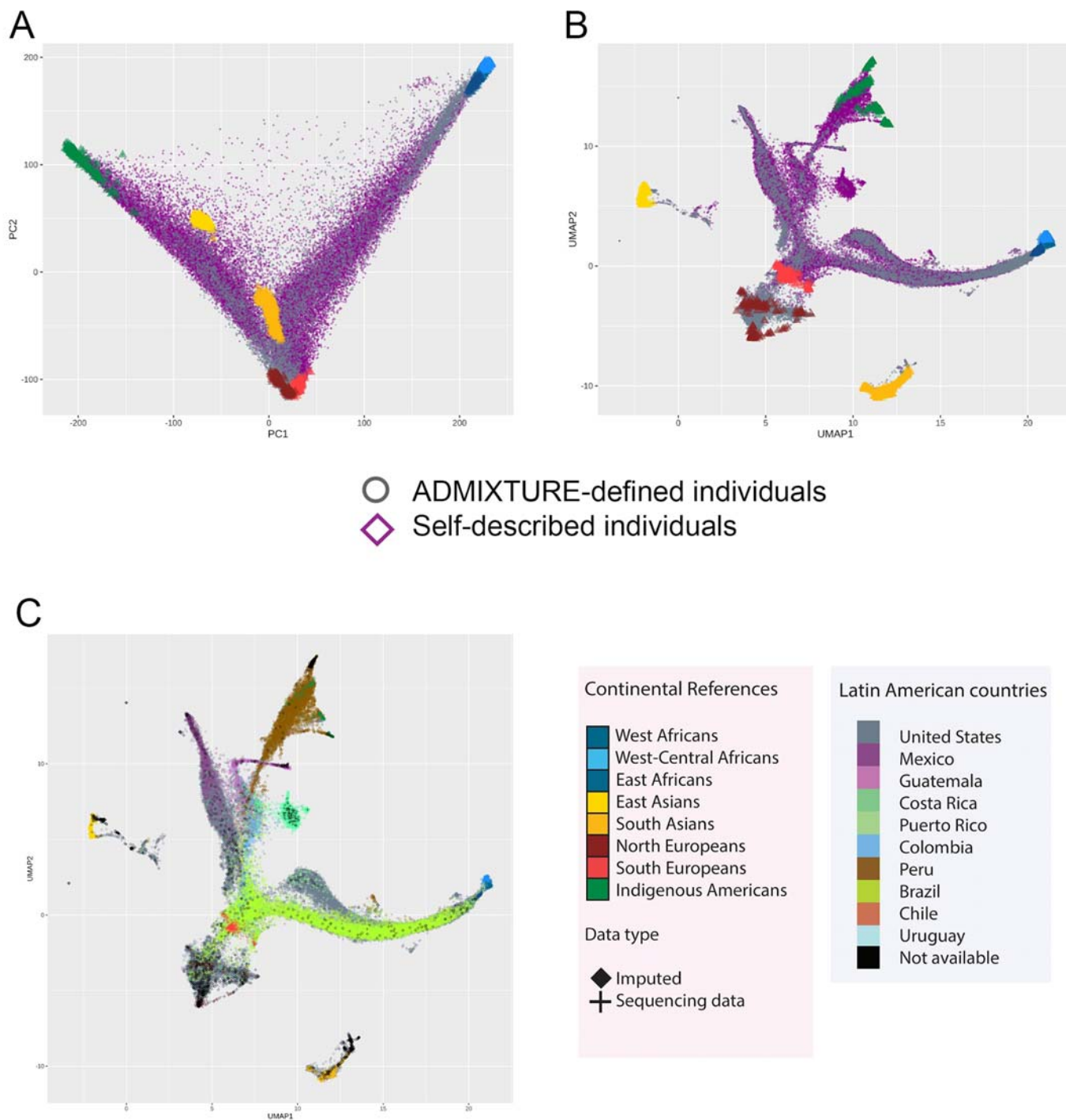


772  
773  
774  
775

**Figure S3. Principal component (PC) analysis GLADdb and ancestral reference groups individuals.** Plot shows the relationships between GLADdb individuals with different data types: Imputed (diamond) and sequencing data (cross).



776  
777

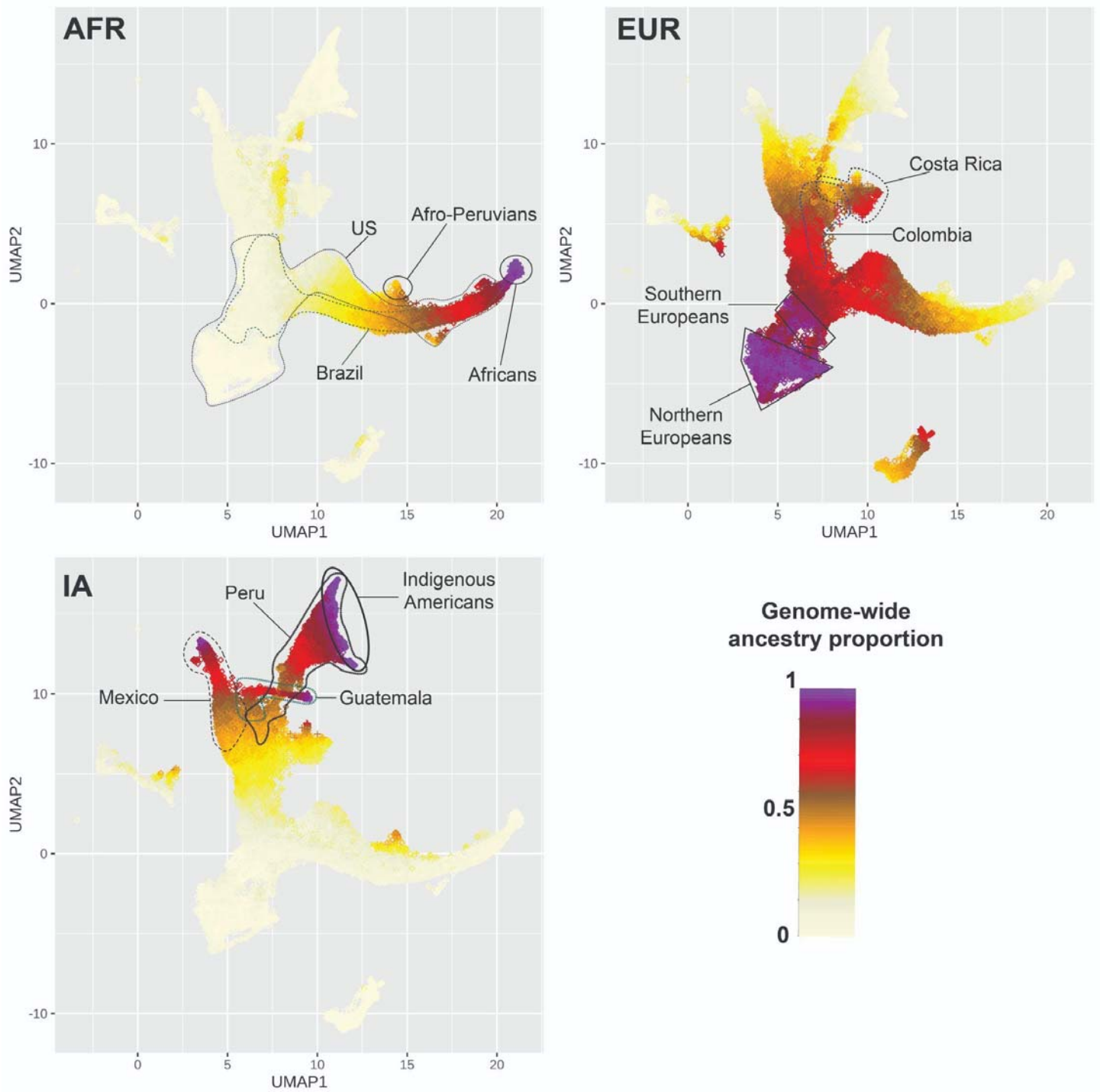


778  
779  
780  
781  
782  
783  
784  
785  
786

**Figure S4. Population structure analysis GLADdb and ancestral reference groups individuals.** Principal component (A) and UMAP (B) analyses showing the relationship between self described and ADMIXTURE-defined individuals in GLADdb. C) UMAP analysis of GLADdb individuals with different data types: Imputed (diamond) and sequencing data (cross).



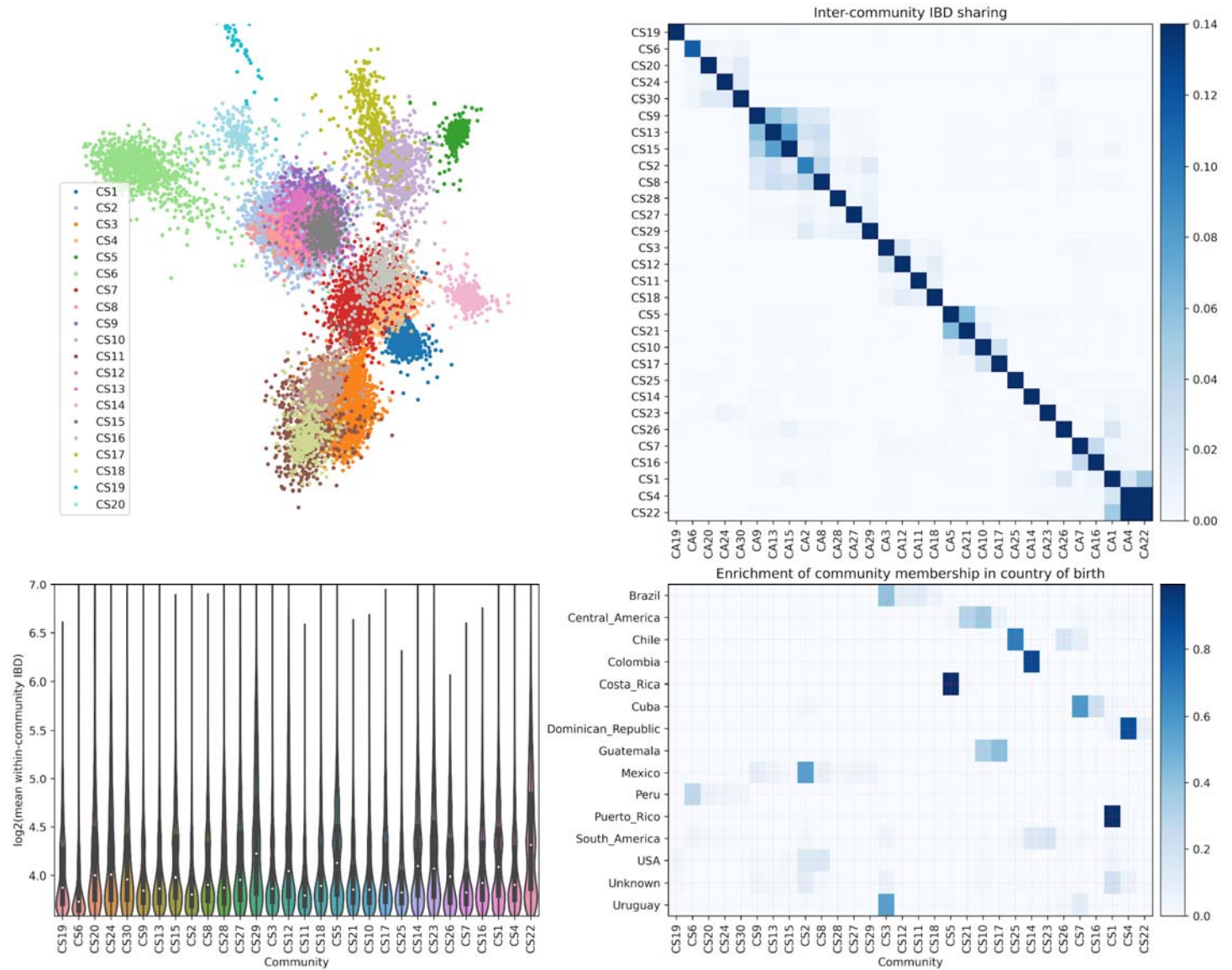
787



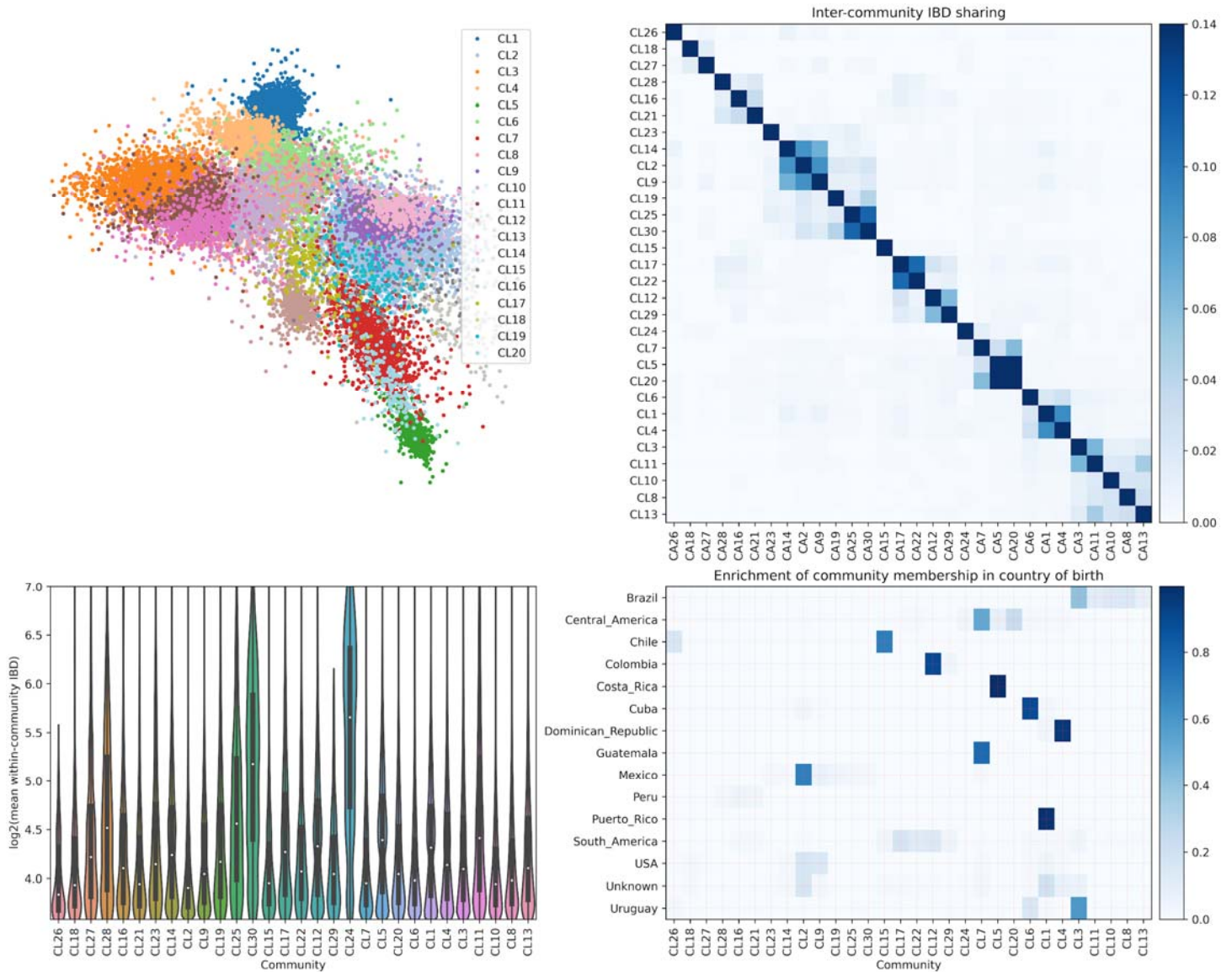
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799

**Figure S5. Genome-wide ancestry clines projected on UMAP analysis.** Continental ancestry clines based on ancestry proportions inferred by ADMIXTURE for African (AFR), European (EUR) and Indigenous American (IA) ancestries in GLADdb individuals.

300



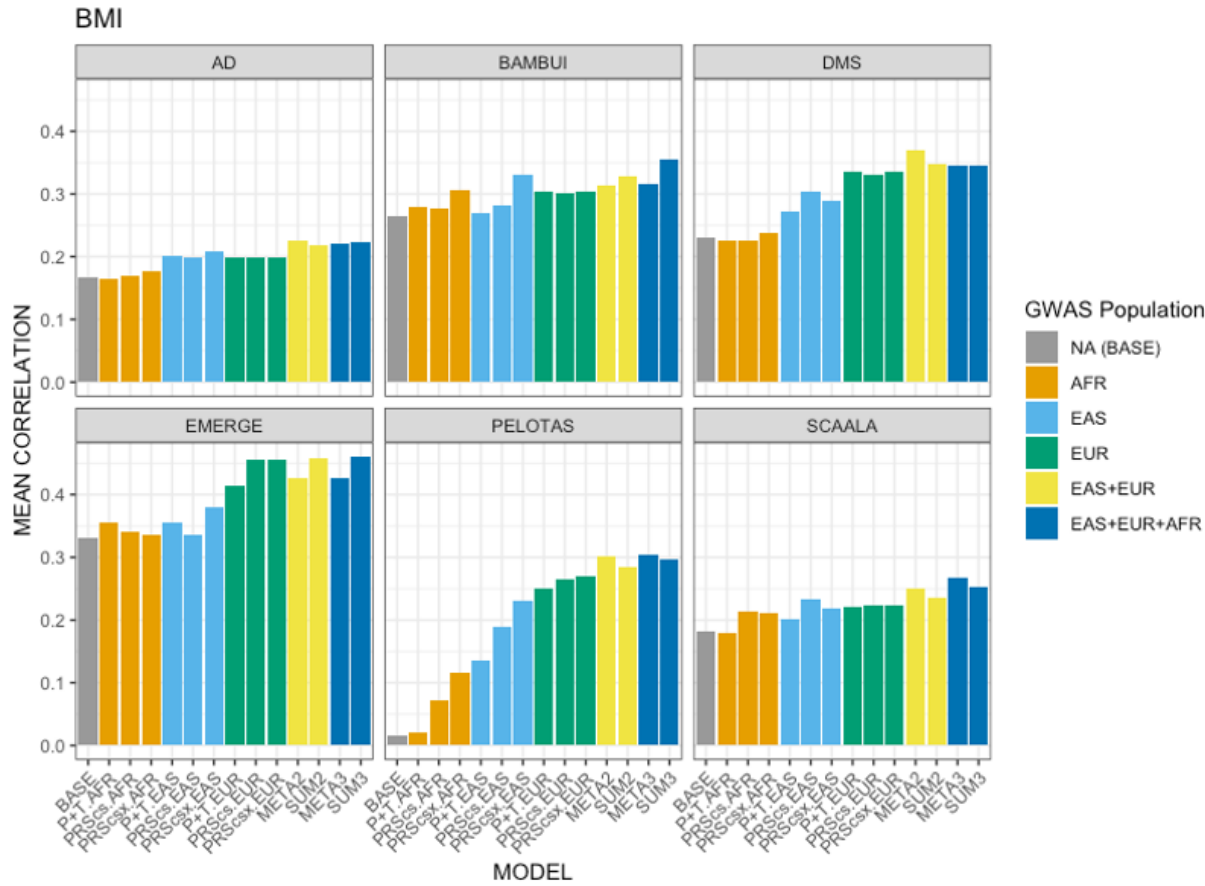
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313



314  
 315  
 316  
 317 **Figure S6B. IBD network community detection using IBD segments greater than 9.3cM.** This interval was selected to explore the  
 318 population dynamics after the colonial times. A) Top 20 IBD network communities visualized using Fruchterman-Reingold layout  
 319 algorithm<sup>82</sup>. For visualization purposes, only individuals with connections > 30 are included in the layout calculation. The community  
 320 labels, such as CL1 and CL2, are named according to the IBD version used and the rank of the community sizes, with CL1 representing  
 321 the largest community when using large IBD segments (> 9.3cM). B) IBD sharing among top 30 inferred communities (ordered by  
 322 agglomerative clustering; the same order was followed in C and D). C) Distribution of IBD shared among individuals in each community.  
 323 D) Enrichment of IBD community membership in the country of origin (i.e., proportions of community labels for individuals born in a given  
 324 country).  
 325  
 326  
 327  
 328  
 329  
 330  
 331  
 332  
 333  
 334  
 335  
 336

337  
338  
339

A.



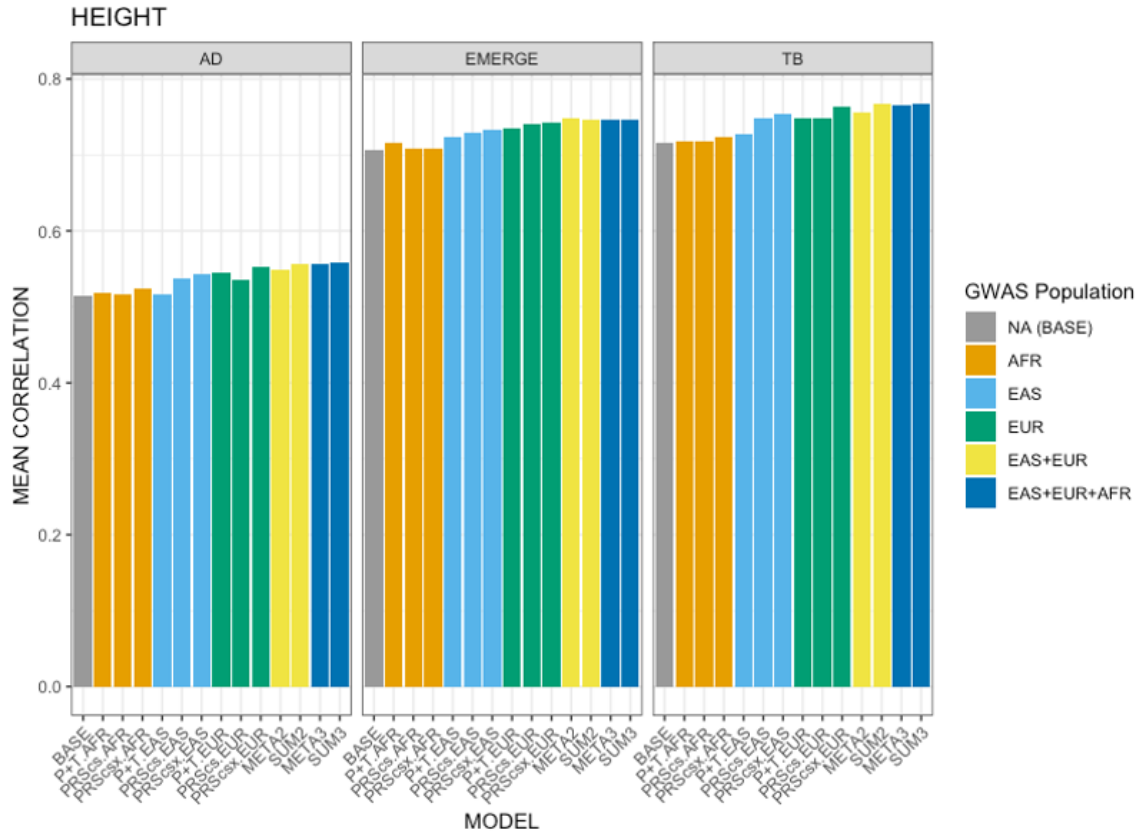
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366

**Figure S7A. Predictive Performance as measured by the mean correlation of the trait with the prediction.** A: Predictive performance for BMI. B: Predictive performance for height. C: Predictive performance for T2D. AD: Columbia University Study of Caribbean Hispanics and Late Onset Alzheimer's disease (phs000496), DMS: Slim Initiative in Genomic Medicine for the Americas (SIGMA): Diabetes in Mexico Study (phs001388), EMERGE: eMERGE Network Phase III: HRC Imputed Array Data (phs001584), TB: Early Progression to Active Tuberculosis in Peruvians (phs002025), and EPIGEN-Brasil (Bambui, Pelotas, and SCAALA).



367  
368  
369

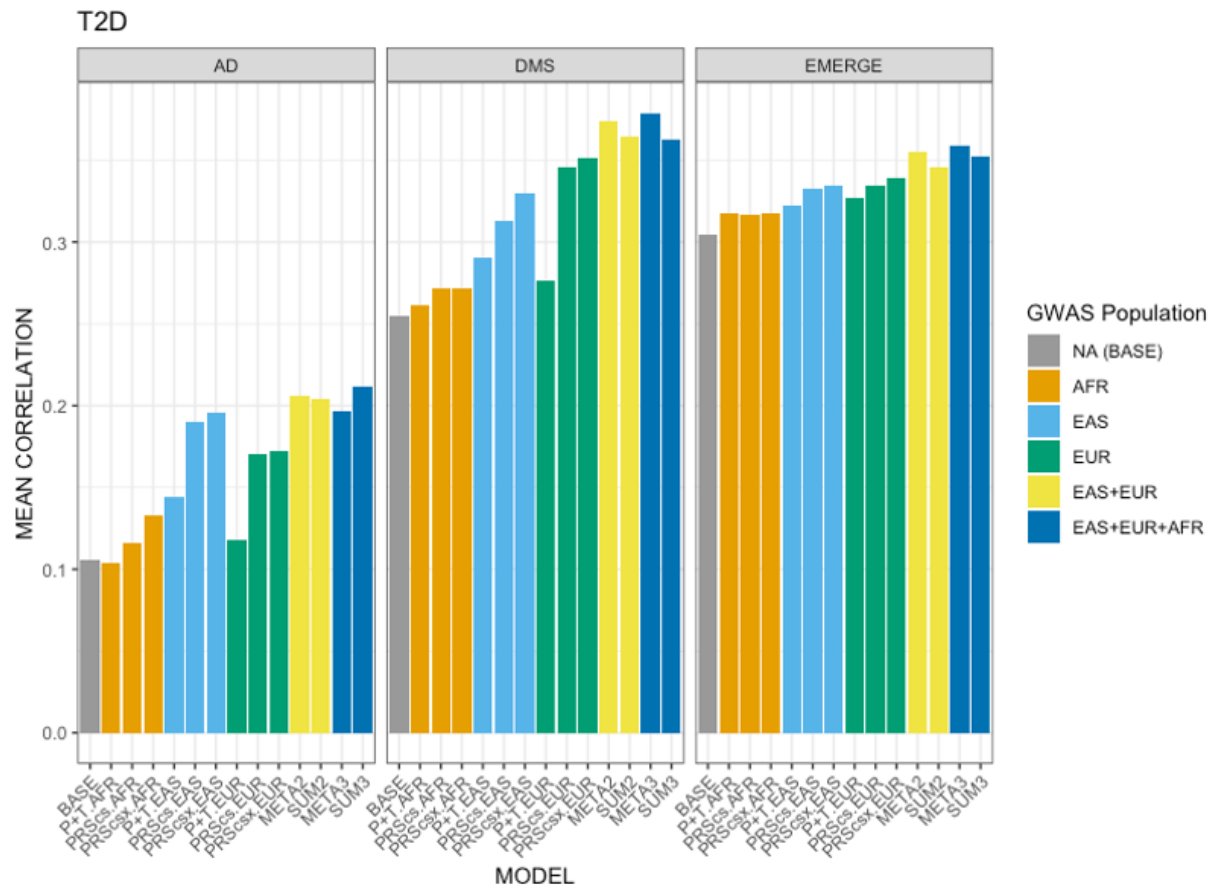
B.



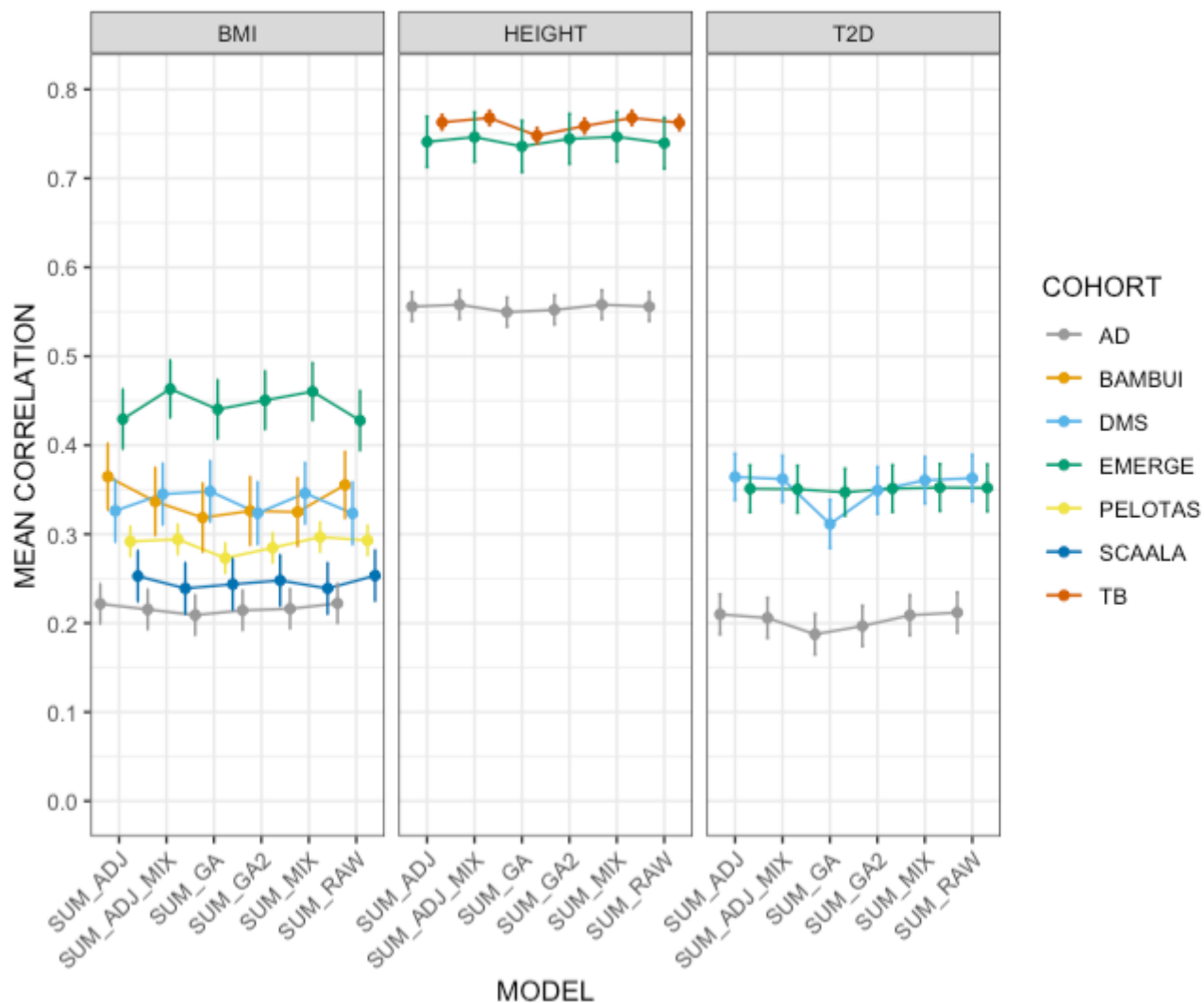
370  
371  
372  
373  
374  
375  
376

**Figure S7B. Predictive Performance by the mean correlation of the trait with the prediction.** A: Predictive performance for BMI. B: Predictive performance for height. C: Predictive performance for T2D. AD: Columbia University Study of Caribbean Hispanics and Late Onset Alzheimer's disease (phs000496), DMS: Slim Initiative in Genomic Medicine for the Americas (SIGMA): Diabetes in Mexico Study (phs001388), EMERGE: eMERGE Network Phase III: HRC Imputed Array Data (phs001584), TB: Early Progression to Active Tuberculosis in Peruvians (phs002025), and EPIGEN-Brasil (Bambui, Pelotas, and SCAALA).

377 C.



378  
 379 **Figure S7C. Predictive Performance as measured by the mean correlation of the trait with the prediction.** A: Predictive  
 380 performance for BMI. B: Predictive performance for height. C: Predictive performance for T2D. AD: Columbia University Study of  
 381 Caribbean Hispanics and Late Onset Alzheimer's disease (phs000496), DMS: Slim Initiative in Genomic Medicine for the Americas  
 382 (SIGMA): Diabetes in Mexico Study (phs001388), EMERGE: eMERGE Network Phase III: HRC Imputed Array Data (phs001584), TB:  
 383 Early Progression to Active Tuberculosis in Peruvians (phs002025), and EPIGEN-Brasil (Bambui, Pelotas, and SCAALA).  
 384  
 385



**Figure S8. PRS linear combination methods across three traits.** Predictive performance of linear combination methods across 3 traits and 7 cohorts. Error bars represent the standard error of the correlation. See methods for model definitions. AD: Columbia University Study of Caribbean Hispanics and Late Onset Alzheimer's disease (phs000496), DMS: Slim Initiative in Genomic Medicine for the Americas (SIGMA): Diabetes in Mexico Study (phs001388), EMERGE: eMERGE Network Phase III: HRC Imputed Array Data (phs001584), TB: Early Progression to Active Tuberculosis in Peruvians (phs002025), and EPIGEN-Brazil (Bambui, Pelotas, and SCAALA).

386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407

308  
309  
310  
311



312  
313 **Figure S9. Screenshots from GLADdb website.** A) Landing page. B) The control matching page, encompassing data preparation and  
314 matching job submission. C-E) Visualization pages, showing respectively all cohorts, selected cohorts, and a zoomed-in view with a  
315 highlighted individual.  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353



## References

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

1. Manichaikul, A. *et al.* Population Structure of Hispanics in the United States: The Multi-Ethnic Study of Atherosclerosis. *PLoS Genet.* **8**, e1002640 (2012).
2. Plecher, H. Latin America - Statistics & Facts. *Statista* <https://www.statista.com/topics/3287/latin-america/> (2019).
3. Noe-Bustamante, L., Hugo Lopez, M. & Manuel Krogstad, J. U.S. Hispanic population surpassed 60 million in 2019, but growth has slowed. *Pew Research Center: U.S. Hispanic population surpassed 60 million in 2019, but growth has slowed* <https://www.pewresearch.org/fact-tank/2020/07/07/u-s-hispanic-population-surpassed-60-million-in-2019-but-growth-has-slowed/> (2020).
4. Mills, M. C. & Rahal, C. The GWAS Diversity Monitor tracks diversity by disease in real time. *Nat. Genet.* **52**, 242–243 (2020).
5. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**, 290–299 (2021).
6. Harris, D. N. *et al.* Evolutionary genomic dynamics of Peruvians before, during, and after the Inca Empire. *Proc. Natl. Acad. Sci.* **115**, E6526–E6535 (2018).
7. Borda, V. *et al.* The genetic structure and adaptation of Andean highlanders and Amazonians are influenced by the interplay between geography and culture. *Proc. Natl. Acad. Sci.* **117**, 32557–32565 (2020).
8. Loesch, D. P. *et al.* Characterizing the Genetic Architecture of Parkinson’s Disease in Latinos. *Ann. Neurol.* **90**, 353–365 (2021).
9. SIGMA Type 2 Diabetes Consortium *et al.* Association of a low-frequency variant in HNF1A with type 2 diabetes in a Latino population. *JAMA* **311**, 2305–2314 (2014).
10. Pino-Yanes, M. *et al.* Genome-wide association study and admixture mapping reveal new loci associated with total IgE levels in Latinos. *J. Allergy Clin. Immunol.* **135**, 1502–1510 (2015).
11. Moreno-Estrada, A. *et al.* The genetics of Mexico recapitulates Native American substructure and affects biomedical traits. *Science* **344**, 1280–1285 (2014).
12. Moreno-Estrada, A. *et al.* Reconstructing the Population Genetic History of the Caribbean. *PLoS Genet.* **9**, e1003925 (2013).
13. Kehdy, F. S. G. *et al.* Origin and dynamics of admixture in Brazilians and its effect on the pattern of deleterious mutations. *Proc. Natl. Acad. Sci.* **112**, 8696–8701 (2015).
14. Adhikari, K. *et al.* A GWAS in Latin Americans highlights the convergent evolution of lighter skin pigmentation in Eurasia. *Nat. Commun.* **10**, 358 (2019).
15. Bonfante, B. *et al.* A GWAS in Latin Americans identifies novel face shape loci, implicating VPS13B and a Denisovan introgressed region in facial variation. *Sci. Adv.* **7**, eabc6160 (2021).
16. Chacón-Duque, J.-C. *et al.* Latin Americans show wide-spread Converso ancestry and imprint of local Native ancestry on physical appearance. *Nat. Commun.* **9**, 5388 (2018).
17. Naslavsky, M. S. *et al.* Whole-genome sequencing of 1,171 elderly admixed individuals from Brazil. *Nat. Commun.* **13**, 1004 (2022).
18. Franceschini, N. *et al.* Variant Discovery and Fine Mapping of Genetic Loci Associated with Blood Pressure Traits in Hispanics and African Americans. *PLOS ONE* **11**, e0164132 (2016).
19. Sofer, T. *et al.* Admixture mapping in the Hispanic Community Health Study/Study of Latinos reveals regions of genetic associations with blood pressure traits. *PLOS ONE* **12**, e0188400 (2017).
20. Qi, Q. *et al.* Genetics of Type 2 Diabetes in U.S. Hispanic/Latino Individuals: Results From the Hispanic Community Health Study/Study of Latinos (HCHS/SOL). *Diabetes* **66**, 1419–1425 (2017).
21. Graff, M. *et al.* Genetic architecture of lipid traits in the Hispanic community health study/study of Latinos. *Lipids Health Dis.* **16**, 200 (2017).
22. Saccone, N. L. *et al.* Genome-Wide Association Study of Heavy Smoking and Daily/Nondaily Smoking in the Hispanic Community Health Study/Study of Latinos (HCHS/SOL). *Nicotine Tob. Res.* **20**, 448–457 (2018).
23. Justice, A. E. *et al.* Genome-wide association study of body fat distribution traits in Hispanics/Latinos from the HCHS/SOL. *Hum. Mol. Genet.* **30**, 2190–2204 (2021).
24. Kerr, K. F. *et al.* Genome-wide association study of heart rate and its variability in Hispanic/Latino cohorts. *Heart Rhythm* **14**, 1675–1684 (2017).
25. Cade, B. E. *et al.* Genetic Associations with Obstructive Sleep Apnea Traits in Hispanic/Latino Americans. *Am. J. Respir. Crit. Care Med.* **194**, 886–897 (2016).

- 007 26. Sofer, T. *et al.* Variants Associated with the Ankle Brachial Index Differ by Hispanic/Latino Ethnic Group: a  
008 genome-wide association study in the Hispanic Community Health Study/Study of Latinos. *Sci. Rep.* **9**,  
009 11410 (2019).
- 010 27. Ziyatdinov, A. *et al.* Genotyping, sequencing and analysis of 140,000 adults from the Mexico City  
011 Prospective Study. (2022) doi:10.1101/2022.06.26.495014.
- 012 28. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of  
013 Complex Diseases of Middle and Old Age. *PLOS Med.* **12**, e1001779 (2015).
- 014 29. Wall, J. D. *et al.* The GenomeAsia 100K Project enables genetic discoveries across Asia. *Nature* **576**, 106–  
015 111 (2019).
- 016 30. Nagai, A. *et al.* Overview of the BioBank Japan Project: Study design and profile. *J. Epidemiol.* **27**, S2–S8  
017 (2017).
- 018 31. Gouveia, M. H. *et al.* Origins, Admixture Dynamics, and Homogenization of the African Gene Pool in the  
019 Americas. *Mol. Biol. Evol.* **37**, 1647–1656 (2020).
- 020 32. Luisi, P. *et al.* Fine-scale genomic analyses of admixed individuals reveal unrecognized genetic ancestry  
021 components in Argentina. *PLOS ONE* **15**, e0233808 (2020).
- 022 33. Nagar, S. D. *et al.* Genetic ancestry and ethnic identity in Ecuador. *Hum. Genet. Genomics Adv.* **2**, 100050  
023 (2021).
- 024 34. Powell, K. The broken promise that undermines human genome research. *Nature* **590**, 198–201 (2021).
- 025 35. Artomov, M., Loboda, A. A., Artyomov, M. N. & Daly, M. J. *A platform for case-control matching enables  
026 association studies without genotype sharing.* <http://biorxiv.org/lookup/doi/10.1101/470450> (2018)  
027 doi:10.1101/470450.
- 028 36. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction across Diverse  
029 Populations. *Am. J. Hum. Genet.* **100**, 635–649 (2017).
- 030 37. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat.  
031 Genet.* **51**, 584–591 (2019).
- 032 38. Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded 1000 Genomes Project  
033 cohort including 602 trios. *Cell* **185**, 3426–3440.e19 (2022).
- 034 39. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated  
035 individuals. *Genome Res.* **19**, 1655–1664 (2009).
- 036 40. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).
- 037 41. Browning, S. R. *et al.* Ancestry-specific recent effective population size in the Americas. *PLOS Genet.* **14**,  
038 e1007385 (2018).
- 039 42. Mooney, J. A. *et al.* Understanding the Hidden Complexity of Latin American Population Isolates. *Am. J.  
040 Hum. Genet.* **103**, 707–726 (2018).
- 041 43. Ongaro, L. *et al.* The Genomic Impact of European Colonization of the Americas. *Curr. Biol.* **29**, 3974–  
042 3986.e4 (2019).
- 043 44. Szpiech, Z. A., Blant, A. & Pemberton, T. J. GARLIC: Genomic Autozygosity Regions Likelihood-based  
044 Inference and Classification. *Bioinformatics* **33**, 2059–2062 (2017).
- 045 45. Baharian, S. *et al.* The Great Migration and African-American Genomic Diversity. *PLOS Genet.* **12**,  
046 e1006059 (2016).
- 047 46. Devlin, B. & Roeder, K. Genomic Control for Association Studies. *Biometrics* **55**, 997–1004 (1999).
- 048 47. Dadd, T., Weale, M. E. & Lewis, C. M. A critical evaluation of genomic control methods for genetic  
049 association studies. *Genet. Epidemiol.* **33**, 290–298 (2009).
- 050 48. Brown, D. W., Myers, T. A. & Machiela, M. J. PCAmatchR: a flexible R package for optimal case–control  
051 matching using weighted principal components. *Bioinformatics* **37**, 1178–1181 (2021).
- 052 49. Munkres, J. Algorithms for the Assignment and Transportation Problems. *J. Soc. Ind. Appl. Math.* **5**, 32–38  
053 (1957).
- 054 50. Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for biobank-scale data. *GigaScience*  
055 **8**, giz082 (2019).
- 056 51. Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A. & Smoller, J. W. Polygenic prediction via Bayesian regression and  
057 continuous shrinkage priors. *Nat. Commun.* **10**, 1776 (2019).
- 058 52. Ruan, Y. *et al.* Improving polygenic prediction in ancestrally diverse populations. *Nat. Genet.* **54**, 573–580  
059 (2022).
- 060 53. Fonseca, L. *et al.* Diversity matters: opportunities in the study of the genetics of psychotic disorders in low-

- 061 and middle-income countries in Latin America. *Braz. J. Psychiatry* **43**, 631–637 (2021).
- 062 54. Durand, J. & Massey, D. S. New World Orders: Continuities and Changes in Latin American Migration. *Ann.*  
063 *Am. Acad. Pol. Soc. Sci.* **630**, 20–52 (2010).
- 064 55. Fleisher, B. M. Some Economic Aspects of Puerto Rican Migration to the United States. *Rev. Econ. Stat.* **45**,  
065 245 (1963).
- 066 56. Meléndez Vélez, E. *Sponsored migration: the state and Puerto Rican postwar migration to the United*  
067 *States*. (The Ohio State University Press, 2017).
- 068 57. Mintz, S. W. PUERTO RICAN EMIGRATION: A THREEFOLD COMPARISON. *Soc. Econ. Stud.* **4**, 311–325  
069 (1955).
- 070 58. Souza, Blase Camacho. Trabajo y Tristeza - 'Work and Sorrow': the Puerto Ricans of Hawaii 1900 to 1902.  
071 *Hawaii. J. Hist.* **18**, 156–173 (1984).
- 072 59. Amaral, E. F. Brazil: internal migration. in *The Encyclopedia of Global Human Migration* (ed. Ness, I.) (Wiley,  
073 2013). doi:10.1002/9781444351071.wbeghm075.
- 074 60. Bastian, M. Brazil, Argentina, Uruguay: Historical and political background. in *Media and Accountability in*  
075 *Latin America* 15–62 (Springer Fachmedien Wiesbaden, 2019). doi:10.1007/978-3-658-24787-4\_2.
- 076 61. Elizaincín, A., Behares, L. E. & Barrios, G. *Nos falemo brasileiro: dialectos portugueses en Uruguay*.  
077 (Editorial Amesur, 1987).
- 078 62. Sirugo, G., Williams, S. M. & Tishkoff, S. A. The Missing Diversity in Human Genetic Studies. *Cell* **177**, 26–  
079 31 (2019).
- 080 63. Wang, Y., Tsuo, K., Kanai, M., Neale, B. M. & Martin, A. R. Challenges and Opportunities for Developing  
081 More Generalizable Polygenic Risk Scores. *Annu. Rev. Biomed. Data Sci.* **5**, 293–320 (2022).
- 082 64. Márquez-Luna, C. *et al.* Incorporating functional priors improves polygenic prediction accuracy in UK  
083 Biobank and 23andMe data sets. *Nat. Commun.* **12**, 6052 (2021).
- 084 65. Weissbrod, O. *et al.* Leveraging fine-mapping and multipopulation training data to improve cross-population  
085 polygenic risk scores. *Nat. Genet.* **54**, 450–458 (2022).
- 086 66. Tian, P. *et al.* Multiethnic polygenic risk prediction in diverse populations through transfer learning. *Front.*  
087 *Genet.* **13**, 906965 (2022).
- 088 67. Castro e Silva, M. A. *et al.* Population Histories and Genomic Diversity of South American Natives. *Mol. Biol.*  
089 *Evol.* **39**, msab339 (2022).
- 090 68. Wiesner, G. L. *et al.* Returning Results in the Genomic Era: Initial Experiences of the eMERGE Network. *J.*  
091 *Pers. Med.* **10**, 30 (2020).
- 092 69. Hu, Y. *et al.* Minority-centric meta-analyses of blood lipid levels identify novel loci in the Population  
093 Architecture using Genomics and Epidemiology (PAGE) study. *PLOS Genet.* **16**, e1008684 (2020).
- 094 70. Picard toolkit. *Broad Institute, GitHub repository* (2019).
- 095 71. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage  
096 Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 097 72. Delaneau, O., Zagury, J.-F., Robinson, M. R., Marchini, J. L. & Dermitzakis, E. T. Accurate, scalable and  
098 integrative haplotype estimation. *Nat. Commun.* **10**, 5436 (2019).
- 099 73. Zhou, Y., Browning, S. R. & Browning, B. L. A Fast and Simple Method for Detecting Identity-by-Descent  
100 Segments in Large-Scale Data. *Am. J. Hum. Genet.* **106**, 426–437 (2020).
- 101 74. Browning, B. L. & Browning, S. R. Improving the Accuracy and Efficiency of Identity-by-Descent Detection in  
102 Population Data. *Genetics* **194**, 459–471 (2013).
- 103 75. Zhou, Y., Browning, S. R. & Browning, B. L. IBDkin: fast estimation of kinship coefficients from identity by  
104 descent segments. *Bioinformatics* **36**, 4519–4520 (2020).
- 105 76. Leal, T. P. *et al.* NAToRA, a relatedness-pruning method to minimize the loss of dataset size in genetic and  
106 omics analyses. *Comput. Struct. Biotechnol. J.* **20**, 1821–1828 (2022).
- 107 77. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. (2012) doi:10.48550/ARXIV.1201.0490.
- 108 78. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension  
109 Reduction. (2018) doi:10.48550/ARXIV.1802.03426.
- 110 79. Pemberton, T. J. *et al.* Genomic patterns of homozygosity in worldwide human populations. *Am. J. Hum.*  
111 *Genet.* **91**, 275–292 (2012).
- 112 80. Maples, B. K., Gravel, S., Kenny, E. E. & Bustamante, C. D. RFMix: A Discriminative Modeling Approach for  
113 Rapid and Robust Local-Ancestry Inference. *Am. J. Hum. Genet.* **93**, 278–288 (2013).
- 114 81. Csardi, G. & Nepusz, T. The igraph software package for complex network research. *InterJournal Complex*

- 115       **Systems**, 1695 (2006).
- 116 82. Fruchterman, T. M. J. & Reingold, E. M. Graph drawing by force-directed placement. *Softw. Pract. Exp.* **21**,
- 117       1129–1164 (1991).
- 118 83. Conomos, M. P., Miller, M. B. & Thornton, T. A. Robust Inference of Population Structure for Ancestry
- 119       Prediction and Correction of Stratification in the Presence of Relatedness. *Genet. Epidemiol.* **39**, 276–293
- 120       (2015).
- 121 84. Conomos, M. P., Reiner, A. P., Weir, B. S. & Thornton, T. A. Model-free Estimation of Recent Genetic
- 122       Relatedness. *Am. J. Hum. Genet.* **98**, 127–148 (2016).
- 123 85. MacArthur, J. *et al.* The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS
- 124       Catalog). *Nucleic Acids Res.* **45**, D896–D901 (2017).
- 125 86. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: A Tool for Genome-wide Complex Trait
- 126       Analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 127 87. Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al.* LD Score regression
- 128       distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295
- 129       (2015).
- 130 88. Márquez-Luna, C., Loh, P.-R., South Asian Type 2 Diabetes (SAT2D) Consortium, The SIGMA Type 2
- 131       Diabetes Consortium & Price, A. L. Multiethnic polygenic risk scores improve risk prediction in diverse
- 132       populations. *Genet. Epidemiol.* **41**, 811–823 (2017).
- 133 89. Bitarello, B. D. & Mathieson, I. Polygenic Scores for Height in Admixed Populations. *G3*
- 134       *GenesGenomesGenetics* **10**, 4027–4036 (2020).
- 135 90. Pain, O. *et al.* Evaluation of polygenic prediction methodology within a reference-standardized framework.
- 136       *PLOS Genet.* **17**, e1009021 (2021).
- 137
- 138
- 139
- 140