1      **Understanding the Hidden Complexity of Latin American Population Isolates**

2

3      *Jazlyn A. Mooney[1,12], Christian D. Huber[2,12], Susan Service[3], Jae Hoon Sul[4], Clare D. Marsden[2],*

4      *Zhongyang Zhang[5,6], Chiara Sabatti[7,8], Andrés Ruiz-Linares[9,10], Gabriel Bedoya[11], Costa*

5      *Rica/Colombia Consortium for Genetic Investigation of Bipolar Endophenotypes, Nelson*

6      *Freimer[3], Kirk E. Lohmueller[1,2] ***

7      [1]Department of Human Genetics, University of California Los Angeles, Los Angeles, CA 90095,
8      USA; [2]Department of Ecology & Evolutionary Biology, University of California Los Angeles,
9      Los Angeles, CA 90095, USA; [3]Center for Neurobehavioral Genetics, Semel Institute for
10     Neuroscience and Human Behavior, University of California Los Angeles, Los Angeles, CA
11     90095, USA; [4]Department of Psychiatry and Biobehavioral Sciences, Semel Center for
12     Informatics and Personalized Genomics, University of California Los Angeles, Los Angeles, CA
13     90095, USA; [5]Department of Genetics and Genomic Sciences, Icahn School of Medicine at
14     Mount Sinai, New York, NY 10029, USA; [6]Icahn Institute for Genomics and Multiscale
15     Biology, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA; [7]Department of
16     Biomedical Data Science, Stanford University, Stanford, CA 94305, USA; [8]Department of
17     Statistics, Stanford University, Stanford, CA 94305, USA; [9]Ministry of Education Key
18     Laboratory of Contemporary Anthropology and Collaborative Innovation Center of Genetics and
19     Development, Fudan University, Shanghai 200438, China; [10]Aix-Marseille Univ, CNRS, EFS,
20     ADES, Marseille, France;[11]Genética Molecular (GENMOL), Universidad de Antioquia,
21     Medellín, Colombia; [12]These authors contributed equally to this work.

22
23
24
25
26
27     **Address correspondence to:**
28     *Kirk E. Lohmueller, Ph.D.
29     Department of Ecology and Evolutionary Biology
30     University of California, Los Angeles
31     621 Charles E. Young Drive South
32     Los Angeles, CA 90095-1606
33
34     Office Phone: (310)-825-7636
35     Fax: (310)-206-0484
36     E-mail: klohmueller@ucla.edu
37
38
39
40
41
42

43    **Abstract:**

44

45    Most population isolates examined to date were founded from a single ancestral population.

46    Consequently, there is limited knowledge about the demographic history of admixed population

47    isolates. Here we investigate genomic diversity of recently admixed population isolates from

48    Costa Rica and Colombia and compare their diversity to a benchmark population isolate, the

49    Finnish. These Latin American isolates originated during the 16$^{th}$ century from admixture

50    between a few hundred European males and Amerindian females, with a limited contribution

51    from African founders. We examine whole genome sequence data from 449 individuals,

52    ascertained as families to build mutigenerational pedigrees, with a mean sequencing depth of

53    coverage of approximately 24X. We find that Latin American isolates have increased genetic

54    diversity relative to the Finnish. However, there is an increase in the amount of identity by

55    descent (IBD) segments in the Latin American isolates relative to the Finnish. The increase in

56    IBD segments is likely a consequence of a very recent and severe population bottleneck during

57    the founding of the admixed population isolates. Furthermore, the proportion of the genome that

58    falls within a long run of homozygosity (ROH) in Costa Rican and Colombian individuals was

59    significantly greater than that in the Finnish, suggesting more recent consanguinity in the Latin

60    American isolates relative to that seen in the Finnish. Lastly, we found that recent consanguinity

61    increased the number of deleterious variants found in the homozygous state, which is relevant if

62    deleterious variants are recessive. Our study suggests there is no single genetic signature of a

63    population isolate.

64

65

66   **Introduction:**

67

68   The use of population isolates to map Mendelian and complex diseases has been a key feature of

69   medical genomics. In addition to experiencing the bottleneck involved with the migration out of

70   Africa, some populations underwent subsequent bottlenecks and remained in relative seclusion

71   afterward. These populations formed present-day isolates[1]. There are numerous benefits to use

72   population isolates to map genes underlying disease. First, and perhaps the largest benefit, is the

73   increased homogeneity of genomes in isolates when compared to outbred populations. Second,

74   isolates experience greater genetic drift than the population from which they were founded. Drift

75   can allow disease causing alleles to exist at an appreciable frequency in isolated populations.

76   Third, isolates may be endogamous and the cryptic relatedness of individuals leads to an

77   enrichment in prevalence of the phenotype of interest and an enrichment of homozygous disease

78   variants[1-4]. Lastly, isolates may also have increased cultural and environmental homogeneity,

79   resulting in a reduction of variability in phenotype due to non-genetic sources.

80          Generally, the genomes of population isolates are thought to exhibit several hallmark

81   features of genetic diversity. Due to bottlenecks associated with their founding, it is thought that

82   isolates should carry lower levels of genetic diversity and lower haplotype diversity than closely

83   related non-isolated populations. Drift experienced by isolates is magnified by small population

84   size, which generates more linkage disequilibrium (LD) than in non-isolated populations. In

85   addition to increased LD, individuals from isolated populations tend to share more regions of the

86   genome identical by descent (IBD) due to small population sizes. Further, due to the isolation

87   after founding and recent mating practices, isolates may have larger regions of the genome found

88   in runs of homozygosity (ROHs) as a result of recent inbreeding. Lastly, bottlenecks and

89    inbreeding should impact patterns of deleterious variation[5–7]. Consequently, one would predict

90    that individuals from isolates will have fewer segregating sites, and the remaining deleterious

91    variants will be segregating at a higher frequency[8]. Indeed, genomic studies over the last decade

92    have documented several of these signatures[2,9,10]. However, it is known that not all isolates share

93    the same demographic history. Therefore, it is essential that we understand how the factors

94    shaping genetic variation in a population, are influenced by the unique demographic history of

95    the population.

96         One archetypal human population isolate with a demography that has been extensively

97    studied is the Finnish [2,11–13]. Finland was populated through two separate major migrations. The

98    first wave originated 4000 years ago from the west, and the second wave originated from the

99    southern shores approximately 2000 years ago. There was a subsequent internal migration and

100   expansion around the 16[th] century. Briefly, the small number of founders, relative isolation,

101   serial bottlenecks, and recent expansion in Finland has allowed drift to play a large role in

102   shaping the gene pool of this population. The aforementioned demographic history of Finland

103   has led to an increase in the prevalence of rare heritable Mendelian diseases, which has made this

104   population particularly fruitful for disease gene discovery studies[12,14]. Most of the initial disease

105   gene discovery studies in Finland exploited LD mapping in affected families and well curated

106   genealogical records to identify causal and candidate variants[12]. More recently, it has been

107   possible to apply population-based linkage analyses to identify disease associated variants as an

108   alternative to GWAS (unpublished data)[15] due to the availability of whole genome sequence data

109   in conjunction with extensive electronic health records.

110        A number of studies have shown that disease detection power can be improved by

111   studying population isolates other than the Finnish[10,16–18]. For example, the Greenlandic Inuit

112    (GI) experienced an extreme bottleneck which caused a depletion of rare variants and

113    segregating sites in their genome[18].The remaining segregating variants are maintained at higher

114    allele frequencies and a larger proportion of these SNPs are deleterious when compared to non-

115    isolated populations. Another study on Southeast Asian populations showed similar results.

116    Specifically, South East Asian populations have experienced more severe founder effects than

117    the Finnish[17], thus causing an excess of rare alleles associated with recessive disease. A study of

118    European population isolates compared the isolates with the closest non-isolated population from

119    similar geographic regions[10], and found that the total number of segregating sites was depleted

120    across all isolates relative to the comparison non-isolate. Of the sites that were segregating in

121    isolates, between ~30,000- 122,000 sites existed at an appreciable frequency (MAF > 5.6%),

122    while remaining rare (MAF < 1.4%) in all of the non-isolate population samples. These variants

123    could serve as candidate markers in genome-wide association studies (GWAS) for novel

124    associations and included SNPs that had been previously associated with cardio-metabolic traits.

125            As previously mentioned, population isolates tend to have a depletion of segregating sites

126    and an increased number of homozygous sites in their genomes. Pemberton and colleagues[19,20]

127    demonstrated that levels of homozygosity differ across the globe, and that long ROH, 2Mb or

128    greater, are strongly influenced by very recent demography. A study involving multiple Jewish

129    isolates showed a link between historic consanguinity and the amount of long ROH in the

130    genome[21].The same correlation between long ROH and parental consanguinity was also

131    observed in Middle Eastern isolates[4]. In another study, authors observed an increased count of

132    ROH and total length of the genome within an ROH in Greek isolates from the Pomak villages

133    and the island of Crete relative to a non-isolated Greek cohort[16]. Lastly, a study in Sardinians

134    showed there was sub-structure within the island when comparing the amount of long ROH

135    across sampled populations[22]. Authors identified an enrichment of long ROH in sub-populations

136    that had experienced recent endogamy, while sub-populations that had not experienced such

137    isolation have ROH levels similar to that of non-Sardinian Italians.

138         While there have been many studies of genetic variation in population isolates, the

139    studies described above have focused on populations where the founders all came from the same

140    ancestral population. However, the founders of Latin American population isolates have come

141    from distinct continental populations. We sampled individuals from mountainous regions of

142    Costa Rica and Colombia where geographic barriers resulted in populations remaining isolated

143    since their founding in the 16th and 17th centuries, until the mid-20th century[23]. Both groups share

144    a similar demographic history, having originated primarily from admixture between a few

145    hundred European males and Amerindian females, with a limited contribution from African

146    founders. After the founding event, both populations experienced a subsequent bottleneck and

147    then a recent expansion, within the last 300 years, the expansion increased the population size

148    over 1000-fold since the initial founding event[23]. The effect that admixture has had on overall

149    patterns of genetic variation in isolates remains elusive, and it is unclear whether these

150    populations share the typical genomic signatures seen in population isolates. While the small

151    founding population size could reduce diversity, because the Costa Rican and Colombian isolates

152    were founded from multiple diverse populations, they could potentially have increased in

153    diversity relative to other population isolates. Lastly, the impact of admixture on deleterious

154    variation remains unclear.

155         To better understand patterns of genetic variation in admixed isolated populations, we

156    compared the Colombian and Costa Rican population isolates to a benchmark isolate, the

157    Finnish, as well as other 1000 Genomes Project populations[24]. We observe that relative to the

158    Finnish, Latin American isolates have increased genetic diversity but an excess of IBD segments.

159    Moreover, we detect an increase in the proportion of an individual's genome that falls within a

160    long ROH in Latin American isolates relative to all other sampled populations and an enrichment

161    of deleterious variation within these long ROH. Demographic simulations indicate that the

162    enrichment of long ROH is a consequence of recent inbreeding in Latin American isolates. We

163    corroborate these results by leveraging extended pedigree data. Pedigree inbreeding coefficients

164    explain approximately 21.8% of the amount of an individual's genome within a long ROH. Next,

165    we examine the relationship between the proportion of European, Native American, and African

166    ancestry and the amount of the genome within an ROH, as well as the relationship to an

167    individual's pedigree inbreeding coefficient. To our knowledge, this is the first time long ROH

168    have been examined in admixed isolated populations. Further, we examine demography across

169    both recent and ancient timescales in these isolates. Our work sheds light on how the distinct

170    demographic histories of population isolates affect both genetic diversity and the distribution of

171    deleterious variation across the genome.

172

173    **Methods:**

174

175    *Pedigree Data for Costa Rican and Colombian Individuals*

176

177    Our study included 10 Costa Rican (CR) and 12 Colombian (CO) multi-generational pedigrees

178    ascertained to include individuals affected by Bipolar Disorder 1. More extensive details about

179    the curation of pedigree data and clinical assessments of diagnosis can be found in Fears et al.[25].

180

7

181    *Identifying Unrelated Individuals*

182

183    We defined unrelated individuals as those who are at most third-degree relatives. We chose this

184    threshold of relatedness because the families from CR and CO are known to be cryptically

185    related. We used KING[26] to identify 30 unrelated individuals from CR and CO. 24 of the 30

186    unrelated individuals in the CO are founders in the pedigree and 15 of the 30 unrelated

187    individuals in the CR are founders. The algorithm implemented in KING estimates familial

188    relationships by modeling the genetic distance between a pair of individuals as a function of

189    allele frequency and kinship coefficient, assuming that SNPs are in Hardy-Weinberg

190    equilibrium.

191

192    We also used KING[26] to identify 30 unrelated individuals from the following 1000 Genomes

193    Project[24] populations: Yoruba (YRI), CEPH-European (CEU), Finnish (FIN), Colombian

194    (CLM), Peruvian (PEL), Puerto Rican (PUR), and Mexican from Los Angeles (MXL). We used

195    these 30 unrelated individuals per population for all analyses unless otherwise stated

196    (**Supplementary Figure 1)**.

197

198    *Genotype Data Processing*

199

200    We generated a joint variant call file (VCF) containing single nucleotide polymorphisms (SNP)

201    from two separate data sets. The first data set contained 210 whole genome sequences sampled

202    from the aforementioned 1000 Genomes Project populations[24]. The second data set contained

203    449 whole genome sequences from Costa Rican and Colombian individuals. Variants in the

204     second data set were called following the GATK best practices pipeline[27] with the

205     HaplotypeCaller of GATK. All multi-allelic SNVs and variants that failed Variant Quality Score

206     Recalibration were removed. Genotypes with genotype quality score $\leq 20$ were set to missing.

207     Further quality control on variants was performed using a logistic regression model that was

208     trained to predict the probability of each variant having good or poor sequencing quality.

209     Individuals with poor sequencing quality and possible sample mix-ups were removed, and all

210     sequenced individuals had high genotype concordance rate between whole genome sequences

211     and genotypes from microarray data. All sequenced individuals had consistency between the

212     reported sex and sex determined from X chromosome and also between empirical estimates of

213     kinship and theoretical estimates. More information on sequencing and quality control

214     procedures is discussed in Sul et al. 2018 (unpublished data).

215          We used the following protocol to merge these two datasets. First, we used guidelines

216     from the 1000 Genomes Project strict mask to filter the Costa Rican and Colombian VCFs as

217     well as the 1000 Genomes Project VCFs. Then, we used GATK to remove sites from both sets of

218     VCFs that were not bi-allelic SNPs or monomorphic. Next, we merged the 1000 Genomes

219     Project VCFs with the Costa Rican and Colombian VCFs into a single joint-VCF for each

220     chromosome. We only used autosomes for our analyses. Lastly, we filtered the merged joint-

221     VCF to only contain sites that were present in at least 90% of individuals. There were a total of

222     57,597,196 SNPs and 1,891,453,144 monomorphic sites in the final data set. We ensured that the

223     merged data sets were comparable by examining the number of derived putatively neutral alleles

224     across the 30 unrelated individuals in all sampled populations, and finding few differences

225     between populations, which is consistent with theory[8] (**Supplementary Figure 2**).

226

9

227    *Calculating Genetic Diversity*

228

229    We computed two measures of genetic diversity from sites called across all 30 unrelated

230    individuals from each population: pi ($\pi$) and Watterson's Theta ($\theta_w$). The average number of

231    pairwise differences per site ($\pi$) was calculated across the genome as:

232    $$\pi = \frac{n}{n-1} \frac{\sum_{i=1}^{L} 2p_i(1-p_i)}{L},$$

233    where $n$ is the total number of chromosomes sampled, $p$ is the frequency of a given allele, and $L$

234    is the length in base pairs of the sampled region. Watterson's Theta, was computed by counting

235    the number of segregating sites and dividing by Watterson's constant, or the $n$-1 harmonic

236    number[28].

237

238    *Site Frequency Spectrum (SFS)*

239

240    Site frequency spectra were generated using the 30 unrelated individuals from each population,

241    SNPs with missing data were removed from these analyses. There was a total of 16 SNPs out of

242    the 57,597,196 SNPs that were removed due to missing data.

243

244    *Linkage Disequilibrium Decay*

245

246    We calculated LD between pairs of SNPs for all unrelated individuals. First, we applied a filter

247    to remove sites that were not at a frequency of at least 10% across all populations. Next, pairwise

248    $r^2$ values were calculated using VCFTools[29]. SNP pairs were then binned according to physical

249    distance (bp) between each other and $r^2$ was averaged within each bin.

10

250

251    *Identifying Identity by Descent Segments*

252

253    To detect regions of the genome that have shared IBD segments between pairs of individuals, we

254    first removed singleton SNPs in each population since singletons are not informative about

255    shared IBD. Then, we called IBD segments using IBDSeq[30]. IBDSeq is a likelihood-based

256    method that is designed to detect IBD segments in unphased sequence data. We chose to use

257    IBDSeq because other methods that require computational phasing could be biased when applied

258    to Latin American population isolates, as they do not have a publicly available reference

259    population to aid in phasing. We compared IBDSeq to two well-known phasing methods

260    Beagle[31] and GERMLINE[32] to determine whether it was feasible to use IBDSeq on an admixed

261    population (**Supplementary Figure 3**). Beagle[33] produced the shortest IBD segments while

262    GERMLINE produced the longest IBD segments. IBDSeq produced segments with a length

263    distribution similar to what we observed in Beagle, though the average segment length was

264    slightly larger, which we expected given that IBDSeq was created to call longer segments that

265    would have previously been broken up when using Beagle for phasing[30]. We used the default

266    parameters for IBDSeq.

267        Next, we filtered the pooled IBD segments to remove artifacts. First, we calculated the

268    physical distance spanned by each IBD segment. Then, we totaled the number of SNPs that fell

269    within each segment. We observed an appreciable number of IBD segments that were extremely

270    long but sparsely covered by SNPs (**Supplementary Figure 4**). IBD segments were removed if

271    the proportion of the IBD segment covered by SNPs was not within one standard deviation of the

272    mean proportion covered across all IBD segments (**Supplementary Figure 4**). Strong deviations

11

273   from the mean could indicate that the IBD segment spans a region of the genome with low

274   mappability, and we are only calling the SNPs at the outer ends of the segment. Therefore, the

275   true segment length might be much shorter than what is being calculated by IBDSeq. Lastly, we

276   converted from physical distance to genetic distance using the deCODE genetic map[34].

277

278   *Enrichment analyses of IBD segments*

279

280   To determine whether certain populations contain more IBD segments than others, we followed

281   the IBD score procedure outlined in Nakatsuka and collegues[17]. A population's IBD score was

282   calculated by computing the total length of all IBD segments between 3 and 20 cM. The score

283   difference is the difference between the query population's IBD score and the Finnish IBD score.

284   The score ratio is the ratio of each population's IBD score relative to the Finnish IBD score. The

285   significance of enrichment relative to the Finnish was evaluated using a permutation test for each

286   population, where IBD segment length was held fixed and labels of the two populations were

287   permuted. We recalculated the score on a total of 10,000 permutations to generate a null-

288   distribution of scores for each isolate.

289

290   *Estimating Effective Population Size*

291

292   We used the output files from IBDSeq to estimate the recent effective population size from the

293   30 unrelated individuals from each sampled population. We estimated effective population size

294   by using the default settings in IBDNe[35]. We set the minimal IBD segment length equal to 2cM

295   since that is the suggested setting when using sequence data.

12

296

297   *Identifying Runs of Homozygosity*

298

299   Runs of homozygosity were identified for each individual using VCFTools, which implements

300   the procedure from Auton et al. 2009[36].  Next, we examined the number of callable sites that lie

301   within each ROH. We found that there was a bi-modal distribution of coverage for ROH, where

302   some ROH appeared to contain almost no callable sites, while others had much higher coverage.

303   We only kept ROHs that were at least 2Mb in length, which we called long runs of

304   homozygosity, and were at least 60% covered by callable sites. (**Supplementary Figure 5)**.

305

306   *Calculating Inbreeding Coefficients*

307

308   SNP-based inbreeding coefficients were calculated using VCFTools[29]. VCFTools calculates the

309   inbreeding coefficient $F$ per individual using the equation $F = \frac{O-E}{N-E}$, where $O$ is the observed

310   number of homozygotes, $E$ is the expected number of homozygotes (given population allele

311   frequency), and $N$ is the total number of genotyped loci.

312          Pedigree-based inbreeding coefficients were computed using the R package kinship2[37].

313

314   *Demographic Simulations*

315

316   In order to investigate how aspects of the population history affect current day genetic diversity

317   in Latin American isolate populations, we simulated genetic variation data using the forward

318   simulation software SLiM 4.2.2[38]. We simulated a sequence length of 10Mb under uniform

13

319    recombination rate of $1\text{x}10^{-8}$ crossing-over events per chromosome per base position per

320    generation and under a mutation rate of $1.5\text{x}10^{-8}$ mutations per chromosome per base position

321    per generation. Every simulation contained intergenic, intronic, and exonic regions, but only

322    nonsynonymous new mutations experienced natural selection in accordance with the distribution

323    of selection coefficients estimated in Kim et al. 2017[39]. Within coding sequences, we set

324    nonsynonymous and synonymous mutations to occur at a ratio of 2.31:1[39,40]. The chromosomal

325    structure of each simulation was randomly generated, following the specification in the SLiM

326    4.2.2 manual (7.3), which is modeled after the distribution of intron and exon lengths in Deutsch

327    and Long[41].

328        We assumed an effective population size in the ancestral African population of 10,000

329    individuals, and a reduction in size to 2,000 individuals, starting 50,000 years ago, reflecting the

330    colonization of the European, Asian, and American continents. The population then recovers to a

331    size of 10,000 individuals 5,000 years ago. The colonization bottleneck is assumed to occur 500

332    years ago by an admixture event with a European population (70% admixture proportion) and is

333    followed by an immediate reduction in population size to 1,000 individuals. The recent

334    expansion in population size is modeled by an increase in population size to 10,000 individuals

335    200 years ago. We simulated data with recent inbreeding and without recent inbreeding. In the

336    former case, inbreeding started at the time of the European colonization 500 years ago and

337    continues until the present. Inbreeding is implemented with the "mateChoice" function in SLiM.

338    Because SLiM's pedigree track function is only valid for at most second-degree related

339    individuals, 50% of the time, mating occurs randomly. However, in the remaining cases, mating

340    occurs between close relatives with a relatedness coefficient bigger than 0.25. This produces

341    levels of consanguinity similar to those seen empirically as measured by F (see Results). Finally,

14

342    we sampled a total of 60 random individuals and calculated summary statistics on the sample

343    data. The simulation script can be found on GitHub (see Web Resources).

344

345    *Annotation of Variants*

346

347    The ancestral allele was determined using the 6-primate EPO alignment (see Web Resources)

348    and we restricted to only those sites called with the highest confidence. After filtering,

349    54,049,081 SNPs remained.

350

351    Subsequently, exonic SNPs were annotated using the SeattleSeq Annotation website (see Web

352    Resources). A total of 693,301 SNPs were annotated as either nonsynonymous or synonymous.

353    We further classified these sites as either putatively neutral or deleterious using Genomic

354    Evolutionary Rate Profiling (GERP) scores[42]. A GERP score less than two was considered as

355    putatively neutral and a GERP score greater than 4 was considered as putatively deleterious,

356    totaling 404,302 classified SNPs.

357

358    *Counting Deleterious Variants*

359

360    We used three different statistics to count the number of deleterious mutations per individual.

361    First, we tabulated the number of deleterious variants (the number of heterozygous plus the

362    number homozygous derived genotypes). Second, we counted the total number of derived

363    deleterious alleles (the number of heterozygous genotypes plus twice the number of homozygous

15

364    derived genotypes). Third, we computed the total number of derived deleterious homozygous

365    genotypes.

366

367    *Testing for an enrichment of deleterious variation in ROHs*

368

369    We were interested in whether there is an enrichment of nonsynonymous or loss-of-function

370    mutations in ROH over non-ROH regions for the three different ways of counting deleterious

371    variants outlined above. To account for differences in neutral variation, we standardized by

372    synonymous variation, which is assumed to be neutral. Then, we calculated the ratio of

373    nonsynonymous over synonymous variation in ROH regions divided by the ratio of

374    nonsynonymous over synonymous variation outside of ROH. We computed significance using a

375    permutation test, where the position of each SNP and its annotation as synonymous versus

376    nonsynonymous was fixed and the positions of the vector of ROH annotations were randomly

377    placed throughout the genome. Thus, the frequency distribution of synonymous and

378    nonsynonymous SNPs, as well as the total amount of ROH and non-ROH annotations, is kept

379    constant when compared to the unpermuted data. We recalculated the ratio for a total of 10,000

380    permutations to form a null-distribution of ratios and then computed significance.

381

382    *Calculating Ancestry Proportions*

383

384    We estimated genome-wide ancestry proportions in members of the CR and CO pedigrees using

385    ADMIXTURE[43] (v1.22). Then the genome-wide ancestry proportions from ADMIXTURE were

386    used as the prior in local ancestry analysis using LAMP[44]. We generated estimates for all 838

16

387    pedigree members with SNP array genotype data. Detailed information on the SNP array data

388    can be found in Pagani et al.[45]. The reference populations were the CEU (n=112) and YRI

389    (n=113) from HapMap[46,47], as well as 52 Native American samples from Central or South

390    America. The Native American samples are the Chibchan-speaking subset of those used in Reich

391    et al.[48], selected to originate from geographical regions relevant to CR/CO and to have virtually

392    no European or African admixture. In total the admixture analysis used 57,180 LD-pruned SNPs

393    and 1115 individuals.

394

395    *Accounting for Relatedness*

396

397    We tested for correlations among several quantities computed for each individual in the Latin

398    American population isolates. Because some of these individuals are closely related, the data

399    points in our linear regression are no longer independent. Therefore, we implemented the R-

400    package GenABEL[49] to incorporate kinship when performing statistical tests for our correlations.

401    We used the polygenic_hglm() function where the *formula* input was the equation for our linear

402    model of interest and the *kinship.matrix* input was a kinship matrix computed from our pedigree

403    computed using kinship2[37]. Our input took the following form: ($F_{PED}$ ~ Length of genome in

404    ROH, kin = kinshipMatrix, data = df).

405

406    **Results**

407

408    *Genetic Variation in Population Isolates*

409

410 We first compared levels of genetic diversity in a sample of 30 unrelated individuals across the

411 1000 Genomes populations[24] and the CO and CR isolates. We split the genome into several

412 different genomic regions and in each region summarized genetic variation using both the

413 average number of pairwise differences ($\pi$) and Watterson's theta ($\theta_w$) (**Figure 1A and B**).

414 Overall, we found differences in diversity across the functional category of sequence studied in

415 all populations, with coding regions exhibiting the lowest diversity and intergenic regions the

416 highest. These patterns are consistent with the role of purifying selection affecting coding

417 diversity[39]. However, if we look genome-wide or focus on intronic regions we see intermediate

418 levels of diversity (**Supplementary Table 1 and 2)**. We suspect that these categories are more

419 strongly influenced by linked selection[50–52].

420 As we are interested in the role of demography in shaping genetic diversity, we focused

421 on comparisons of intergenic levels of diversity as those are most likely to be neutrally evolving

422 (**Figure 1A and B**). Overall, the YRI had the highest level of diversity ($\pi \approx 0.0010$; $\theta_w \approx 0.0012$)

423 (**Supplementary Table 1 and 2**). The European populations (CEU and FIN) had lower levels of

424 diversity. The CEU and FIN had similar levels of $\pi$ (approximately 0.0004), despite the FIN

425 being considered an isolated population. However, the FIN had reduced numbers of SNPs as

426 reflected by lower values of $\theta_w$ (CEU $\approx$ 0.0008 & FIN $\approx$ 0.0008). The CO and CR had levels of

427 diversity similar to that of several other Latin American populations in the 1000 Genomes

428 Project (CLM and MXL). We found no clear pattern of the population isolates (FIN, CO, CR)

429 having lower diversity than their most similar non-isolated population. Instead, diversity levels

430 tended to be higher across all the sampled Latin American populations (CLM, CO, CR, MXL,

431 and PUR) when compared to the European populations. One exception to this pattern is the PEL

432 population, who had the lowest neutral levels of diversity ($\pi \approx 0.0007$; $\theta_w \approx 0.0007$).

18

433     Next, we examined the proportional site frequency spectrum (SFS; **Figure 1C**). Latin

434     American populations had the highest proportion of singletons. The CO and CR had similar

435     proportions of singletons when compared to other 1000 Genomes Project Latin American

436     populations. Conversely, the FIN had the lowest proportion of singletons in comparison to all

437     sampled populations. The depletion of singletons relative to common variation supports the

438     presence of a stronger founder effect during the FIN population history[13].

439     We also examined patterns of linkage disequilibrium (LD), since LD is affected by

440     population size and recent bottlenecks[53,54]. **Figure 1D** shows the mean decay of $r^2$ with physical

441     distance over 2Mb intervals across the genome in each population. We found that the YRI had

442     the lowest levels of LD for each bin of physical distance, and the PEL formed the upper bound of

443     the LD decay curves. The remaining Latin American populations (PUR, MXL, CLM, CO, CR)

444     clustered together, close to the YRI, while the CEU and FIN are shifted toward higher values,

445     like those seen in the PEL.

446     The FIN were previously shown to have more extensive haplotype blocks in their

447     genome in comparison to the Latin American isolates[9]. In line with these findings, we observed

448     faster LD decay in the Latin American isolates relative to the FIN. When considering pairs of

449     SNPs 150kb or more apart, rates of LD decay become quite similar across all the sampled

450     populations. Analogous to other diversity statistics, LD in the CO and CR closely resembled

451     those of non-isolated Latin American populations. Once again, we found there is no clear pattern

452     of having lower diversity or more LD that holds across all the population isolates (FIN, CO, CR)

453     when compared to their most similar non-isolated population.

454

455     *Latin American isolates carry more IBD segments than Finnish*

19

456

457    Next, we used IBD sharing between pairs of individuals to gain insight about more recent

458    demographic events within populations (**Figure 2**). We compared the amount of IBD within each

459    population by computing an IBD score. Each population's IBD score was calculated by totaling

460    the length of IBD segments between 3cM and 20cM. We expressed IBD scores for each

461    population as the ratio of the IBD score for a given population relative to the IBD score in the

462    FIN (**Figure 2A**). We also tabulated the total count of IBD segments for each population. The

463    CEU showed the lowest number of both called IBD segments and the lowest IBD score relative

464    to the FIN (p-value = 0.0001). Latin American populations formed the upper bounds of both total

465    IBD segments called and IBD enrichment scores (**Figure 2A**). The PUR had the largest number

466    of IBD segments (1402) and had a 2.1-fold increase in IBD score relative to the FIN (p-value <

467    $1x10^{-4}$). The CO and CR isolates had a 1.8-fold and 2-fold increase in their IBD scores relative to

468    the FIN (p-value < $1x10^{-4}$), as well as carrying more IBD segments than the FIN (**Figure 2B and**

469    **2C**). However, there were some Latin American populations that exhibited depletions in both

470    IBD segments and IBD scores relative to the FIN. The MXL and PEL have the lowest number of

471    IBD segments for the Latin American populations. Previous work has shown that a larger

472    effective population size in admixed populations likely drove the depletion of IBD segments in

473    these two Latin American populations[55].

474

475    *Inferring the Demographic History of Latin American Isolates*

476

477    We next leveraged the patterns of IBD described above to estimate the effective population size

478    using IBDNe[35] on the 30 unrelated individuals from each population (**Figure 3**). The use of only

479    30 unrelated individuals caused limitations for accurate estimation of $N_e$ (see Discussion), but

480    the demographic history of the population was robust to the number of individuals used. First,

481    we found that recent demography differs vastly between the European populations (FIN and

482    CEU). In general, CEU experienced population expansions over much of their demographic

483    history. It was only in the most recent generations that they experienced a decrease in $N_e$. The

484    FIN, on the other hand, have experienced a long population decline since their founding,

485    approximately 4000 years ago, followed by a recent population expansion.

486        When analyzing the Latin American isolates, we detected a recent bottleneck,

487    approximately 500 years ago (**Figure 3**). This bottleneck could correspond to the recorded

488    bottleneck that followed the founding of these populations, and it appears to be much shorter and

489    less severe than the bottleneck seen in the FIN. The strength and duration of bottlenecks varied

490    across each of the Latin American populations. For example, we observed a more severe

491    bottleneck in the CR, CO, CLM, and PUR than in PEL or MXL. However, we detected a

492    subsequent period of growth across all populations following the bottleneck. The rate of growth

493    differed across each population, and the PEL appeared to be growing at a much more rapid rate

494    than any of the other Latin American populations.

495

496    *Exploring Recent Consanguinity*

497

498    Isolated populations may have experienced recent consanguinity. To test for this, we began by

499    examining SNP-based inbreeding coefficients ($F_{SNP}$) (**Supplementary Figure 6**). YRI

500    individuals had the lowest median inbreeding coefficients and the CO and CR isolates had the

501    highest median inbreeding coefficients. Further, the CO and CR also had the highest maximum

21

502     $F_{SNP}$ values in the entire sample of unrelated individuals from any population (**Supplementary**

503     **Figure 6**). Median levels of $F_{SNP}$ in the CEU suggested that they are more inbred than the FIN,

504     which may be a result of how 1000 Genomes samples were selected. The PEL had the largest

505     variance in $F_{SNP}$ across any of the sampled populations.

506          Next, we examined patterns of long runs (>2Mb, see Methods) of homozygosity (ROH),

507     since ROH have been linked to recent consanguinity[20,21,56–58]. The YRI and CEU had the lowest

508     amount of their genome contained within an ROH (**Figure 4A**). The FIN had higher median

509     amounts of their genome within an ROH in comparison to the CEU. Latin American isolates had

510     the highest median amount of the genome contained within an ROH. Specifically, the CR had

511     the highest median at 21.7 Mb. Further, the Latin American isolates also had the greatest

512     variance in the amount of the genome contained within an ROH. For example, one of the CO

513     individuals had approximately 230 Mb of her/his genome contained in long ROHs.

514          As expected, we found that the amount of the genome contained in a long ROH strongly

515     correlated with an individual's $F_{SNP}$ (CO: $R^2 = 0.8060$, p –value = 1.1 x $10^{-11}$; CR: $R^2 = 0.7740$,

516     p-value = 9.5 x $10^{-11}$; FIN: $R^2 = 0.1288$, p-value = 0.03) (**Figure 4B-4D**). Indeed, individuals

517     with higher values of $F_{SNP}$ tended to have more of their genome within an ROH. Further, the

518     individual with the highest $F_{SNP}$ (0.133) also had the largest amount of his/her genome in long

519     ROHs (230Mb).

520          The total number of ROH segments per individual followed a similar pattern as the total

521     amount of genome within an ROH (**Supplementary Figure 7**). For example, in populations with

522     low values of $F_{SNP}$, ROH segments were not frequent. One YRI individual and three CEU

523     individuals carried a ROH >4Mb, whereas more than 50% of CO and CR individuals carried an

22

524    ROH >4Mb. Additionally, the longest ROHs identified (>20MB), only occurred in Latin

525    American populations, where there were the largest values of $F_{SNP}$ (**Supplementary Figure 7**).

526          Importantly, the FIN individuals had significantly fewer ROH segments than the CO and

527    CR, and most individuals had an $F_{SNP}$ close to 0 (**Figure 4**). The Latin American isolates had the

528    most ROH in comparison to any other sampled population, as well as the largest values of $F_{SNP}$

529    (**Figure 4**).

530

531    *Determining the Mechanisms that Generate Runs of Homozygosity*

532

533    In principle, ROHs can be generated either by recent consanguinity over the last few generations,

534    or by older historical processes, such as bottlenecks[19,56,58–60]. Based on both historical data[23] and

535    inference from IBDNe analyses, Latin American population isolates show evidence of recent

536    population bottlenecks. Therefore, we used two complementary strategies to test whether recent

537    consanguinity or bottlenecks drove the observed increase in ROHs in the Latin American

538    isolates. First, we used the extensive pedigree data for 449 sequenced individuals to calculate a

539    pedigree inbreeding coefficient ($F_{PED}$) for each individual (**Figure 5**). Most individuals had a

540    $F_{PED}$ of 0. However, there were several individuals with values of $F_{PED}$ as high as 0.07 in CR and

541    0.06 in CO. We observed a significant correlation between $F_{SNP}$ and $F_{PED}$ ($R^2$=0.1520 and p-

542    value < $2 \times 10^{-16}$), even after accounting for the non-independence of individuals based on their

543    kinship (**Figure 5A**; see Methods). These correlations suggest that the recent consanguinity

544    captured within the last few generations in the pedigree was likely sufficient to drive the increase

545    in ROHs in the CO and CR populations. $F_{SNP}$ was a substantially better predictor of the amount

546    of an individual's genome that falls within a ROH ($R^2 = 0.7540$ and p-value < $2 \times 10^{-16}$), than

23

547     $F_{PED}$ ($R^2$ = 0.2180 and p-value < 2 x $10^{-16}$ ) (**Figure 5B and C**) likely due to the fact that $F_{SNP}$

548     captured distant background relatedness within the population as well as the realized level of

549     consanguinity, rather than the expected value[61]. Further, because the pedigrees were ascertained

550     and analyzed separately, connections between pedigrees were not accounted for in $F_{PED}$, but were

551     likely captured by $F_{SNP.}$

552          As a second approach to determine the mechanism driving the increase in ROHs in the

553     CO and CR populations, we conducted forward in time demographic simulations. We simulated

554     a 10Mb region under a demographic model that reflected changes in effective population size

555     during the human expansion across the European, Asian and American continents, as well as the

556     more recent bottleneck during the Spanish colonization about 500 years ago (**Figure 5D**; see

557     Methods). We compared simulations with no inbreeding to simulations with recent inbreeding.

558     Consanguinity in the populations was modeled to begin 500 years ago, and simulated individuals

559     had an inbreeding coefficient of about 0.075. This level of inbreeding was comparable to the

560     level of inbreeding in some of the CO and CR individuals, based on calculations using pedigree

561     data (see Methods). Our simulations suggested that the recent population bottleneck caused by

562     the Spanish colonization was not capable of generating the large amounts of the genome within

563     an ROH (>2Mb) that we observed for some of the individuals (**Figure 5D**). Only when

564     simulating recent inbreeding could levels of the genome in an ROH comparable to that we

565     observed be generated. Thus, recent inbreeding was paramount for generating the long ROH that

566     we observed in the CO and CR isolates.

567

568     *Local Ancestry*

569

24

570    Since the Latin American isolates originated from an admixture event between Native

571    Americans, Africans, and Europeans, we tested for a correlation between $F_{PED}$ and the proportion

572    of European, African, or Native American ancestry (**Supplementary Figure 8**). We used the

573    entire sequenced Costa Rican and Colombian data set (n=449) for the local ancestry analyses and

574    accounted for relatedness of individuals in all the following reported p-values (see Methods). We

575    found that European ancestry was positively correlated with $F_{PED}$ (p-value = 0.0052) while

576    Native American ancestry was negatively correlated with $F_{PED}$ (p-value = 0.0245). African

577    ancestry was also negatively correlated with $F_{PED}$ (p-value = 0.0496).

578            Then, we asked if the proportion of ancestry was correlated with the amount of the

579    genome within an ROH (**Supplementary Figure 8**). The correlation between ancestry and

580    amount of the genome within an ROH followed the same trend as the correlation between

581    ancestry and $F_{PED}$. Native American ancestry and African ancestry are negatively correlated with

582    the amount of the genome within a long ROH (p-value = 3.91 x $10^{-14}$ and p-value = 6.76 x $10^{-07}$,

583    respectively). European ancestry was positively correlated with the amount of an individual's

584    genome within an ROH (p-value < 2 x $10^{-16}$).

585

586    *Recent Consanguinity is Correlated with an Increase of Deleterious Variation*

587

588    It is well known that demography impacts patterns of deleterious variation in populations[5,8,20,62–

589    66]. Thus, we compared patterns of putatively deleterious variation in the CO and CR to those in

590    the FIN. Variants were classified as putatively deleterious or putatively neutral using GERP

591    scores (see Methods). Recall that we consider three ways of counting deleterious variants in the

592    genome of an individual: first, counting the number of heterozygous genotypes plus twice the

25

593    number of homozygous derived genotypes (i.e. the total number of derived deleterious alleles),

594    second, counting the number of heterozygous and homozygous derived genotypes (counting

595    variants), and third, counting only the number of homozygote derived genotypes (counting

596    homozygotes). The first quantity is most relevant if deleterious alleles are additive, while the

597    third is most relevant if they are recessive. First, we looked at absolute counts of derived

598    deleterious variation across isolates (**Supplementary Figure 9**). Then, we used linear regression

599    to test if there was a relationship between the amount of an individuals' genome in an ROH and

600    the number of nonsynonymous sites in the genome for each counting method (**Figure 6**).

601        The FIN carried approximately 1% more derived deleterious nonsynonymous alleles per

602    individual than CO and CR (p-value = 0.0007; p-value = 0.0013). However, there was no

603    significant difference in the number of putatively neutral synonymous derived alleles per

604    individual. These results suggest that the difference seen for putatively deleterious variants is not

605    driven by data artifacts (**Supplementary Figure 9**), and the FIN indeed have a slightly higher

606    additive genetic load than the CO or CR. Turning to the number of variants per individual, FIN

607    individuals carried significantly more deleterious nonsynonymous variants than the CR but not

608    the CO (p-value = 0.0110). However, CO and CR did not differ significantly in the number of

609    deleterious variants carried per individual (**Supplementary Figure 9**). When we examined

610    neutral synonymous variants, CO had significantly more variants than either FIN or CR (p-value

611    = $8.56^{-06}$; p-value = 0.0054, respectively). Finally, when counting the number of homozygous

612    derived genotypes, we found that the FIN carried 3.3% more deleterious variants in the

613    homozygous state per individual than CO but not the CR (p-value = 0.0003) (**Supplementary**

614    **Figure 9**). Additionally, the FIN carried significantly more neutral homozygous genotypes per

615    individual than either population (CO p-value = $1.01 \times 10^{-05}$; CR p-value = $6.96 \times 10^{-05}$). The

26

616    increased deleterious and neutral variation in homozygous form is an expected consequence of

617    the long-term bottleneck that the FIN experienced during their founding.

618         We next tested whether the amount of the genome in an individual contained within a

619    ROH was correlated with the number of nonsynonymous mutations carried by the individual.

620    Counting nonsynonymous (NS) or synonymous (SYN) allele copies did not show any correlation

621    with the amount of an individuals' genome that falls within an ROH for the CR or FIN (**Figure**

622    **6A and D; Supplementary Figures 10-13**). However, in the CO, as the amount of the genome

623    within an ROH increased, individuals tended to carry more NS alleles, though this correlation

624    was strongly driven by a single individual, who also had the highest $F_{SNP}$ and $F_{PED}$ ($R^2 = 0.2393$;

625    p-value = 0.0036; **Supplementary Figure 10**). Importantly, the number of SYN alleles per

626    individual was not correlated with the amount of the genome in an ROH (p-value = 0.2261).

627         When counting variants per individual, we observed a significant negative correlation

628    with the amount of an individuals' genome that falls within an ROH in the Latin American

629    isolates (**Figure 6B and E; Supplementary Figures 10-12**). The negative correlation is a result

630    of heterozygous sites being lost when an ROH is formed due to inbreeding. Conversely, when

631    counting homozygous genotypes per individual, we observed a significant positive correlation

632    with the amount of an individual's genome that falls within an ROH in both the Latin American

633    isolates and FIN (**Figure 6C and F; Supplementary Figures 10-13**). Homozygous genotypes

634    were the only statistic that correlated significantly with the amount of the genome in an ROH

635    across all isolated population for both SYN and NS sites. We observed a stronger correlation

636    between the number of NS homozygous genotypes and the amount of an individual's genome

637    within an ROH in the Latin American isolates ($R^2 = 0.5000$ (CO) & $R^2 = 0.2165$ (CR); p-value =

638    $7.546^{-06}$(CO) and p-value = 0.0059(CR)) compared to the FIN ($R^2 = 0.1130$ and p-value =

27

639 0.0389) (**Supplementary Figures 10-13)**. This pattern exists because the majority of CO and CR

640 individuals carried larger proportion of their genome within an ROH while the FIN individuals

641 do not harbor many ROH.

642  We next asked whether there was an enrichment or depletion of NS variants relative to

643 SYN variants within versus outside of an ROH using a permutation test on the three different

644 counting approaches (see Methods). When variants or allele copies were counted, none of the

645 populations produced significant results (**Table 1**). When homozygous genotypes were counted,

646 ROHs in the MXL and CR were enriched for homozygous NS genotypes relative to SYN

647 homozygous genotypes (p-value = 0.0052 and p-value = 0.0169) (**Table 1**). Additionally, if we

648 pooled the CR and CO populations, we also observed a significant enrichment of deleterious

649 variation within an ROH compared to non-ROH regions of the genome (p-value = 0.0011).

650  We tested whether $F_{SNP}$ was correlated with the amount of deleterious variation per

651 individual. We only used isolates for these regressions, because we are particularly interested in

652 how recent consanguinity affected deleterious variation in the genome. We observed the exact

653 same pattern with $F_{SNP}$ as with ROH (**Supplementary Figure 14**). Briefly, counting NS or SYN

654 allele copies did not show any correlation with $F_{SNP}$ for the CR or FIN, but there was a

655 significant correlation with NS allele copies in CO (**Supplementary Figure 14; Supplementary**

656 **Figures 15-17**). Counting NS and SYN variants per individual produced a significant negative

657 correlation with $F_{SNP}$ in the Latin American isolates (**Supplementary Figure 12;**

658 **Supplementary Figures 15-17**). Counting the number of NS and SYN homozygous genotypes

659 per individual was positively correlated with $F_{SNP}$ in the both Latin American isolates and FIN

660 (**Supplementary Figure 14; Supplementary Figures 15-17**). Again, counting homozygotes

661 was the only method with significant results across all isolated populations for both SYN and NS

662    variants. The ability to recapitulate the pattern we observed in ROH using $F_{SNP}$ was reassuring

663    and adds further support to the strong relationship between recent consanguinity and ROH.

664         Lastly, because we had multi-generational pedigrees for the Latin American isolates, we

665    examined the correlation between putatively deleterious variation and recent consanguinity as

666    measured by ($F_{PED}$). All the following reported p-values account for kinship (see Methods).

667    When we pooled the CO and CR individuals together, we did not observe any relationship

668    between counting derived deleterious allele copies and $F_{PED}$ after correcting for kinship (**Figure**

669    **7A**). Moreover, we observed a negative correlation between $F_{PED}$ and the number of deleterious

670    variants per individual ($R^2 = 0.0375$, p-value $= 6.02$ x $10^{-06}$). The number of neutral variants per

671    individual was also negatively correlated with $F_{PED}$ (p-value $= 2.26$ x $10^{-10}$) (**Figure 7B**). Finally,

672    we observed a positive correlation between $F_{PED}$ and derived deleterious homozygotes ($R^2 =$

673    $0.0575$, p-value $= 1.0$ x $10^{-06}$) as well as between $F_{PED}$ and the number of neutral derived

674    homozygotes per individual. (p-value $= 1.03$ x $10^{-08}$) (**Figure 7C**). These results suggest that

675    recent consanguinity during the last few generations has increased the number of derived

676    deleterious homozygous genotypes in these two populations.

677

678    **Discussion:**

679

680    Here we present the first comprehensive study of genetic diversity, demographic history,

681    identity-by-descent, runs of homozygosity, and deleterious mutations in multiple admixed

682    isolated populations. We show that admixture sufficiently increases genetic diversity of the

683    Colombian and Costa Rican isolates, such that each isolate has diversity levels comparable to a

684    non-isolated population. However, we still observe characteristics in the Latin American isolates

29

685    that are hallmarks of an archetypal isolate, such as: an excess of IBD segments, cryptic

686    relatedness within the population, and an enrichment of long ROH. Further, we demonstrate that

687    long ROHs contain an enrichment of deleterious variants carried in the homozygous state, which

688    has potential implications for fitness and disease risk.

689         Taken together our results support historical data which states that a recent admixture

690    event, within the last 500 years, founded the Colombian and Costa Rican population isolates. A

691    bottleneck corresponding to the Spanish Settlement, followed the founding event and then each

692    population has increased in size until the present day[23,67]. We see evidence of these processes in

693    the inference of demography from IBD patterns. The IBDNe inference shows a severe

694    population bottleneck occurring approximately 500 years ago, coinciding with the historical

695    record[23]. Importantly, the bottleneck experienced in the Latin American isolates was not as

696    prolonged as that experienced by the Finnish. Further, the Finnish bottleneck occurred thousands

697    of years ago. The difference in bottleneck timescales likely accounts for some portion of the

698    higher genetic diversity observed in Latin American population isolates in comparison to the

699    Finnish. In other words, the bottlenecks captured by IBDNe in the Latin Americans are too

700    recent to markedly impact levels of genetic diversity. Further, the admixture process experienced

701    by the Latin American isolates could increase levels of genetic diversity, especially because

702    some individuals have appreciable levels of African ancestry[50]. We see little difference in

703    patterns of genetic variation in the 1000 Genomes Colombian samples and the Colombian

704    sample studied in this project. The Latin American isolates occupy areas that were considered as

705    being geographically isolated at the time of sampling -- the Central Valley of Costa Rica and the

706    department of Antioquia in Colombia[23], while the 1000 Genomes CLM sample was taken from

707    Medellín which is included within the Antioquia region[68–71]. Thus, there is likely some amount

30

708    of shared demography between the 1000 Genome CLM and our isolated Colombia population.

709    However, it is worth noting that the individuals in the Latin American isolates are from pedigrees

710    ascertained for Bipolar Disorder 1, rather than a random sample from the area.

711        Our results beg the question, what constitutes a population isolate? For example, is it a

712    requirement that population isolates have low genetic diversity relative to the source population?

713    Under this definition, the Latin American population isolates would not qualify as population

714    isolates. The bottleneck in the Costa Ricans and Colombians seems to have had little effect on

715    their genetic diversity, as their diversity levels are similar to non-isolated Latin American

716    populations. The Finnish, on the other hand, experienced a long-term bottleneck that has resulted

717    in a depletion of segregating sites, and of the remaining segregating sites, there is an enrichment

718    of deleterious variants relative to non-isolated populations[2,13], and would clearly qualify as an

719    isolate. However, if one measures isolation based on IBD we see that there is an enrichment of

720    IBD segments in the Latin American isolates relative to the Finnish. Further, looking at ROH,

721    Latin American individuals from population isolates have a larger burden of ROH than Finnish,

722    thus increasing the chances of identifying more shared genomic regions in the Latin American

723    isolates than the Finnish. By this metric, the Latin American population isolates would certainly

724    qualify as a population isolate. Thus, both the Costa Rican and Colombian populations and the

725    Finnish are isolates but in different ways. For example, the Costa Ricans and Colombians are

726    historical isolates, meaning these populations are not currently isolated but they exhibit many

727    traits of an isolate, whereas the Finnish are contemporary isolates, meaning the population is still

728    isolated and is the archetypal isolate that one would imagine. Our work also suggests that

729    isolated populations have distinct demographic histories that impact genetic variation in different

730    ways; and it is critical that researchers study and quantify the consequences of demography in

731    each population.

732        We find that Latin American isolates have the largest ROH burden in comparison to any

733    other sampled population. Our work corroborate results from a recent review on ROH where

734    authors state that populations with small $N_e$ and recent consanguinity will harbor the largest

735    amount of ROH[72]. Because previous research has shown a strong correlation between recent

736    inbreeding, quantified by both $F_{SNP}$ and $F_{PED}$, and long runs of homozygosity, we were

737    particularly interested in the mechanism behind the generation of long ROH[19–21,56–58,73]. We used

738    simulations to test which demographic scenarios could produce long ROH (**Figure 5**). These

739    simulations and availability of extended pedigree data were crucial, because the $F_{SNP}$ metric can

740    also be influenced by a recent bottleneck. Thus, having $F_{PED}$ available allowed us to test whether

741    the correlation we observed between $F_{SNP}$ and ROH was a consequence of a recent bottleneck or

742    recent consanguinity. If small population size or admixture was responsible for generating the

743    ROHs, these processes would not be reflected in $F_{PED}$. Thus, we would not expect to find a

744    correlation between $F_{PED}$ and the amount of the genome in ROHs. The fact that we observe a

745    correlation between $F_{PED}$ and the amount of the genome in ROH suggests that recent

746    consanguinity (as measured by $F_{PED}$) is related to the extent of long ROHs in the genome.

747    Further, our simulations show that neither admixture nor a recent population bottleneck could

748    generate the high levels of long ROH that are observed in some individuals. It was only when we

749    incorporated inbreeding into the simulation that levels of ROH comparable to what we observed

750    in our data were produced. Thus, both lines of evidence suggest that the Latin American

751    population isolates have experienced more recent consanguinity than other population isolates,

752    like the Finnish. Further, in Finland it has previously been shown that the frequency of

753    consanguinity, due to first-cousin marriages, is quite low and the best predictors of these unions

754    were socio-economic class and ethnicity, rather than geographic barriers or population density[74].

755    On the other hand, for the two Latin American isolates consanguinity could be a consequence of

756    increased geographic barriers preventing movement of individuals over more dispersed areas. It

757    is also important to point out that it is unclear the extent to which ascertaining individuals from

758    large pedigrees may impact the number of ROHs in our sample. Thus, the finding of an increase

759    in ROHs may not be generalizable to Colombian and Costa Rican populations as a whole.

760    However, we observed a similar pattern of increased ROH in the CLM, which suggests that the

761    pedigree ascertainment of the CO and CR may not be generating the increase in ROHs.

762        We also tested how recent consanguinity affects deleterious variation in the genome.

763    When counting homozygous derived deleterious genotypes, we found a positive correlation

764    between the number of nonsynonymous homozygous genotypes and the amount of an

765    individual's genome within an ROH (**Figure 6**). Further, we observed an enrichment for

766    nonsynonymous homozygous derived genotypes relative to synonymous homozygous derived

767    genotypes within ROHs versus the rest of the genome (**Table 1**). This enrichment can be a result

768    of nonsynonymous mutations generally segregating at lower frequency and typically being

769    carried as a single copy in an individual. When an ROH is formed, the chromosome that was

770    carrying the mutation is copied, thus allowing the mutation to increase the number of

771    homozygotes within the ROH[19,20]. Since long ROH are a product of recent consanguinity, and

772    these populations have experienced recent consanguinity, we see a corresponding increase in the

773    burden of deleterious variants in the genomes of Costa Rican and Colombian isolates. Since we

774    are more likely to see deleterious variants in the homozygous form in areas of the genome that

775    fall within an ROH, our work is particularly relevant for alleles associated with recessive

776    diseases. Lastly, we provide a mechanism for how recent consanguinity can reduce fitness in

33

777   natural populations[75–77]. Specifically, if gene-knockouts and deleterious mutations tend to be

778   recessive[40,78–82], as suggested by several studies, then recent consanguinity will increase the

779   number of homozygous derived deleterious variants carried by an individual in a long ROH, thus

780   leading to a reduction of fitness in the sampled population[6].

781         Utilizing estimated ancestry proportions from across the genome, we tested for a

782   correlation between an individual's ancestry and the amount of their genome that falls within an

783   ROH. To our knowledge, this is the first time that the relationship between proportion of

784   ancestry and the amount of the genome within an ROH has been examined (**Supplementary**

785   **Figure 8**). We found a positive correlation between the proportion of European ancestry and the

786   amount of an individual's genome within a run of homozygosity. These results are consistent

787   with the Latin American isolates originating from a small number of European founders, which

788   would decrease genetic diversity and increase homozygosity for those areas of the genome

789   containing European haplotypes. We observed a negative correlation between Native American

790   ancestry and the amount of the genome contained within an ROH (**Supplementary Figure 8**).

791   This finding appears to be at odds with previous research[19,72] that detected the opposite pattern.

792   Some of this difference may be due to distinct sampling strategies of the Native American source

793   population in our study compared to previous work. The reference Native American population

794   we used was composed of Chibchan-speaking individuals from Reich et al.[48]. Chibchan-speaking

795   populations inherited their Native American ancestry from admixture between Southern and

796   Northern American lineages, the necessity of  admixture was particularly apparent in the Cabecar

797   of Costa Rica[48]. Because our reference Native American population is admixed, and Native

798   American populations tend to be small, it is likely that drift has affected different alleles in

799   source populations that formed the current Chibchan-speaking populations. The Chibchan-

800    speaking populations may have more diversity, fewer fixed homozygous sites, than previously

801    sampled Native American populations which could explain the negative correlation we observed

802    between ancestry and ROH.

803         While we were able to capture evidence of recent bottlenecks and expansions within

804    Latin American isolates using IBDNe[35] (**Figure 3**), our demographic inferences have some

805    limitations. For example, the current estimates of $N_e$ are unrealistically large. This inflation may

806    be due to low sample size, since we only used 30 individuals, or it may be a result of applying

807    IBDNe to admixed populations. IBDNe was designed to be applied to un-admixed, randomly

808    mating populations. Thus, one area of future research could be exploring the influence of

809    admixture on IBD and developing methods to infer demography using IBD patterns in admixed

810    populations. Interestingly, in our study, the populations with the highest IBD scores were

811    admixed (PUR, CO, CR, and CLM). Furthermore, because IBD segments may contain useful

812    information for identifying regions of the genome that contain disease associated mutations,

813    especially within individuals with the highest amounts of consanguinity, it may be useful to

814    deconvolute ancestry for each segment when identifying disease associated mutations because

815    disease prevalence may differ in each parental population.

816         Population isolates have frequently been used for mapping Mendelian disease genes[17,83–

817    88] and studying complex diseases[16,89–95]. Isolates are thought to be beneficial in comparison to

818    non-isolated populations because of their increased homogeneity of the gene pool, disease

819    causing alleles potentially existing at an appreciable frequency due to drift, possible enrichment

820    in prevalence of the phenotype of interest[1–4], and a likely reduction in the variability of

821    phenotypes. Our work shows that the genetic diversity and genomic background of population

822    isolates varies immensely. Therefore, it is imperative that we understand the unique genetic

823     diversity belonging to each population isolate. Researchers should adapt their study design to

824     integrate the demographic history of the population, to better leverage the power of the unique

825     genetic features of the population of interest. For example, if we knew beforehand that there was

826     a history of consanguineous unions within the study population, researchers could target ROH

827     for disease mapping. This method has previously been used to identify human knockouts,

828     discover novel loci associated with disease, and understand gene function[95–98]. Further,

829     populations with large amounts of ROH could help us better understand disease architecture

830     because ROH may harbor more recessive mutations that do not have full penetrance, since the

831     prevalence of recessively acting variants in ROH is enriched relative to non-ROH portions of the

832     genome. Most importantly, our work highlights the importance of understanding the

833     demographic history of isolated populations, as differences in demographic history will greatly

834     impact patterns of genetic variation in isolates.

835
836

**Supplemental Data:**

Supplemental data include seventeen figures and two tables.

**Conflicts of Interest:**

**Web Resources**:

6-primate EPO alignment: ftp://ftp.ensembl.org/pub/release-75/fasta/ancestral_alleles/
ADMIXTURE: https://www.genetics.ucla.edu/software/admixture/download.html
IBDSeq: http://faculty.washington.edu/browning/ibdseq.html
IBDNe: http://faculty.washington.edu/browning/ibdne.html#download
KING (version 2.1): http://people.virginia.edu/~wc9c/KING/history.htm
LAMP: http://lamp.icsi.berkeley.edu/lamp/
PLINK: http://www.cog-genomics.org/plink2
ROH simulation script: https://github.com/LohmuellerLab/ROH_Latin_American_Isolates
SeattleSeq Annotation website: http://snp.gs.washington.edu/SeattleSeqAnnotation138/
SLIM: https://messerlab.org/slim/
VCFTools: http://vcftools.sourceforge.net/downloads.html
GATK: https://software.broadinstitute.org/gatk/download/archive

**References:**

1. Peltonen, L., Palotie, A., and Lange, K. (2000). Use of population isolates for mapping complex traits. Nat. Rev. Genet. *1*, 182–190.

2. Lim, E.T., Würtz, P., Havulinna, A.S., Palta, P., Tukiainen, T., Rehnström, K., Esko, T., Mägi, R., Inouye, M., Lappalainen, T., et al. (2014). Distribution and medical impact of loss-of-function variants in the Finnish founder population. PLoS Genet *10*, e1004494.

3. Gudbjartsson, D.F., Helgason, H., Gudjonsson, S.A., Zink, F., Oddson, A., Gylfason, A., Besenbacher, S., Magnusson, G., Halldorsson, B.V., and Hjartarson, E. (2015). Large-scale whole-genome sequencing of the Icelandic population. Nat. Genet. *47*, 435.

4. Scott, E.M., Halees, A., Itan, Y., Spencer, E.G., He, Y., Azab, M.A., Gabriel, S.B., Belkadi, A., Boisson, B., Abel, L., et al. (2016). Characterization of Greater Middle Eastern genetic variation for enhanced disease gene discovery. Nat. Genet.

879   5. Lohmueller, K.E., Indap, A.R., Schmidt, S., Boyko, A.R., Hernandez, R.D., Hubisz, M.J.,
880   Sninsky, J.J., White, T.J., Sunyaev, S.R., and Nielsen, R. (2008). Proportionally more deleterious
881   genetic variation in European than in African populations. Nature *451*, 994.

882   6. Charlesworth, D., and Willis, J.H. (2009). The genetics of inbreeding depression. Nat. Rev.
883   Genet. *10*, 783.

884   7. Lohmueller, K.E. (2014). The Impact of Population Demography and Selection on the Genetic
885   Architecture of Complex Traits. PLOS Genet. *10*, e1004379.

886   8. Simons, Y.B., Turchin, M.C., Pritchard, J.K., and Sella, G. (2014). The deleterious mutation
887   load is insensitive to recent population history. Nat. Genet. *46*, 220–224.

888   9. Service, S., DeYoung, J., Karayiorgou, M., Roos, J.L., Pretorious, H., Bedoya, G., Ospina, J.,
889   Ruiz-Linares, A., Macedo, A., Palha, J.A., et al. (2006). Magnitude and distribution of linkage
890   disequilibrium in population isolates and implications for genome-wide association studies. Nat.
891   Genet. *38*, 556–560.

892   10. Xue, Y., Mezzavilla, M., Haber, M., McCarthy, S., Chen, Y., Narasimhan, V., Gilly, A.,
893   Ayub, Q., Colonna, V., Southam, L., et al. (2017). Enrichment of low-frequency functional
894   variants revealed by whole-genome sequencing of multiple isolated European populations. Nat.
895   Commun. *8*, 15927.

896   11. Kittles, R.A., Perola, M., Peltonen, L., Bergen, A.W., Aragon, R.A., Virkkunen, M.,
897   Linnoila, M., Goldman, D., and Long, J.C. (1998). Dual origins of Finns revealed by Y
898   chromosome haplotype variation. Am. J. Hum. Genet. *62*, 1171–1179.

899   12. Peltonen, L., Jalanko, A., and Varilo, T. (1999). Molecular genetics the Finnish disease
900   heritage. Hum. Mol. Genet. *8*, 1913–1923.

901   13. Wang, S.R., Agarwala, V., Flannick, J., Chiang, C.W.K., Altshuler, D., Flannick, J.,
902   Manning, A., Hartl, C., Agarwala, V., Fontanillas, P., et al. (2014). Simulation of Finnish
903   Population History, Guided by Empirical Genetic Data, to Assess Power of Rare-Variant Tests in
904   Finland. Am. J. Hum. Genet. *94*, 710–720.

905   14. De La Chapelle, A., and Wright, F.A. (1998). Linkage disequilibrium mapping in isolated
906   populations: the example of Finland revisited. Proc. Natl. Acad. Sci. *95*, 12416–12423.

907   15. Martin, A.R., Karczewski, K.J., Kerminen, S., Kurki, M., Sarin, A.-P., Artomov, M.,
908   Eriksson, J.G., Esko, T., Genovese, G., Havulinna, A.S., et al. (2017). Haplotype sharing
909   provides insights into fine-scale population history and disease in Finland. BioRxiv 200113.

910   16. Panoutsopoulou, K., Hatzikotoulas, K., Xifara, D.K., Colonna, V., Farmaki, A.-E., Ritchie,
911   G.R.S., Southam, L., Gilly, A., Tachmazidou, I., Fatumo, S., et al. (2014). Genetic
912   characterization of Greek population isolates reveals strong genetic drift at missense and trait-
913   associated variants. Nat. Commun. *5*, 5345.

914    17. Nakatsuka, N., Moorjani, P., Rai, N., Sarkar, B., Tandon, A., Patterson, N., Bhavani, G.S.,
915    Girisha, K.M., Mustak, M.S., and Srinivasan, S. (2017). The promise of discovering population-
916    specific disease-associated genes in South Asia. Nat. Genet. *49*, 1403.

917    18. Pedersen, C.-E.T., Lohmueller, K.E., Grarup, N., Bjerregaard, P., Hansen, T., Siegismund,
918    H.R., Moltke, I., and Albrechtsen, A. (2017). The Effect of an Extreme and Prolonged
919    Population Bottleneck on Patterns of Deleterious Variation: Insights from the Greenlandic Inuit.
920    Genetics *205*, 787–801.

921    19. Pemberton, T.J., Absher, D., Feldman, M.W., Myers, R.M., Rosenberg, N.A., and Li, J.Z.
922    (2012). Genomic patterns of homozygosity in worldwide human populations. Am. J. Hum.
923    Genet. *91*, 275–292.

924    20. Pemberton, T.J., and Szpiech, Z.A. (2018). Relationship between Deleterious Variation,
925    Genomic Autozygosity, and Disease Risk: Insights from The 1000 Genomes Project. Am. J.
926    Hum. Genet. *0*,.

927    21. Kang, J.T., Goldberg, A., Edge, M.D., Behar, D.M., and Rosenberg, N.A. (2016).
928    Consanguinity Rates Predict Long Runs of Homozygosity in Jewish Populations. Hum. Hered.
929    *82*, 87–102.

930    22. Gaetano, C.D., Fiorito, G., Ortu, M.F., Rosa, F., Guarrera, S., Pardini, B., Cusi, D., Frau, F.,
931    Barlassina, C., Troffa, C., et al. (2014). Sardinians Genetic Background Explained by Runs of
932    Homozygosity and Genomic Regions under Positive Selection. PLOS ONE *9*, e91237.

933    23. Carvajal-Carmona, L.G., Ophoff, R., Hartiala, J., Molina, J., Leon, P., Ospina, J., Bedoya,
934    G., Freimer, N., and Ruiz-Linares, A. (2003). Genetic demography of Antioquia (Colombia) and
935    the central valley of Costa Rica. Hum. Genet. *112*, 534–541.

936    24. Consortium, 1000 Genomes Project (2015). A global reference for human genetic variation.

937    25. Fears, S.C., Kremeyer, B., Araya, C., Araya, X., Bejarano, J., Ramirez, M., Castrillón, G.,
938    Gomez-Franco, J., Lopez, M.C., and Montoya, G. (2014). Multisystem component phenotypes of
939    bipolar disorder for genetic investigations of extended pedigrees. JAMA Psychiatry *71*, 375–387.

940    26. Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.-M. (2010).
941    Robust relationship inference in genome-wide association studies. Bioinformatics *26*, 2867–
942    2873.

943    27. DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl, C., Philippakis,
944    A.A., Del Angel, G., Rivas, M.A., and Hanna, M. (2011). A framework for variation discovery
945    and genotyping using next-generation DNA sequencing data. Nat. Genet. *43*, 491.

946    28. Watterson, G.A. (1975). On the number of segregating sites in genetical models without
947    recombination. Theor. Popul. Biol. *7*, 256–276.

948   29. Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., Handsaker,
949   R.E., Lunter, G., Marth, G.T., and Sherry, S.T. (2011). The variant call format and VCFtools.
950   Bioinformatics *27*, 2156–2158.

951   30. Browning, B.L., and Browning, S.R. (2013). Detecting Identity by Descent and Estimating
952   Genotype Error Rates in Sequence Data. Am. J. Hum. Genet. *93*, 840–851.

953   31. Browning, S.R., and Browning, B.L. (2007). Rapid and accurate haplotype phasing and
954   missing-data inference for whole-genome association studies by use of localized haplotype
955   clustering. Am. J. Hum. Genet. *81*, 1084–1097.

956   32. Gusev, A., Lowe, J.K., Stoffel, M., Daly, M.J., Altshuler, D., Breslow, J.L., Friedman, J.M.,
957   and Pe'er, I. (2009). Whole population, genome-wide mapping of hidden relatedness. Genome
958   Res. *19*, 318–326.

959   33. Browning, B.L., and Browning, S.R. (2013). Improving the accuracy and efficiency of
960   identity-by-descent detection in population data. Genetics *194*, 459–471.

961   34. Kong, A., Gudbjartsson, D.F., Sainz, J., Jonsdottir, G.M., Gudjonsson, S.A., Richardsson, B.,
962   Sigurdardottir, S., Barnard, J., Hallbeck, B., and Masson, G. (2002). A high-resolution
963   recombination map of the human genome. Nat. Genet. *31*, 241.

964   35. Browning, S.R., and Browning, B.L. (2015). Accurate non-parametric estimation of recent
965   effective population size from segments of identity by descent. Am. J. Hum. Genet. *97*, 404–418.

966   36. Auton, A., Bryc, K., Boyko, A.R., Lohmueller, K.E., Novembre, J., Reynolds, A., Indap, A.,
967   Wright, M.H., Degenhardt, J.D., Gutenkunst, R.N., et al. (2009). Global distribution of genomic
968   diversity underscores rich complex history of continental human populations. Genome Res. *19*,
969   795–803.

970   37. Sinnwell, J.P., Therneau, T.M., and Schaid, D.J. (2014). The kinship2 R package for
971   pedigree data. Hum. Hered. *78*, 91–93.

972   38. Haller, B.C., and Messer, P.W. (2016). SLiM 2: flexible, interactive forward genetic
973   simulations. Mol. Biol. Evol. *34*, 230–240.

974   39. Kim, B.Y., Huber, C.D., and Lohmueller, K.E. (2017). Inference of the distribution of
975   selection coefficients for new nonsynonymous mutations using large samples. Genetics *206*,
976   345–361.

977   40. Huber, C.D., Kim, B.Y., Marsden, C.D., and Lohmueller, K.E. (2017). Determining the
978   factors driving selective effects of new nonsynonymous mutations. Proc. Natl. Acad. Sci. *114*,
979   4465–4470.

980   41. Long, M., and Deutsch, M. (1999). Association of intron phases with conservation at splice
981   site sequences and evolution of spliceosomal introns. Mol. Biol. Evol. *16*, 1528–1534.

982  42. Cooper, G.M., Stone, E.A., Asimenos, G., Green, E.D., Batzoglou, S., and Sidow, A. (2005).
983  Distribution and intensity of constraint in mammalian genomic sequence. Genome Res. *15*, 901–
984  913.

985  43. Alexander, D.H., Novembre, J., and Lange, K. (2009). Fast model-based estimation of
986  ancestry in unrelated individuals. Genome Res. *19*, 1655–1664.

987  44. Sankararaman, S., Sridhar, S., Kimmel, G., and Halperin, E. (2008). Estimating local
988  ancestry in admixed populations. Am. J. Hum. Genet. *82*, 290–303.

989  45. Pagani, L., Clair, P.A.S., Teshiba, T.M., Fears, S.C., Araya, C., Araya, X., Bejarano, J.,
990  Ramirez, M., Castrillón, G., and Gomez-Makhinson, J. (2016). Genetic contributions to circadian
991  activity rhythm and sleep pattern phenotypes in pedigrees segregating for severe bipolar
992  disorder. Proc. Natl. Acad. Sci. *113*, E754–E761.

993  46. Consortium, I.H. (2003). The international HapMap project. Nature *426*, 789.

994  47. Consortium, I.H. (2007). A second generation human haplotype map of over 3.1 million
995  SNPs. Nature *449*, 851.

996  48. Reich, D., Patterson, N., Campbell, D., Tandon, A., Mazieres, S., Ray, N., Parra, M.V.,
997  Rojas, W., Duque, C., and Mesa, N. (2012). Reconstructing native American population history.
998  Nature *488*, 370.

999  49. Aulchenko, Y.S., Ripke, S., Isaacs, A., and Van Duijn, C.M. (2007). GenABEL: an R library
1000  for genome-wide association analysis. Bioinformatics *23*, 1294–1296.

1001  50. Lohmueller, K.E., Albrechtsen, A., Li, Y., Kim, S.Y., Korneliussen, T., Vinckenbosch, N.,
1002  Tian, G., Huerta-Sanchez, E., Feder, A.F., and Grarup, N. (2011). Natural selection affects
1003  multiple aspects of genetic variation at putatively neutral sites across the human genome. PLoS
1004  Genet. *7*, e1002326.

1005  51. Cai, J.J., Macpherson, J.M., Sella, G., and Petrov, D.A. (2009). Pervasive hitchhiking at
1006  coding and regulatory sites in humans. PLoS Genet. *5*, e1000336.

1007  52. Hernandez, R.D., Kelley, J.L., Elyashiv, E., Melton, S.C., Auton, A., McVean, G., Sella, G.,
1008  and Przeworski, M. (2011). Classic selective sweeps were rare in recent human evolution.
1009  Science *331*, 920–924.

1010  53. Stumpf, M.P., and Goldstein, D.B. (2003). Demography, recombination hotspot intensity,
1011  and the block structure of linkage disequilibrium. Curr. Biol. *13*, 1–8.

1012  54. Pritchard, J.K., and Przeworski, M. (2001). Linkage disequilibrium in humans: models and
1013  data. Am. J. Hum. Genet. *69*, 1–14.

1014  55. Gravel, S., Zakharia, F., Moreno-Estrada, A., Byrnes, J.K., Muzzio, M., Rodriguez-Flores,
1015  J.L., Kenny, E.E., Gignoux, C.R., Maples, B.K., Guiblet, W., et al. (2013). Reconstructing

1016    Native American Migrations from Whole-Genome and Whole-Exome Data. PLOS Genet. *9*,
1017    e1004023.

1018    56. Szpiech, Z.A., Xu, J., Pemberton, T.J., Peng, W., Zöllner, S., Rosenberg, N.A., and Li, J.Z.
1019    (2013). Long runs of homozygosity are enriched for deleterious variation. Am. J. Hum. Genet.
1020    *93*, 90–102.

1021    57. McQuillan, R., Leutenegger, A.-L., Abdel-Rahman, R., Franklin, C.S., Pericic, M., Barac-
1022    Lauc, L., Smolej-Narancic, N., Janicijevic, B., Polasek, O., and Tenesa, A. (2008). Runs of
1023    homozygosity in European populations. Am. J. Hum. Genet. *83*, 359–372.

1024    58. Kirin, M., McQuillan, R., Franklin, C.S., Campbell, H., McKeigue, P.M., and Wilson, J.F.
1025    (2010). Genomic runs of homozygosity record population history and consanguinity. PloS One
1026    *5*, e13996.

1027    59. Pemberton, T.J., and Szpiech, Z.A. (2018). Relationship between Deleterious Variation,
1028    Genomic Autozygosity, and Disease Risk: Insights from The 1000 Genomes Project. Am. J.
1029    Hum. Genet. *0*,.

1030    60. Ceballos, F.C., Joshi, P.K., Clark, D.W., Ramsay, M., and Wilson, J.F. (2018). Runs of
1031    homozygosity: windows into population history and trait architecture. Nat. Rev. Genet.

1032    61. Kardos, M., Taylor, H.R., Ellegren, H., Luikart, G., and Allendorf, F.W. (2016). Genomics
1033    advances the study of inbreeding depression in the wild. Evol. Appl. *9*, 1205–1218.

1034    62. Kimura, M., Maruyama, T., and Crow, J.F. (1963). The mutation load in small populations.
1035    Genetics *48*, 1303–1312.

1036    63. Ohta, T. (1973). Slightly deleterious mutant substitutions in evolution. Nature *246*, 96.

1037    64. Hodgkinson, A., Casals, F., Idaghdour, Y., Grenier, J.-C., Hernandez, R.D., and Awadalla, P.
1038    (2013). Selective constraint, background selection, and mutation accumulation variability within
1039    and between human populations. BMC Genomics *14*, 495.

1040    65. Peischl, S., Dupanloup, I., Kirkpatrick, M., and Excoffier, L. (2013). On the accumulation of
1041    deleterious mutations during range expansions. Mol. Ecol. *22*, 5972–5982.

1042    66. Fu, W., Gittelman, R.M., Bamshad, M.J., and Akey, J.M. (2014). Characteristics of Neutral
1043    and Deleterious Protein-Coding Variation among Individuals and Populations. Am. J. Hum.
1044    Genet. *95*, 421–436.

1045    67. Escamilla, M.A., Spesny, M., Reus, V.I., Gallegos, A., Meza, L., Molina, J., Sandkuijl, L.A.,
1046    Fournier, E., Leon, P.E., Smith, L.B., et al. (1996). Use of linkage disequilibrium approaches to
1047    map genes for bipolar disorder in the Costa Rican population. Am. J. Med. Genet. *67*, 244–253.

1048    68. Wang, S., Ray, N., Rojas, W., Parra, M.V., Bedoya, G., Gallo, C., Poletti, G., Mazzotti, G.,
1049    Hill, K., and Hurtado, A.M. (2008). Geographic patterns of genome admixture in Latin American
1050    Mestizos. PLoS Genet. *4*, e1000037.

69. Bedoya, G., Montoya, P., García, J., Soto, I., Bourgeois, S., Carvajal, L., Labuda, D., Alvarez, V., Ospina, J., and Hedrick, P.W. (2006). Admixture dynamics in Hispanics: a shift in the nuclear genetic ancestry of a South American population isolate. Proc. Natl. Acad. Sci. *103*, 7234–7239.

70. Safford, F., and Palacios, M. (2002). Colombia: Fragmented land, divided society (Oxford University Press, USA).

71. Carvajal-Carmona, L.G., Soto, I.D., Pineda, N., Ortíz-Barrientos, D., Duque, C., Ospina-Duque, J., McCarthy, M., Montoya, P., Alvarez, V.M., and Bedoya, G. (2000). Strong Amerind/white sex bias and a possible Sephardic contribution among the founders of a population in northwest Colombia. Am. J. Hum. Genet. *67*, 1287–1295.

72. Ceballos, F.C., Joshi, P.K., Clark, D.W., Ramsay, M., and Wilson, J.F. (2018). Runs of homozygosity: windows into population history and trait architecture. Nat. Rev. Genet.

73. Li, L.-H., Ho, S.-F., Chen, C.-H., Wei, C.-Y., Wong, W.-C., Li, L.-Y., Hung, S.-I., Chung, W.-H., Pan, W.-H., Lee, M.-T.M., et al. (2006). Long contiguous stretches of homozygosity in the human genome. Hum. Mutat. *27*, 1115–1121.

74. Jorde, L.B., and Pitkänen, K.J. (1991). Inbreeding in Finland. Am. J. Phys. Anthropol. *84*, 127–139.

75. Wright, S. (1984). Evolution and the genetics of populations, volume 3: experimental results and evolutionary deductions (University of Chicago press).

76. Charlesworth, B., and Charlesworth, D. (1999). The genetic basis of inbreeding depression. Genet. Res. *74*, 329–340.

77. Wang, J., Hill, W.G., Charlesworth, D., and Charlesworth, B. (1999). Dynamics of inbreeding depression due to deleterious mutations in small populations: mutation parameters and inbreeding rate. Genet. Res. *74*, 165–178.

78. Balick, D.J., Do, R., Cassa, C.A., Reich, D., and Sunyaev, S.R. (2015). Dominance of deleterious alleles controls the response to a population bottleneck. PLoS Genet. *11*, e1005436.

79. Mukai, T., Chigusa, S.I., Mettler, L.E., and Crow, J.F. (1972). Mutation rate and dominance of genes affecting viability in Drosophila melanogaster. Genetics *72*, 335–355.

80. Simmons, M.J., and Crow, J.F. (1977). Mutations affecting fitness in Drosophila populations. Annu. Rev. Genet. *11*, 49–78.

81. Phadnis, N., and Fry, J.D. (2005). Widespread correlations between dominance and homozygous effects of mutations: implications for theories of dominance. Genetics *171*, 385–392.

82. Agrawal, A.F., and Whitlock, M.C. (2011). Inferences about the distribution of dominance drawn from yeast gene knockout data. Genetics *187*, 553–566.

43

1086  83. Myerowitz, R., and Costigan, F.C. (1988). The major defect in Ashkenazi Jews with Tay-
1087  Sachs disease is an insertion in the gene for the alpha-chain of beta-hexosaminidase. J. Biol.
1088  Chem. *263*, 18587–18589.

1089  84. Hästbacka, J., de la Chapelle, A., Mahtani, M.M., Clines, G., Reeve-Daly, M.P., Daly, M.,
1090  Hamilton, B.A., Kusumi, K., Trivedi, B., and Weaver, A. (1994). The diastrophic dysplasia gene
1091  encodes a novel sulfate transporter: positional cloning by fine-structure linkage disequilibrium
1092  mapping. Cell *78*, 1073–1087.

1093  85. Aittomäki, K., Lucena, J.D., Pakarinen, P., Sistonen, P., Tapanainen, J., Gromoll, J.,
1094  Kaskikari, R., Sankila, E.-M., Lehväslaiho, H., and Engel, A.R. (1995). Mutation in the follicle-
1095  stimulating hormone receptor gene causes hereditary hypergonadotropic ovarian failure. Cell *82*,
1096  959–968.

1097  86. Ruiz-Perez, V.L., Ide, S.E., Strom, T.M., Lorenz, B., Wilson, D., Woods, K., King, L.,
1098  Francomano, C., Freisinger, P., and Spranger, S. (2000). Mutations in a new gene in Ellis-van
1099  Creveld syndrome and Weyers acrodental dysostosis. Nat. Genet. *24*, 283.

1100  87. Verhoeven, K., Villanova, M., Rossi, A., Malandrini, A., De Jonghe, P., and Timmerman, V.
1101  (2001). Localization of the gene for the intermediate form of Charcot-Marie-Tooth to
1102  chromosome 10q24. 1-q25. 1. Am. J. Hum. Genet. *69*, 889–894.

1103  88. Valente, E.M., Bentivoglio, A.R., Dixon, P.H., Ferraris, A., Ialongo, T., Frontali, M.,
1104  Albanese, A., and Wood, N.W. (2001). Localization of a novel locus for autosomal recessive
1105  early-onset parkinsonism, PARK6, on human chromosome 1p35-p36. Am. J. Hum. Genet. *68*,
1106  895–900.

1107  89. McInnes, L.A., Reus, V.I., Barnes, G., Charlat, O., Jawahar, S., Lewitzky, S., Yang, Q.,
1108  Duong, Q., Spesny, M., and Araya, C. (2001). Fine-scale mapping of a locus for severe bipolar
1109  mood disorder on chromosome 18p11. 3 in the Costa Rican population. Proc. Natl. Acad. Sci.
1110  *98*, 11485–11490.

1111  90. Ober, C., Tan, Z., Sun, Y., Possick, J.D., Pan, L., Nicolae, R., Radford, S., Parry, R.R.,
1112  Heinzmann, A., and Deichmann, K.A. (2008). Effect of variation in CHI3L1 on serum YKL-40
1113  level, risk of asthma, and lung function. N. Engl. J. Med. *358*, 1682–1691.

1114  91. Stacey, S.N., Manolescu, A., Sulem, P., Thorlacius, S., Gudjonsson, S.A., Jonsson, G.F.,
1115  Jakobsdottir, M., Bergthorsson, J.T., Gudmundsson, J., Aben, K.K., et al. (2008). Common
1116  variants on chromosome 5p12 confer susceptibility to estrogen receptor–positive breast cancer.
1117  Nat. Genet. *40*, 703–706.

1118  92. Stacey, S.N., Sulem, P., Jonasdottir, A., Masson, G., Gudmundsson, J., Gudbjartsson, D.F.,
1119  Magnusson, O.T., Gudjonsson, S.A., Sigurgeirsson, B., and Thorisdottir, K. (2011). A germline
1120  variant in the TP53 polyadenylation signal confers cancer susceptibility. Nat. Genet. *43*, 1098.

1121  93. Gudmundsson, J., Sulem, P., Gudbjartsson, D.F., Masson, G., Agnarsson, B.A.,
1122  Benediktsdottir, K.R., Sigurdsson, A., Magnusson, O.T., Gudjonsson, S.A., and Magnusdottir,

44

1123    D.N. (2012). A study based on whole-genome sequencing yields a rare variant at 8q24 associated
1124    with prostate cancer. Nat. Genet. *44*, 1326.

1125    94. Tachmazidou, I., Dedoussis, G., Southam, L., Farmaki, A.-E., Ritchie, G.R., Xifara, D.K.,
1126    Matchan, A., Hatzikotoulas, K., Rayner, N.W., and Chen, Y. (2013). A rare functional
1127    cardioprotective APOC3 variant has risen in frequency in distinct population isolates. Nat.
1128    Commun. *4*, 2872.

1129    95. Saleheen, D., Natarajan, P., Armean, I.M., Zhao, W., Rasheed, A., Khetarpal, S.A., Won, H.-
1130    H., Karczewski, K.J., O'Donnell-Luria, A.H., and Samocha, K.E. (2017). Human knockouts and
1131    phenotypic analysis in a cohort with a high rate of consanguinity. Nature *544*, 235.

1132    96. Botstein, D., and Risch, N. (2003). Discovering genotypes underlying human phenotypes:
1133    past successes for mendelian disease, future approaches for complex disease. Nat. Genet. *33*,
1134    228.

1135    97. Lencz, T., Lambert, C., DeRosse, P., Burdick, K.E., Morgan, T.V., Kane, J.M., Kucherlapati,
1136    R., and Malhotra, A.K. (2007). Runs of homozygosity reveal highly penetrant recessive loci in
1137    schizophrenia. Proc. Natl. Acad. Sci. *104*, 19942–19947.

1138    98. Mezzavilla, M., Vozzi, D., Badii, R., Alkowari, M.K., Abdulhadi, K., Girotto, G., and
1139    Gasparini, P. (2015). Increased rate of deleterious variants in long runs of homozygosity of an
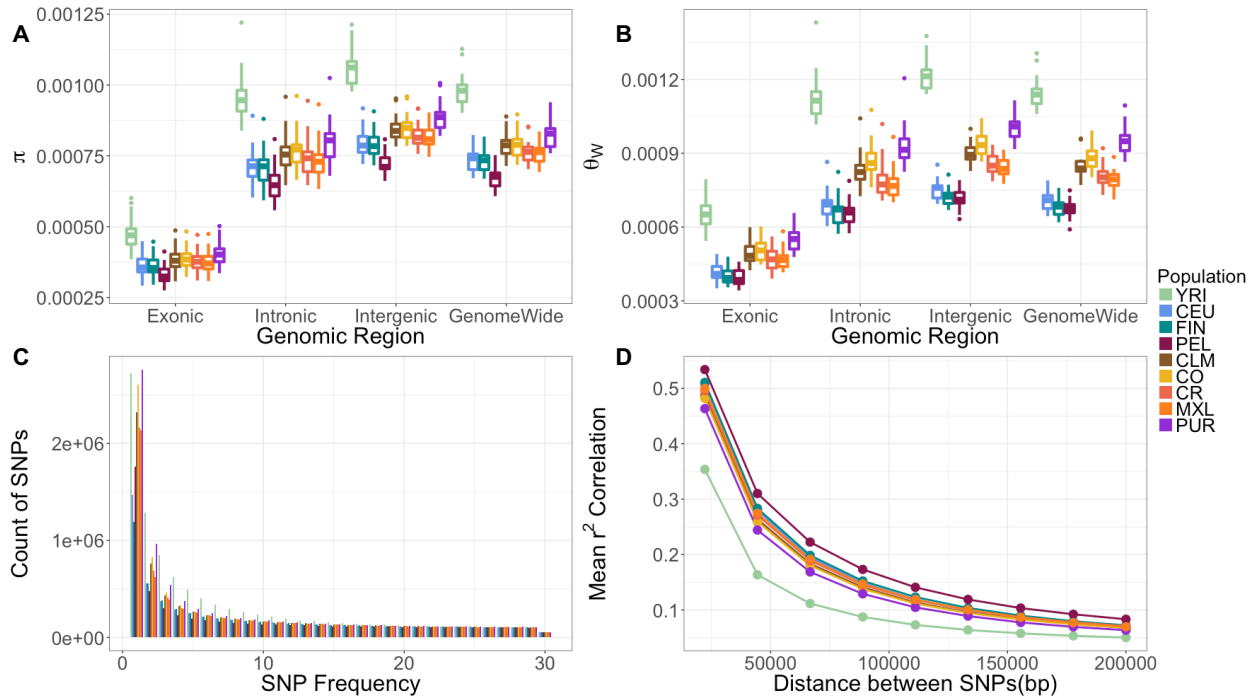1140    inbred population from Qatar. Hum. Hered. *79*, 14–19.

1141

1142 **Figures**
1143



1144
1145
1146 **Figure 1. Patterns of genetic variation in the Colombian and Costa Rican populations**
1147 **compared to the 1000 Genomes populations.**
1148 (A) Diversity measured using the average pairwise differences between sequences, $\pi$. (B)
1149 Diversity measured using the number segregating sites, Watterson's theta ($\theta_W$). (C) The site
1150 frequency spectrum for each population. (D) Average LD ($r^2$) between pairs of SNPs. All
1151 statistics were calculated using 30 unrelated individuals per population (see Methods). Box plots
1152 in (A) and (B) show the distribution over 22 autosomes. YRI: Yoruba 1000 Genomes; CEU:
1153 Ceph-European 1000 Genomes; FIN: Finnish 1000 Genomes; PEL: Peruvian 1000 Genomes;
1154 CLM: Colombian 1000 Genomes; CO: Colombia; CR: Costa Rica; MXL: Mexican from Los
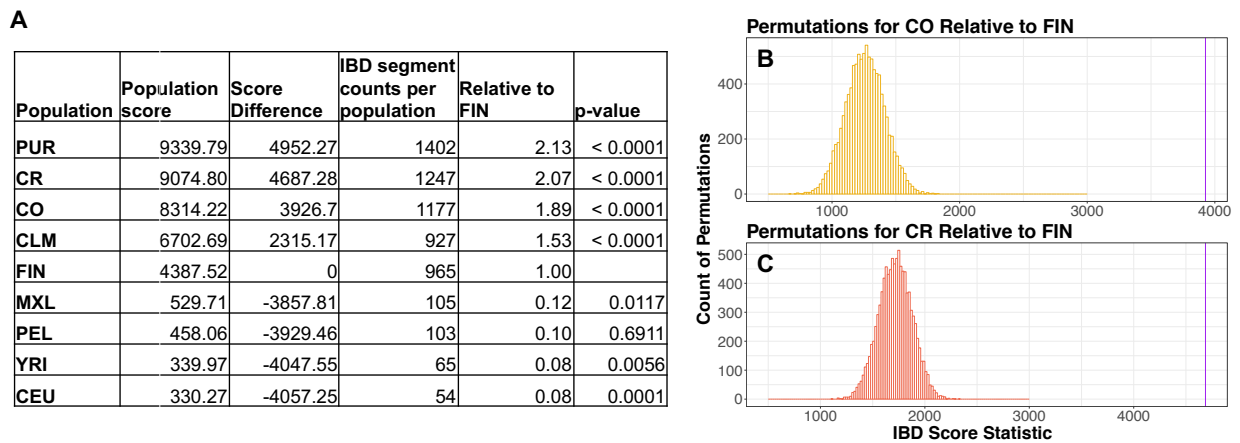1155 Angeles 1000 Genomes; and PUR: Puerto Rican 1000 Genomes.
1156

46

**A**

| Population | Population score | Score Difference | IBD segment counts per population | Relative to FIN | p-value |
|---|---|---|---|---|---|
| PUR | 9339.79 | 4952.27 | 1402 | 2.13 | < 0.0001 |
| CR | 9074.80 | 4687.28 | 1247 | 2.07 | < 0.0001 |
| CO | 8314.22 | 3926.7 | 1177 | 1.89 | < 0.0001 |
| CLM | 6702.69 | 2315.17 | 927 | 1.53 | < 0.0001 |
| FIN | 4387.52 | 0 | 965 | 1.00 | |
| MXL | 529.71 | -3857.81 | 105 | 0.12 | 0.0117 |
| PEL | 458.06 | -3929.46 | 103 | 0.10 | 0.6911 |
| YRI | 339.97 | -4047.55 | 65 | 0.08 | 0.0056 |
| CEU | 330.27 | -4057.25 | 54 | 0.08 | 0.0001 |



Permutations for CO Relative to FIN

Permutations for CR Relative to FIN

**Figure 2**. **Latin American population isolates (CR and CO) have significantly more identity by descent (IBD) segments relative to the Finnish (FIN).**
IBDSeq was used to generate IBD segments for the 30 unrelated individuals in each population. (A) Population score was calculated by summing all IBD segments between 3cM and 20cM for each population. Score difference is the population score minus the FIN population score. IBD enrichment for each population score is reported as relative to the FIN (i.e. FIN score is 1.0). (B&C) Histogram of 10,000 permutation tests of Colombia ($p < 1.0\ e^{-04}$) and Costa Rica ($p < 1.0\ e^{-04}$) population scores versus Finnish score. The observed score for each population is demarcated by the purple line. Population abbreviations are as in **Figure 1**.
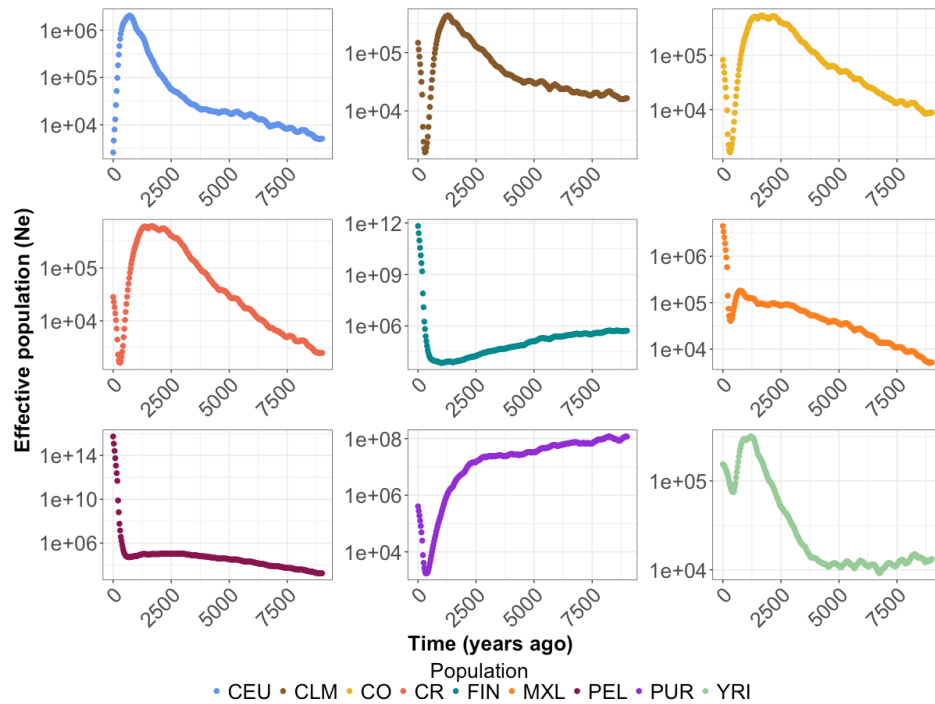
**Figure 3**. **Recent effective population size differs across populations.** IBDNe[35] (see Methods) was used to infer effective population size (Ne) over the last 9000 years for each population. Note the FIN shows a long slow decline followed by recent growth. The CO and CR show sharp bottlenecks approximately 500 years ago followed by recent growth. Population abbreviations are as in **Figure 1**.
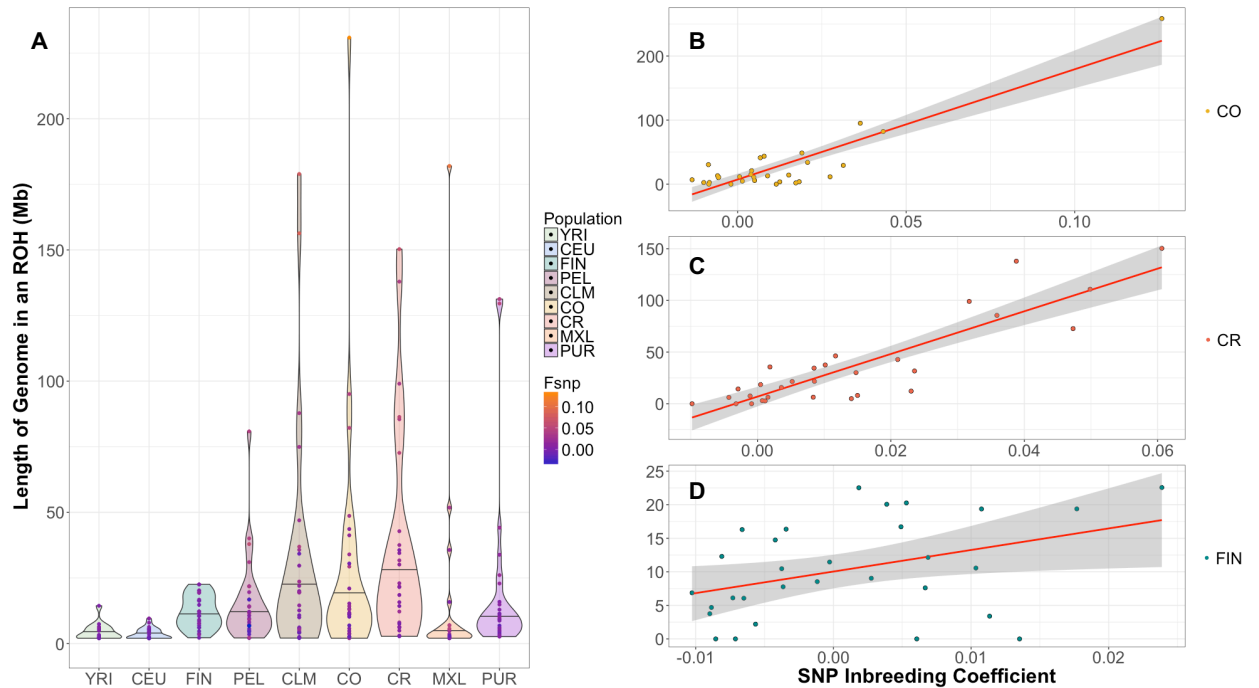
1177
1178
1179 **Figure 4**. **Length of the genome in a run of homozygosity (ROH) varies across populations**
1180 **and correlates with SNP inbreeding coefficient.** The length of the genome in an ROH was
1181 calculated for each unrelated individual ($n$=30 per population) by summing the physical distance
1182 (Mb) of each ROH >2Mb. (A) The length of the genome in an ROH varies by population. The
1183 black line within the violin marks the median. $F_{SNP}$ for each individual was overlaid within the
1184 ROH violin plot. A blue hue indicates the lowest $F_{SNP}$ and orange indicates the highest $F_{SNP}$. (B)
1185 Length of the genome in an ROH is strongly correlated with $F_{SNP}$ in Colombians, ($R^2 = 0.8060$, p
1186 –value = $1.1 \times 10^{-11}$). (C) Length of the genome in an ROH is strongly correlated with $F_{SNP}$ in
1187 Costa Ricans, ($R^2 = 0.7740$, p-value = $9.5 \times 10^{-11}$). (D) Length of the genome in an ROH is
1188 positively correlated with $F_{SNP}$ in Finnish, ($R^2 = 0.1288$, p-value = 0.03). Population
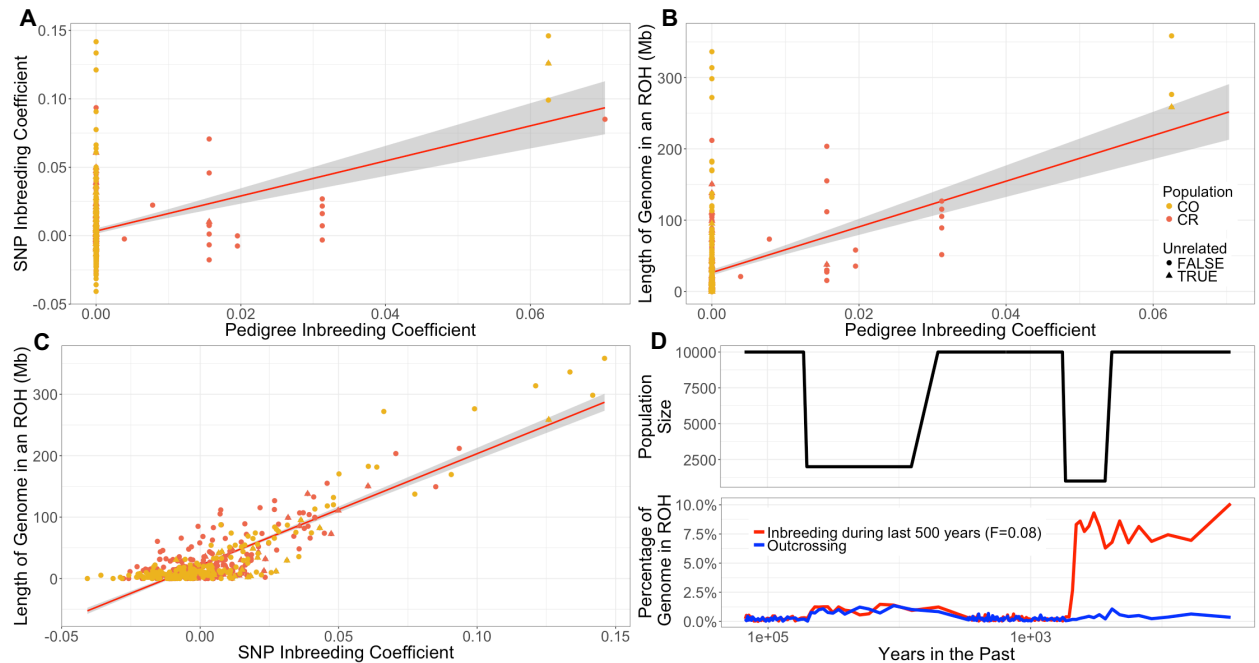1189 abbreviations are as in **Figure 1**.
1190

49

**Figure 5. Recent consanguinity creates ROH in Costa Rica and Colombia.**
Triangles represent the individuals that were sampled in the unrelated data set ($n$=30). (A) $F_{SNP}$ is correlated with the pedigree inbreeding coefficient ($F_{PED}$; $R^2$=0.1520, p-value < 2 x$10^{-16}$) in the full data. (B) The length of the genome in an ROH is correlated with $F_{PED}$ ($R^2$ = 0.2180, p-value < 2 x $10^{-16}$). (C) The length of the genome in an ROH is correlated with $F_{SNP}$ ($R^2$ = 0.7540, p-value < 2 x $10^{-16}$). (D) Forward simulations show that recent consanguinity during the last 500 years can generate ROHs while bottlenecks cannot. Top panel shows the changes in population size used in the simulations. Bottom panel shows how the percent of the simulated in genome within an ROH changes over time. Population abbreviations are as in **Figure 1**.
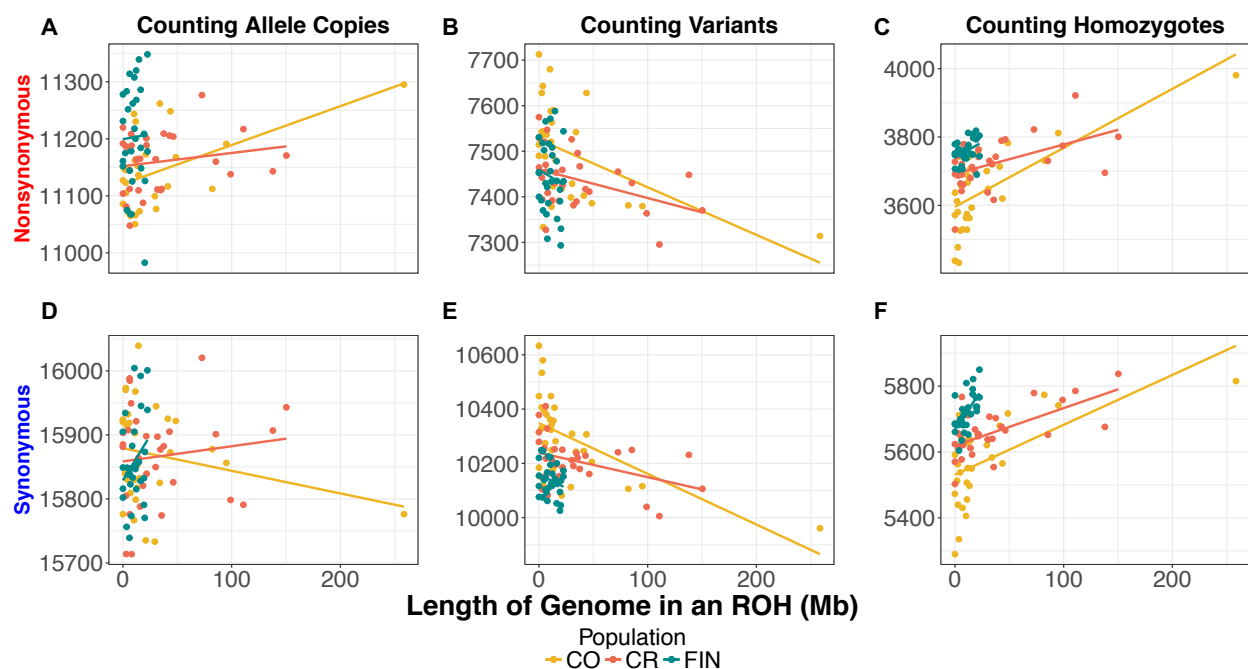
50

1203
1204
1205 **Figure 6. The correlation between ROH and nonsynonymous variation in the Colombian,**
1206 **Costa Rican, and Finnish samples.** The count of nonsynonymous and synonymous mutations
1207 per individual as a function of the length of the genome in an ROH in the Colombia (CO), Costa
1208 Rican (CR) and Finnish (FIN) populations: (A) Number of nonsynonymous alleles per
1209 individual. (B) Number of nonsynonymous variants per individual. (C) Number of homozygous
1210 nonsynonymous genotypes per individual. (D) The number of synonymous alleles per individual.
1211 (E) The number of synonymous variants per individual. (F) The number of homozygous
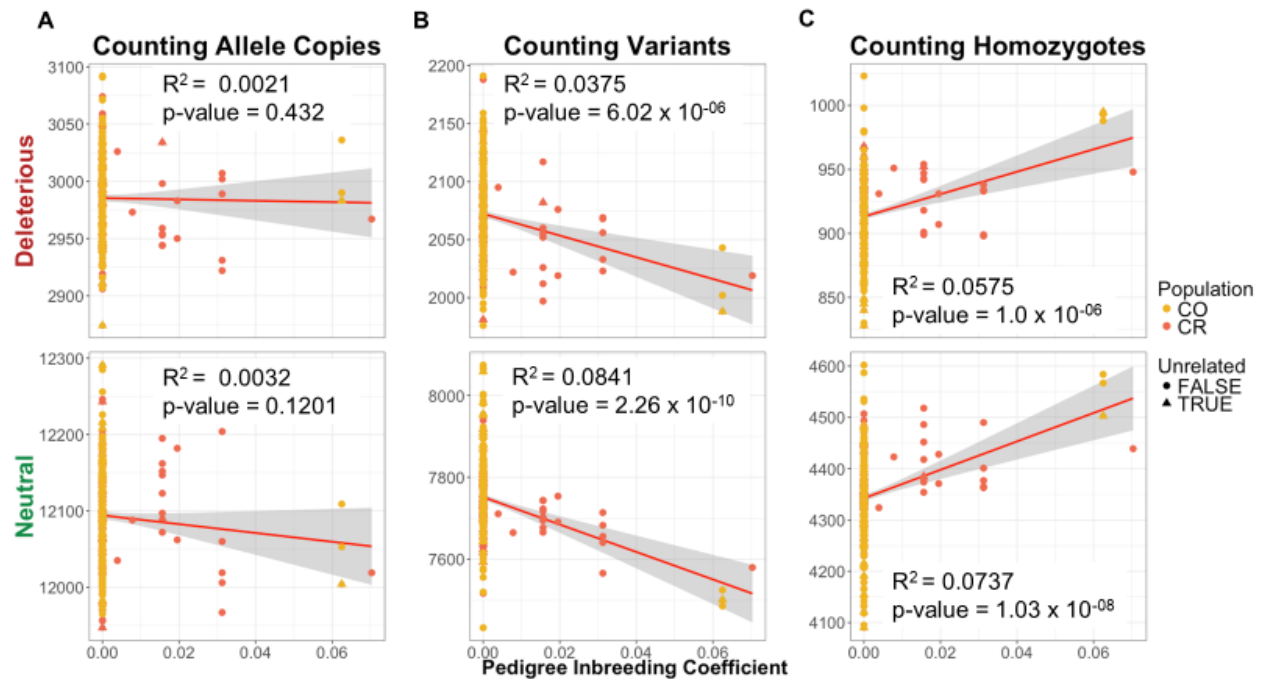1212 synonymous genotypes per individual. Population abbreviations are as in **Figure 1**.
1213

51

**Figure 7. Pedigree inbreeding coefficient ($F_{PED}$) is correlated with deleterious variation.** Triangles represent the individuals that were sampled in the unrelated data set (*n*=30). Variants were predicted as either putatively deleterious (nonsynonymous) SNPs or putatively neutral (synonymous) SNPs using GERP[42]. Correlation between $F_{PED}$ and the number of mutations per individual in Colombians and Costa Ricans. (A) Number of derived alleles per individual. (B) Number of variants per individual. (C) Number of homozygous derived genotypes per individual. The first row depicts the correlation between deleterious sites using each counting method and $F_{PED}$ for sequenced individuals from Latin American isolates. The second row depicts the correlation between neutral sites using each counting method and $F_{PED}$ in the same individuals. Population abbreviations are as in **Figure 1**.

1227 **Table 1. Enrichment of nonsynonymous homozygous derived genotypes within ROHs**

1228

| Population | Allele Copies Odds Ratio | Allele Copies p-value | Variants Odds Ratio | Variants p-value | Homozygotes Odds Ratio | Homozygotes p-value |
|---|---|---|---|---|---|---|
| YRI | 1.059 | 0.664 | 1.048 | 0.762 | 1.129 | 0.417 |
| CEU | 1.203 | 0.105 | 1.208 | 0.138 | 1.252 | 0.082 |
| FIN | 0.937 | 0.324 | 0.92 | 0.265 | 1.003 | 0.957 |
| PEL | 0.986 | 0.797 | 0.972 | 0.638 | 1.038 | 0.54 |
| CLM | 0.99 | 0.755 | 0.964 | 0.337 | 1.066 | 0.097 |
| CO | 1.008 | 0.828 | 0.985 | 0.714 | 1.074 | 0.097 |
| CR | 1.015 | 0.607 | 0.991 | 0.806 | 1.085 | **0.0169*** |
| CO & CR | 1.025 | 0.283 | 1.002 | 0.936 | 1.088 | **0.0011*** |
| MXL | 1.112 | 0.053 | 1.089 | 0.169 | 1.19 | **0.005*** |
| PUR | 0.981 | 0.635 | 0.965 | 0.411 | 1.047 | 0.301 |

1229

1230 The table summarizes the results of our enrichment analyses for each population sampled as well
1231 as a combined super-population of Colombians and Costa Ricans (COCR). Odds ratios were
1232 calculated as the ratio nonsynonymous variants relative to synonymous variants within versus
1233 outside of an ROH for each counting method. Asterisks are used to indicate significant p-values
1234 after permutation test was conducted (see **Methods**). Population abbreviations are as in **Figure**
1235 **1**.

1236