Edinburgh Research Explorer

# Enrichment of low-frequency functional variants revealed by whole-genome sequencing of multiple isolated European populations

**Enrichment of low-frequency functional variants revealed by whole-genome sequencing of multiple isolated European populations**

Yali Xue[1*], Massimo Mezzavilla[1,2*], Marc Haber[1*], Shane McCarthy[1*], Yuan Chen[1], Vagheesh Narasimhan[1], Arthur Gilly[1], Qasim Ayub[1], Vincenza Colonna[1,3], Lorraine Southam[1,4], Christopher Finan[1], Andrea Massaia[1,5], Himanshu Chheda[6], Priit Palta[6,7], Graham Ritchie[1,8,9], Jennifer Asimit[1], George Dedoussis[10], Paolo Gasparini[11], Aarno Palotie[1,6,12-16], Samuli Ripatti[1,6,17], Nicole Soranzo[1,18], Daniela Toniolo[19], James F. Wilson[9,20], Richard Durbin[1], Chris Tyler-Smith[1], Eleftheria Zeggini[1]

[1]The Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambs. CB10 1SA, UK.
[2]Institute for Maternal and Child Health, IRCCS Burlo Garofolo, University of Trieste, 34137 Trieste, Italy.
[3]Consiglio Nazionale delle Ricerche, Istituto di Genetica e Biofisica "Adriano Buzzati-Traverso", via Pietro Castellino 111, 80131 Napoli, Italy.
[4]Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK.
[5] National Heart and Lung Institute, Imperial College London, London SW7 2AZ, UK
[6]Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Tukholmankatu 8, 00290 Helsinki, Finland.
[7]Estonian Genome Center, University of Tartu, 23B Riia Street, 51010 Tartu, Estonia.
[8]European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambs. CB10 1SD, UK.
[9]MRC Human Genetics Unit, MRC IGMM, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh EH4 2XU, UK.
[10]Department of Nutrition and Dietetics, Harokopio University Athens, Athens, Eleftheriou Venizelou 70, Kallithea 176 76, Greece.
[11]Medical Genetics, DSM, University of Trieste and IRCCS (Istituto di Ricovero e Cura a Carattere Scientifico) Burlo Garofolo Children Hospital, Via dell'Istria, 65, 34137, Trieste, Italy.
[12]Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA 02114, USA.
[13]Program in Medical and Population Genetics, The Broad Institute of MIT and Harvard, Cambridge, MA 02114, USA.
[14]The Stanley Center for Psychiatric Research, The Broad Institute of MIT and Harvard, Cambridge, MA 02114, USA.
[15]Psychiatric & Neurodevelopmental Genetics Unit, Department of Psychiatry, Massachusetts General Hospital, Boston, MA 02114, USA.
[16]Department of Neurology, Massachusetts General Hospital, Boston, MA 02114, USA.
[17]Department of Public Health, University of Helsinki, Helsinki FI-00014, Finland.
[18]Department of Haematology, University of Cambridge, Cambridge CB2 0XY, UK.

46    [19]Division of Genetics and Cell Biology, San Raffaele Scientific Institute, via Olgettina
47    60, 20132, Milan, Italy.
48    [20]Usher Institute of Population Health Sciences and Informatics, University of
49    Edinburgh, Teviot Place, Edinburgh EH8 9AG, Scotland, UK.
50
51    * Equal contribution to the work
52
53    Correspondence should be addressed to Y. X. (ylx@sanger.ac.uk) or E. Z.
54    (Eleftheria@sanger.ac.uk).
55
58

## Abstract

**The genetic features of isolated populations can boost power in complex-trait association studies, and an in-depth understanding of how their genetic variation has been shaped by their demographic history can help leverage these advantageous characteristics. Here, we perform a comprehensive investigation using 3059 newly-generated low-depth whole-genome sequences from eight European isolates and two matched general populations, together with published data from the 1000 Genomes Project and UK10K. Sequencing data give deeper and richer insights into population demography and genetic characteristics than genotype-chip data, distinguishing related populations more effectively and allowing their functional variants to be studied more fully. We demonstrate relaxation of purifying selection in the isolates, leading to enrichment of rare and low-frequency functional variants, using novel statistics, *DVxy* and *SVxy*. We also develop an isolation-index (*Isx*) that predicts the overall level of such key genetic characteristics and can thus help guide population choice in future complex-trait association studies.**

## Introduction

Population variation in disease susceptibility has been shaped by environment, demography and evolutionary history. Isolated populations (isolates) have generally experienced bottlenecks and strong genetic drift, so by chance some deleterious rare variants have increased in frequency while some neutral rare variation is lost, both helpful characteristics for the discovery of novel rare variant signals underpinning complex traits[1-3]. Studies to date have focused on individual isolates and have identified several disease-associated signals[4-12]. However, isolates differ in the time when they became isolated, their initial population size, the level of gene flow from outside and other historical demographic factors, and consequently also differ in their power for association studies[2]. We thus generate and analyze low-depth (4x-10x) whole-genome sequences (WGS) from eight cohorts drawn from isolated European populations and compare each isolate with the closest non-isolated (general) population, for which we also generate or access WGS data. We then investigate empirically how these historical differences influence the population-genetic properties of isolates, and frame these insights in terms of their consequences for study design in complex trait association studies.

## Results

**Samples, sequencing and QC.** The dataset includes newly-generated low-depth (4x-10x) WGS from eight cohorts drawn from isolated European populations: one each from Kuusamo in Finland (FIK) and Crete in Greece (GRM[13]), four from Friuli-Venezia Giulia in Italy (IF1, IF2, IF3 and IF4[14]), and one each from Val Borbera in

104    Italy (IVB[15]) and the Orkney Islands in the UK (UKO[16]); and the closest non-isolated
105    (general) population: Finland (FIG[9]), Greece (GRG), together with publicly available
106    data for Italy (ITG[17]) and UK (UKG[18]) (Fig. 1a and Supplementary Table 1). We
107    generated a superset of variants called in these cohorts and all 26 population
108    samples in the 1000 Genomes Project Phase 3[17], and performed multi-sample
109    genotype calling across all 9375 samples (3059 from the current study, 2353 from
110    the 1000 Genomes Project Phase 3 release, and 3781 from UK10K). Both individual
111    population and amalgamated genotype call data, which have greater than 99%
112    concordance with genotyping data (Supplementary Table 2), are available to the
113    scientific community (Data availability).
114
115    **General description of the variants in the isolates.** We identified approximately
116    12.2 million (M) variants with minor allele frequency (MAF) ≤2% (rare), 5.5M with
117    MAF >2-≤5% (low-frequency) and 8.3M variants with MAF >5% (common) across
118    the 10 populations newly sequenced here (eight isolates, GRG and FIG). Of these,
119    10.5%, 0.7% and 0.3%, respectively, are novel (Table 1 and Supplementary Table 3).
120    As expected, most of the isolates have lower numbers of variant sites per genome
121    than their closest general population (Supplementary Fig.1, Supplementary Table 5).
122    We find ~188,000-~513,000 variants that are common with MAF >5.5% in each
123    isolate but with MAF ≤1.4% in its closest general population (Table 1); ~30,000-
124    122,000 of these per isolate have frequency ≤1.4% in all the general samples
125    studied, among which ~150-~700 in coding regions and ~500-~2800 genome-wide
126    are deleterious (Supplementary Table 4).  These common and low-frequency
127    variants are thus useful markers for whole-genome association studies in these
128    populations and some of them (if absent from the general population) could
129    potentially lead to novel association signals. They include known examples such as
130    rs76353203 (R19X) in *APOC3* in GRM, which is associated with high-density
131    lipoprotein and triglyceride levels[6].
132
133    **Population-genetic analyses in the isolates.** Previous population-genetic studies
134    of isolates have, with some exceptions[11,19], been based on common variants found
135    on genotyping arrays, and have illustrated general characteristics such as low
136    genetic diversity and longer shared haplotypes[9,13-15,19,20]. Rare variants discovered
137    from sequencing are on average more recent in origin than common variants[21] and
138    therefore more powerful for distinguishing closely-related populations and more
139    informative about recent demographic history. We find that isolates are, as expected,
140    genetically close to their matched general population in principal component
141    analyses (PCA), ADMIXTURE[22] and TreeMix[23] using common variants (Fig. 1b,
142    Supplementary Figs. 2-5 and Supplementary Table 6), but PCA using rare and low-
143    frequency variants, as found previously[24], distinguishes them more clearly from the
144    general population and also from other isolates, particularly among the Italian
145    samples (Fig. 1c, Supplementary Fig. 2). The majority of sharing of variants present
146    just twice across all samples of 36 individuals from each population ($f_2$ variants[21])
147    takes place within the same population, and the isolates generally share more with
148    their closest general population than with other populations. This latter trend,
149    however, is not apparent for IF1-IF4, who show little sharing with any other

4

150 population, pointing to a greater level of isolation and lower level of gene flow with
151 their general population (Fig. 1d, lower triangle and Supplementary Fig. 7), which is
152 confirmed by *f3*-statistics[25] comparing with a worldwide population panel of HGDP-
153 CEPH samples using common SNPs (Supplementary Fig. 6). $f_3$-$f_{10}$ variant sharing
154 demonstrates sharing by ITG and IVB with both Greek and UK populations (Fig. 1d,
155 upper triangle and Supplementary Fig. 7), potentially indicative of their more
156 ancient heritage.
157
158 **Population demographic history.** All populations studied here, both isolates and
159 general, appear to have shared a comparable effective population size (*Ne*)history
160 before 20 thousand years ago (KYA) based on the multiple sequentially Markovian
161 coalescent (MSMC) method[26] (Supplementary Fig. 9). The isolates diverged from
162 their general populations within the last ~5000 years based on LD estimations[27]
163 (Supplementary Table 7 and Supplementary Fig. 8) and yet had sharp decreases in
164 their population sizes in more recent times as estimated using inferred long
165 segments of identity by descent (IBD)[28] (Fig. 1e, f and Supplementary Fig.10).
166 Different isolates also split from their respective general populations at different
167 times. For example, IF1-IF4 split from ITG ~4-5 KYA, while most other isolates split
168 from their general populations within the last ~1,000 years (Supplementary Table
169 7).
170
171 The different demographic histories of different isolates should lead to different
172 genetic characteristics. To summarize these features in a single quantitative
173 measure that can be calculated from genotype data, as well as sequence data, we
174 developed an isolation index (*Isx*) which combines information on the divergence
175 time from the general population (*Tdg*), *Ne* and migration rate (*M*), such that early-
176 divergence-time isolates with small *Ne* and low *M* have a high *Isx* value (Fig. 2a and
177 Supplementary Fig. 11). The different isolates show different *Isx* values: IF1, IF2, IF3
178 and IF4 have the highest, while IVB has the lowest (Supplementary Table 8). *Isx*
179 values are highly correlated with other population-genetic characteristics (e.g. Fig.
180 2b, c, Supplementary Table 11), such as genome-wide pairwise $F_{ST}$ between isolates
181 and their matching general population (reflecting the genetic drift of the isolates)
182 (Supplementary Fig. 12), the total length and number of runs of homozygosity
183 (ROH) (Supplementary Fig. 13), inbreeding coefficient (F) (Supplementary Fig. 14)
184 and length of LD (Supplementary Figs. 15-16 and Supplementary Table 9, 10). All
185 these characteristics are correlated, but the pairwise correlation coefficients show
186 that *Isx* is a slightly better overall predictor of the other measures than any single
187 existing measure (Fig. 2c, Supplementary Fig. 17 and Supplementary Table 11);
188 moreover, it is potentially more robust to confounding factors as it is calculated
189 from three demographic parameters, while the others are all based on single
190 measurements.
191
192 **Purifying selection analyses.** Several lines of evidence suggest relaxed purifying
193 selection in the isolates due to their reduced *Ne*, although as expected we do not
194 detect substantially increased genetic load per genome using the *Rxy* statistic[29]
195 based on all of the variants in the genomes (Fig. 3a and Supplementary Table 12).

196 First, we see different levels of enrichment of low-frequency functional variants in
197 isolates (Fig. 3b and c, Supplementary Tables 13 and14, Supplementary Figs. 18a)
198 quantified by a new statistic, *DVxy-coding*, developed here (DV: drifted variants).
199 *DVxy-coding* measures the ratio of functional coding variants (missense plus loss-of-
200 function (LoF)) in isolates compared to the closest general population (and vice-
201 versa), adjusted for the corresponding ratios of intergenic variants in order to
202 correct for the effect of genetic drift. We applied this only to a subclass of DVs,
203 defined as low-frequency (2-5%, the best choice according to the sample size we
204 have) in any isolate, yet at least three-fold higher than in the closest general
205 population (and vice versa). We find that *DVxy-coding* is >1 in all isolates and <1 in
206 all general populations (Fig. 3c, Supplementary Fig. 18a and Supplementary Table
207 13). We also calculated a similar *DVxy-wg* statistic by stratifying whole-genome
208 variants according to their combined annotation dependent depletion (CADD) score
209 (0-5, neutral variants; 5-10, mildly deleterious; 10-20, deleterious; and >20, highly
210 deleterious; these cut-off choices balance the number of variants in each bin to allow
211 us comparable statistical power among all bins, although the conclusions are robust
212 to the particular cut-off values chosen and different bins (Supplementary Figs. 18b
213 and Supplementary Fig. 19)). The *DVxy-wg* values are differentiated for variants
214 with CADD score of 10-20 and significantly so (assessed using the jack-knife
215 bootstrap method) for ones with CADD scores >20, with *DVxy-wg* values >1 in all
216 isolates and <1 in all general populations (Fig. 3b, Supplementary Fig. 18b and
217 Supplementary Table 14). This demonstrates enrichment of low-frequency
218 functional variants, both coding and genome-wide with CADD score >10, in the
219 isolated populations. Moreover, both *DVxy-coding* and *DVxy-wg* values are
220 correlated with *Isx*, suggesting that different isolation characteristics lead to
221 different levels of enrichment of functional variants.
222
223 We also investigated the relaxation of purifying selection by assessing functional
224 (missense) singleton variants (SV) pooled for all of the genes that have at least one
225 singleton missense or synonymous variant in a pair of populations (one isolate and
226 its general population), correcting with pooled synonymous variants (*SVxy* statistic,).
227 We find a substantial deviation from 1 for functional singletons in all of the isolates
228 (Fig. 3d and Supplementary Table 15), with *SVxy* values positively correlating with
229 *Isx* (Fig. 2c and Supplementary Fig. 20). We also find that the proportion of relaxed
230 essential genes[30] with *SVxy* >1 in isolates is significantly higher than in the general
231 population (Supplementary Table 15). Such rare and low-frequency drifted
232 functional variants, measured by both DVxy and *DVxy*, are particularly relevant for
233 boosting the power of association studies[6].
234
235 **Positive selection analyses.** We do not find convincing evidence for positive
236 selection in any isolate using deltaDAF[31], PCAdapt[32] or SDS[33], although we do
237 identify some highly differentiated variants (Supplementary Fig. 21 and
238 Supplementary Tables 16,17), including in the protein-coding genes *ALK*, *SPNS2*,
239 *SLC39A11* and *ACSS2*, which can nevertheless be accounted for by drift.
240 Interestingly, we also find six highly-differentiated variants shared between
241 different isolates from Italy, IF2, IF3 and IF4, but interpret them as likely to result

6

242  from drift or positive selection for the ancestral allele in the ITG (Supplementary
243  Table 17). We find that the SDS method has little power in our samples because of
244  their small size, and failed to detect selection even at the lactose tolerance SNP in
245  the UKO, a known strong signal of recent selection (Supplementary Fig. 22).
246
247

## Discussion

249
250  Isolated populations have special characteristics that can be leveraged to increase
251  the power of association studies, as several previous studies have shown[19,34].
252  Nevertheless, only a small proportion of functional variants have increased in
253  frequency in any one isolate, so multiple isolates must be investigated to reveal the
254  full diversity of associated variants. Here, we probed an extended allele frequency
255  spectrum of variants potentially underpinning human complex disease through the
256  analysis of whole-genome sequence data in multiple isolates matched to nearby
257  non-isolated populations, capturing common, low-frequency and rare variants. We
258  quantified different levels of isolation resulting from different demographic
259  histories and have demonstrated that the *Isx* statistic, calculated even from SNP-chip
260  data, reliably captures these relevant features. This study provides a systematic
261  evaluation of the genetic characteristics of multiple European isolates and for the
262  first time empirically demonstrates enrichment of rare functional variants across
263  multiple isolates. With the advent of large-scale whole-genome sequencing, studies
264  in isolates are poised to continue as major contributors to our understanding of
265  complex disease etiology.
266
267

## Methods

269
270  **Dataset and variant calling:** The dataset includes 3059 whole-genome low-depth
271  sequences generated at The Wellcome Trust Sanger Institute using the Illumina
272  Genome Analyzer II and Illumina HiSeq 2000 platforms, as well as 100 high-depth
273  sequences from the Illumina HiSeq X Ten (Fig. 1a and Supplementary Table 1).
274  Informed consent was obtained from all subjects and the study was approved by the
275  HMDMC (Human Materials and Data Management Committee) of the Welcome Trust
276  Sanger Institute. The multi-sample genotype calling across all of the low-coverage
277  sequencing data from the current study, as well as 2353 from the 1000 Genomes
278  Project Phase 3 release, and 3781 from UK10K (a total of 9375) was performed with
279  the defined site selection criteria (Supplementary Note). Genotype likelihoods were
280  calculated with samtools/bcftools (0.2.0-rc9) and then genotypes were called and
281  phased using Beagle v4 (r1274)[35]. We assessed the performance of the genotype
282  calling from the low coverage data using the available genotype chip data for a
283  subset of the cohorts consisting of 4665 individuals, and calculated the discordance
284  rates on chromosome 20 separately for the categories REF-REF, REF-ALT and ALT-
285  ALT.
286

287 The sample sizes are very different across these collections, and we used three
288 different standard-sized subsets of the samples for different analyses: (1) the whole
289 dataset; (2) the sample-size-matched dataset, obtained either by randomly selecting
290 samples from general population to match the isolated population (for example, we
291 randomly select 377 from FIG to match FIK), or by randomly selecting a subset of
292 the isolated population to match the general population (for example, we randomly
293 select 108 IVB to match the general population ITG); (3) the minimum-sample-size
294 dataset of 36 individuals per population. By doing this, we maximize the use of the
295 data for different analyses, and we specify which dataset is used for each analysis.
296 The sequencing depth is also different across different populations, within a 2.5-fold
297 range (apart from GRG, in which variants were called differently, details in
298 Supplementary Notes), and we allowed for these differences when interpreting the
299 results.
300
301 **Variant counts:** We first re-annotated all variants using the Variant Effect Predictor
302 (VEP) annotation from Ensembl 76 with the "- pick" option, which gives one
303 annotation per variant. We then performed variant counting at both the population
304 and individual level, stratifying by functional categories and frequency bins. These
305 counts were either plotted in figures or summarized as median values in tables. We
306 carried out these analyses using both the sample-size-matched dataset and the
307 minimum-sample-size dataset.
308
309 **Population-genetic analyses:** We used the whole dataset for the analyses in this
310 section, unless otherwise specified. Principal component analyses (PCAs) were
311 performed separately with common variants or rare variants using EIGENSTRAT
312 v.501[36]. Shared ancestry between the populations studied here was evaluated using
313 ADMIXTURE v1.22[22]. The relationships between the populations studied here,
314 combined with worldwide populations from the HGDP-CEPH panel[37], were also
315 examined using ancestry graph analyses implemented in TreeMix v.1.12[23]. We also
316 used formal test of f3-satitisitcs[25] to investigate population mixture in the history of
317 the populations studied here, as well as worldwide populations from the HGDP-
318 CEPH panel. Rare $f_2$ variants (with only two copies of the alternative allele in the
319 minimum-sample-size dataset) and moderately rare $f_{3-10}$ variants (3-10 copies of the
320 alternative allele in the same dataset) are particularly informative for investigating
321 recent human history[21]. We investigated the sharing pattern of these two types of
322 variant by summing all $f_2$ variants or any random two alleles of the $f_{3-10}$ variants
323 shared by pairs of individuals. We plotted the results as a heat map using the image[1]
324 function from the base R package (https://stat.ethz.ch/R-manual/R-
325 devel/library/graphics/html/image.html). Variants were aggregated by pair of
326 individuals using the 'count' function of the plyr package, then arranged in matrix
327 form and colorized using 'colorRampPalette' from the colorspace package
328 (https://cran.r-project.org/web/packages/colorspace/index.html). Runs of
329 homozygosity (ROH),  inbreeding coefficient (F) as well as the length of LD-blocks
330 were calculated in PLINK, and finally genome-wide $F_{ST}$ values between isolates and
331 their general populations were calculated with the software 4P[38] using the
332 minimum-sample-size dataset.

333
334 **Demographic inferences:** LD-based[39-41] demographic inference was performed in
335 the NeON R package[27] using the minimum-sample-size dataset; the median and
336 confidence interval were estimated using the 50th, 5th and 95th percentiles of the
337 distribution of long-term *Ne* in each time interval. We used the multiple sequentially
338 Markovian coalescent (MSMC) method[26] to infer demographic changes before
339 20,000 years ago using four individual sequences from each population. In order to
340 account for some loss of heterozygous sites in the low-depth data, we used a slow
341 mutation rate of $0.8 \times 10^{-8}$ mutations per nucleotide per generation and a longer
342 generation time of 33 years. We then estimated more recent demographic changes
343 (from the present to ~9,000 years ago) using IBDNe[28] with the minimum-sample-
344 size dataset. We used IBDseq[42] to detect IBD segments in sequence data from
345 chromosome 2 in all populations. We then used IBDNe with the default parameters
346 and a minimum IBD segment length of 2 centiMorgan (cM) units. We assumed a
347 generation time of 29 years.
348
349 **Isolation index:** In order to quantify the different isolation levels of different
350 isolates, we developed an index that combines three demographic parameters: (a)
351 *Tdg*, (b) *Ne*, and (c) the level of private isolate ancestry (*M*). We call this estimate the
352 Isolation index (*Isx*). It is defined as:
353

$$Isx = \frac{\log(Tdg(100 * M)^2)}{\log(Ne)}$$

354
355 Both *Tdg* and *Ne* were inferred from the LD-based method using the NeON R
356 package[27]. *M* is difficult to estimate directly from SNP genotype data, so here we
357 estimated the difference of shared ancestral components between an isolate and its
358 general population from ADMIXTURE analysis. We ran ADMIXTURE with only one
359 isolate and it closest general population using K=2. We then estimated the difference
360 in the means of ancestry between the isolate and its general population. The *M*
361 parameter was defined as Delta Ancestry.
362
363 **_Rxy_ analysis:** Rxy statistics[29] between each pair of populations (an isolate and its
364 closest general population) for different functional categories were calculated using
365 the matched-sample-size data for missense and LoF variants, including stop gain,
366 splice donor and acceptor variants, using synonymous variants as controls (we did
367 not use intragenic variants as control because of the ascertainment in the ITG which
368 has high-depth exome sequences and low depth for the rest of the genome). We also
369 calculated *Rxy* statistics for variants with CADD scores[43] greater than 10 and 20,
370 using variants with CADD scores less than 5 as controls. The mean and standard
371 deviation for each Rxy value were obtained from 100 bootstraps.
372
373 **_DVxy_ analysis:** A new statistic, *DVxy*, was developed to quantify the enrichment of
374 low-frequency functional variants in the isolates using both the matched-sample-
375 size and minimum-sample-size datasets. It calculates the proportion of functional

9

376  variants in each isolate compared with its general population, correcting for genetic
377  drift at the same time. We calculated $DVxy$ specifically for the subset of variants with
378  DAF 2-5% in the isolate, and at least three times lower in its closest general
379  population, or vice-versa. We called these variants "drifted variants" (DV). $DVxy$ was
380  calculated for both coding regions and whole genomes.
381
382  For coding variants, we defined missense or missense plus LoF variants as
383  functional variants.  We counted the number of functional DVs and neutral
384  (intergenic) DVs in each isolate (population $x$) and the corresponding general
385  population (population $y$). The ratio between the fraction of DV variants from the
386  isolated population (corrected by the count of intergenic variants) and the
387  corresponding fraction of DV variants from its general population was defined as
388  the $DVxy$ statistic. If $DVxy$ is equal to 1, there is no enrichment for the functional DVs
389  in the isolate; less than 1 indicates depletion, and greater than 1 indicates
390  enrichment.
391

$$DVxy\_coding = \frac{\dfrac{\%DVx\ missense}{\%DVx\ intergenic}}{\dfrac{\%DVy\ missense}{\%DVy\ intergenic}}$$

392
393  For the whole genome, we used different CADD score cut-offs and bins. We
394  calculated a DV statistic by stratifying the variants according to their CADD scores
395  (0-5, neutral variants; 5-10, mildly deleterious; 10-20, deleterious; and greater than
396  20, highly deleterious) for each isolate and its closest general population. We finally
397  calculated a ratio of the fraction of DV variants (from each class) between the isolate
398  and its general population, and vice-versa. The following formula shows the $DVxy$-
399  $wg$ calculation for variants with CADD score between $i$ (isolate) and $j$ (general
400  population).
401

$$DVxy_{CADD(ij)} = \frac{\%\ DVx\ (CADD\ i-j)}{\%DVy(CADD\ i-j)}$$

402
403  The 95% confidence interval for each calculation was obtained by randomly
404  sampling data from 20 chromosomes 100 times.
405
406  **_SVxy_ analysis:** We further investigated the relaxation of purifying selection in the
407  isolated populations using singleton variants. Here, we also used the minimum-
408  sample-size dataset. Another new statistic, $SVxy$, was developed to measure the ratio
409  of missense vs synonymous singletons per gene in each population, as well as the
410  ratio of the sum of singletons in all genes which have at least one singleton in the
411  pair of the populations (one isolate and one general population). We counted the
412  number of missense singletons and synonymous singletons per gene in each
413  population, and $SVgene$ was calculated as:
414

415     *SVgene* = (SV missense count +1)/ (SV synonymous count +1)

416

417     *SVgene* >1 indicates relaxation of purifying selection; *SVgene* = 1 indicates

418     neutrality; and *SVgene* <1 indicates purifying selection.

419

420     We then divided the gene list into essential genes[30] and non-essential genes (the

421     rest), and calculated a statistic, $G_{SV}$, for each population, defined as:

422

423     $G_{SV}$ = percentage of essential genes with *SVgene* >1/percentage of non-essential

424     genes with *SVgene* >1

425

426     We finally calculated a statistic, *SVxy*, which is the ratio of *SVpop* of each isolate to

427     *SVpop* of its general population. *SVpop* for each isolate and its general population

428     was calculated using all genes which have at least one singleton in the pair of the

429     populations and defined as *SVpop* = Σ (SV missense counts)/Σ(SV synonymous

430     counts).

431

432     We used the same annotation as in the variant counts. We calculated a confidence

433     interval for each estimate using bootstrapping of 80% of the genes 100 times.

434

435     **Correlation analyses**: We calculated pair-wise correlation coefficients between the

436     *Isx* values, population-genetic measurements ROH, F, $F_{ST}$, and number and length of

437     LD blocks, as well as the newly-developed statistics *DVxy* and SVxy using the

438     Pearson correlation in R.

439

440     **Positive selection analyses:** We calculated genome-wide pairwise derived allele

441     frequency differences (deltaDAF) for each pair of populations (an isolate and its

442     general population) as described previously[31] using the matched-sample-size

443     dataset. We also carried out PCAdapt analyses[32] for each pair of populations using

444     the whole dataset. Both analyses look for high derived allele frequency variants in

445     the isolates, and will not be affected by sample size. Finally, we ran the singleton

446     density score (SDS) method[33] using the whole UKO and UKG datasets, which have

447     the largest sample sizes for both isolate and its general population, and thus the

448     greatest power for this method.

449

450     **Data availability:**

451

452     Amalgamated genotype calls across all populations studied are available through

453     the European Genome/Phenome Archive (EGAD00001002014) with Data Access

454     Agreement described in the Supplementary Information.

455

456

457

458

459

**References**

1.  Zeggini, E. Using genetically isolated populations to understand the genomic basis of disease. *Genome Med* **6**, 83 (2014).
2.  Hatzikotoulas, K., Gilly, A. & Zeggini, E. Using population isolates in genetic association studies. *Brief Funct Genomics* **13**, 371-7 (2014).
3.  Zuk, O. *et al.* Searching for missing heritability: designing rare variant association studies. *Proc Natl Acad Sci U S A* **111**, E455-64 (2014).
4.  Pollin, T.I. *et al.* A null mutation in human APOC3 confers a favorable plasma lipid profile and apparent cardioprotection. *Science* **322**, 1702-5 (2008).
5.  Gudmundsson, J. *et al.* A study based on whole-genome sequencing yields a rare variant at 8q24 associated with prostate cancer. *Nat Genet* **44**, 1326-9 (2012).
6.  Tachmazidou, I. *et al.* A rare functional cardioprotective APOC3 variant has risen in frequency in distinct population isolates. *Nat Commun* **4**, 2872 (2013).
7.  Huyghe, J.R. *et al.* Exome array analysis identifies new loci and low-frequency variants influencing insulin processing and secretion. *Nat Genet* **45**, 197-201 (2013).
8.  Li, A.H. *et al.* Analysis of loss-of-function variants and 20 risk factor phenotypes in 8,554 individuals identifies loci influencing chronic disease. *Nat Genet* **47**, 640-2 (2015).
9.  Lim, E.T. *et al.* Distribution and medical impact of loss-of-function variants in the Finnish founder population. *PLoS Genet* **10**, e1004494 (2014).
10. Moltke, I. *et al.* A common Greenlandic TBC1D4 variant confers muscle insulin resistance and type 2 diabetes. *Nature* **512**, 190-3 (2014).
11. Sidore, C. *et al.* Genome sequencing elucidates Sardinian genetic architecture and augments association analyses for lipid and blood inflammatory markers. *Nat Genet* **47**, 1272-81 (2015).
12. Steinthorsdottir, V. *et al.* Identification of low-frequency and rare sequence variants associated with elevated or reduced risk of type 2 diabetes. *Nat Genet* **46**, 294-8 (2014).
13. Panoutsopoulou, K. *et al.* Genetic characterization of Greek population isolates reveals strong genetic drift at missense and trait-associated variants. *Nat Commun* **5**, 5345 (2014).
14. Esko, T. *et al.* Genetic characterization of northeastern Italian population isolates in the context of broader European genetic diversity. *Eur J Hum Genet* **21**, 659-65 (2013).
15. Colonna, V. *et al.* Small effective population size and genetic homogeneity in the Val Borbera isolate. *Eur J Hum Genet* **21**, 89-94 (2013).

503    16.    Vitart, V. *et al.* SLC2A9 is a newly identified urate transporter influencing
504            serum urate concentration, urate excretion and gout. *Nat Genet* **40**, 437-42
505            (2008).
506    17.    The 1000 Genomes Project Consortium. A global reference for human genetic
507            variation. *Nature* **526**, 68-74 (2015).
508    18.    The UK10K Consortium. The UK10K project identifies rare variants in health
509            and disease. *Nature* **526**, 82-90 (2015).
510    19.    Gudbjartsson, D.F. *et al.* Large-scale whole-genome sequencing of the
511            Icelandic population. *Nat Genet* **47**, 435-44 (2015).
512    20.    McQuillan, R. *et al.* Runs of homozygosity in European populations. *Am J Hum
513            Genet* **83**, 359-72 (2008).
514    21.    Mathieson, I. & McVean, G. Demography and the age of rare variants. *PLoS
515            Genet* **10**, e1004528 (2014).
516    22.    Alexander, D.H., Novembre, J. & Lange, K. Fast model-based estimation of
517            ancestry in unrelated individuals. *Genome Res* **19**, 1655-64 (2009).
518    23.    Pickrell, J.K. & Pritchard, J.K. Inference of population splits and mixtures from
519            genome-wide allele frequency data. *PLoS Genet* **8**, e1002967 (2012).
520    24.    O'Connor, T.D. *et al.* Rare variation facilitates inferences of fine-scale
521            population structure in humans. *Mol Biol Evol* **32**, 653-60 (2015).
522    25.    Reich, D., Thangaraj, K., Patterson, N., Price, A.L. & Singh, L. Reconstructing
523            Indian population history. *Nature* **461**, 489-494 (2009).
524    26.    Schiffels, S. & Durbin, R. Inferring human population size and separation
525            history from multiple genome sequences. *Nat Genet* **46**, 919-25 (2014).
526    27.    Mezzavilla, M. & Ghirotto, S. Neon: An R package to estimate human effective
527            population size and divergence time from patterns of linkage disequilibrium
528            between SNPs. *J Comput Sci Syst Biol* **8**, 37-44 (2015).
529    28.    Browning, S.R. & Browning, B.L. Accurate non-parametric estimation of
530            recent effective population size from segments of identity by descent. *Am J
531            Hum Genet* **97**, 404-18 (2015).
532    29.    Do, R. *et al.* No evidence that selection has been less effective at removing
533            deleterious mutations in Europeans than in Africans. *Nat Genet* **47**, 126-31
534            (2015).
535    30.    Blomen, V.A. *et al.* Gene essentiality and synthetic lethality in haploid human
536            cells. *Science* **350**, 1092-6 (2015).
537    31.    Colonna, V. *et al.* Human genomic regions with exceptionally high levels of
538            population differentiation identified from 911 whole-genome sequences.
539            *Genome Biol* **15**, R88 (2014).
540    32.    Duforet-Frebourg, N., Bazin, E. & Blum, M.B.G. Genome scans for detecting
541            local adaptation using a Bayesian factor model. *Mol Biol Evol* **31**, 2483-2495
542            (2014).
543    33.    Field, Y. *et al.* Detection of human adaptation during the past 2000 years.
544            *Science* **354**, 760-764 (2016).
545    34.    Zoledziewska, M., Sidore, C. & Chiang, C.W. Height-reducing variants and
546            selection for short stature in Sardinia. *Nat Genet* **47**, 1352-1356 (2015).

547    35.    Browning, S.R. & Browning, B.L. Rapid and accurate haplotype phasing and
548          missing-data inference for whole-genome association studies by use of
549          localized haplotype clustering. *Am J Hum Genet* **81**, 1084-97 (2007).
550    36.    Price, A.L. *et al.* Principal components analysis corrects for stratification in
551          genome-wide association studies. *Nat Genet* **38**, 904-9 (2006).
552    37.    Li, J.Z. *et al.* Worldwide human relationships inferred from genome-wide
553          patterns of variation. *Science* **319**, 1100-4 (2008).
554    38.    Benazzo, A., Panziera, A. & Bertorelle, G. 4P: fast computing of population
555          genetics statistics from large DNA polymorphism panels. *Ecol Evol* **5**, 172-
556          175 (2014).
557    39.    Hill, W.G. Estimation of effective population size from data on linkage
558          disequilibrium. *Genetical Res.* **38**, 209-216 (1981).
559    40.    Hayes, B.J., Visscher, P.M., McPartlan, H.C. & Goddard, M.E. Novel multilocus
560          measure of linkage disequilibrium to estimate past effective population size.
561          *Genome Res* **13**, 635-43 (2003).
562    41.    Tenesa, A. *et al.* Recent human effective population size estimated from
563          linkage disequilibrium. *Genome Res* **17**, 520-6 (2007).
564    42.    Browning, B.L. & Browning, S.R. Detecting identity by descent and estimating
565          genotype error rates in sequence data. *Am J Hum Genet* **93**, 840-51 (2013).
566
567

## Author contributions

Y.X., C.T.-S., R.D. and E.Z.: design and supervision of the project. G.D., P.G., A.P., S.R., N.S., D.T. and J.F.W.:  population liaison, sampling and DNA provision. N.S., JF.W. and R.D: comments and approval of the manuscript on behalf of the population consortia. Y.X., M.M. and M.H.: statistical method development. M.M., M.H., S.M., V.N., A.G., Q.A., V.C., L.S., C.F., G.R., H.C. and P.R. and J.A.: population-genetic analyses, statistical analyses and data interpretation. Y.C. and A.M.: bioinformatics support. S.M., N.S. and R.D: data processing and QC. Y.X., M.M., M.H., S.M., V.C., C.T-.S. and E.Z.: manuscript drafting. All authors: approval of the final version of the manuscript.

## Competing financial interests

The authors declare no competing financial interests.

## Acknowledgements

**Figure legends**


Fig. 1. General characteristics and demographic history of isolated and matched general populations. a. Geographical locations of samples. The base map was plotted in R using the mapdata package and circles were added using Photoshop. b. PCA using common variants. c. PCA using low-frequency variants. d. Sharing of rare variants within and between populations. Upper left triangle: $f_2$ variants; lower right triangle $f_{3-f10}$ variants. e. Effective population size (Ne) inferred from IBDNe for UKO and UKG during the past 9 KY. f. The lowest Ne inferred by IBDNe for all populations for the past 3KY, plotted as a function of the time at which it occurred.

Fig. 2. Isolation index (*Isx*) and its correlation with other genetic measures. a. Information summarized in *Isx*. b. Example of the correlation between *Isx* and other statistics, here *DVxy-coding*. c. Summary of the correlations between *Isx* and other population-genetic statistics. All the correlation coefficients are high and statistically significant.

Fig. 3. Purifying selection in the isolates and general populations. a. *Rxy*-missense statistic in each isolate, showing no evidence for increased genetic load in the isolates. The mean and standard deviation for each Rxy value from 100 bootstraps are shown. b. *DVxy-wg* (*DVxy*-whole genome) statistic in isolates and general populations, stratified by CADD score, showing enrichment of highly-functional low-frequency variants. c. *DVxy-coding* statistic in isolates and general populations, showing enrichment of low-frequency missense variants in isolates. d. SVxy-missense statistic in each isolate, showing relaxation of purifying selection in isolates in singletons. The standard errors for both *DVxy* and *SVxy* were calculated by randomly sampling data from 20 chromosomes 100 times. All of these analyses are based on the minimum-sample-size dataset (36 individuals from each population).

624 **Tables:**

625

626 **Table 1. Summary of variants discovered in this study**

627

| POP | n | average depth | MAF ≤2% | | MAF >2-≤5% | | MAF >5% | | Novel common SNPs in isolate* | Novel common SNPs in isolate** |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | total | novel % | total | novel % | total | novel % | | |
| FIK | 377 | 4x | 4,066,373 | 10.90 | 1,553,076 | 1.20 | 6,025,077 | 0.70 | 190,527 | 70,579 |
| FIG | 1564 | 6x | 6,548,833 | 11.80 | 1,540,915 | 0.80 | 6,053,704 | 0.70 | n.a. | n.a. |
| GRM | 249 | 4x | 5,129,513 | 7.20 | 1,447,981 | 1.10 | 6,111,923 | 0.80 | 513,272 | 49,884 |
| GRG*** | 99 | 10-30x | 3,757,110 | n.a. | 1,321,955 | n.a. | 5,842,537 | n.a. | n.a. | n.a. |
| IF1 | 60 | 4-10x | 1,456,881 | 1.30 | 1,420,929 | 1.30 | 5,890,714 | 0.80 | 320,191 | 119,157 |
| IF2 | 45 | 4-10x | 1,063,098 | 1.30 | 1,554,145 | 1.00 | 6,001,568 | 0.80 | 273,694 | 94,496 |
| IF3 | 47 | 4-10x | 961,059 | 1.30 | 1,455,284 | 1.10 | 6,068,304 | 0.80 | 299,603 | 107,281 |
| IF4 | 36 | 4-10x | 1,030,673 | 1.30 | 1,124,789 | 1.10 | 6,001,625 | 0.80 | 308,356 | 122,254 |
| IVB | 222 | 6x | 4,857,767 | 1.60 | 1,396,799 | 0.80 | 6,112,476 | 0.80 | 188,972 | 30,284 |
| UKO | 397 | 4x | 5,963,416 | 11.70 | 1,471,782 | 0.80 | 6,047,383 | 0.80 | 193,300 | 36,512 |
| Total | 3096 | | 12,218,797 | 10.50 | 5,503,179 | 0.70 | 8,301,524 | 0.30 | | |

628

629 'Novel' variants are those not found in 1000 Genomes Project Phase 3 or UK10K

630 project. *Variants that are common (minor allele frequency, MAF ≥5.6%, alternative

631 allele count ≥ 4) in an isolated population but not common (MAF <1.4%, alternative

632 allele count ≤ 1) in its closest general population. **Variants that are common (MAF

633 ≥5.6%, alternative allele count ≥ 4) in an isolated population but not (MAF <1.4%,

634 alternative allele count ≤ 1) in *any* of the general populations. ***Different variant

635 calling procedure in this population.

636

637

638

639

**Enrichment of low-frequency functional variants revealed by whole-genome sequencing of multiple isolated European populations**

Yali Xue[1*], Massimo Mezzavilla[1,2*], Marc Haber[1*], Shane McCarthy[1*], Yuan Chen[1], Vagheesh Narasimhan[1], Arthur Gilly[1], Qasim Ayub[1], Vincenza Colonna[1,3], Lorraine Southam[1,4], Christopher Finan[1], Andrea Massaia[1,5], Himanshu Chheda[6], Priit Palta[6,7], Graham Ritchie[1,8,9], Jennifer Asimit[1], George Dedoussis[10], Paolo Gasparini[11], Aarno Palotie[1,6,12-16], Samuli Ripatti[1,6,17], Nicole Soranzo[1,18], Daniela Toniolo[19], James F. Wilson[9,20], Richard Durbin[1], Chris Tyler-Smith[1], Eleftheria Zeggini[1]

[1]The Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambs. CB10 1SA, UK.
[2]Institute for Maternal and Child Health, IRCCS Burlo Garofolo, University of Trieste, 34137 Trieste, Italy.
[3]Consiglio Nazionale delle Ricerche, Istituto di Genetica e Biofisica "Adriano Buzzati-Traverso", via Pietro Castellino 111, 80131 Napoli, Italy.
[4]Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK.
[5] National Heart and Lung Institute, Imperial College London, London SW7 2AZ, UK
[6]Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Tukholmankatu 8, 00290 Helsinki, Finland.
[7]Estonian Genome Center, University of Tartu, 23B Riia Street, 51010 Tartu, Estonia.
[8]European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambs. CB10 1SD, UK.
[9]MRC Human Genetics Unit, MRC IGMM, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh EH4 2XU, UK.
[10]Department of Nutrition and Dietetics, Harokopio University Athens, Athens, Eleftheriou Venizelou 70, Kallithea 176 76, Greece.
[11]Medical Genetics, DSM, University of Trieste and IRCCS (Istituto di Ricovero e Cura a Carattere Scientifico) Burlo Garofolo Children Hospital, Via dell'Istria, 65, 34137, Trieste, Italy.
[12]Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA 02114, USA.
[13]Program in Medical and Population Genetics, The Broad Institute of MIT and Harvard, Cambridge, MA 02114, USA.
[14]The Stanley Center for Psychiatric Research, The Broad Institute of MIT and Harvard, Cambridge, MA 02114, USA.
[15]Psychiatric & Neurodevelopmental Genetics Unit, Department of Psychiatry, Massachusetts General Hospital, Boston, MA 02114, USA.
[16]Department of Neurology, Massachusetts General Hospital, Boston, MA 02114, USA.
[17]Department of Public Health, University of Helsinki, Helsinki FI-00014, Finland.
[18]Department of Haematology, University of Cambridge, Cambridge CB2 0XY, UK.

46 [19]Division of Genetics and Cell Biology, San Raffaele Scientific Institute, via Olgettina
47 60, 20132, Milan, Italy.
48 [20]Usher Institute of Population Health Sciences and Informatics, University of
49 Edinburgh, Teviot Place, Edinburgh EH8 9AG, Scotland, UK.
50
51 * Equal contribution to the work
52
53 Correspondence should be addressed to Y. X. (ylx@sanger.ac.uk) or E. Z.
54 (Eleftheria@sanger.ac.uk).
55
58

## Abstract

**The genetic features of isolated populations can boost power in complex-trait association studies, and an in-depth understanding of how their genetic variation has been shaped by their demographic history can help leverage these advantageous characteristics. Here, we perform a comprehensive investigation using 3059 newly-generated low-depth whole-genome sequences from eight European isolates and two matched general populations, together with published data from the 1000 Genomes Project and UK10K. Sequencing data give deeper and richer insights into population demography and genetic characteristics than genotype-chip data, distinguishing related populations more effectively and allowing their functional variants to be studied more fully. We demonstrate relaxation of purifying selection in the isolates, leading to enrichment of rare and low-frequency functional variants, using novel statistics, *DVxy* and *SVxy*. We also develop an isolation-index (*Isx*) that predicts the overall level of such key genetic characteristics and can thus help guide population choice in future complex-trait association studies.**

## Introduction

Population variation in disease susceptibility has been shaped by environment, demography and evolutionary history. Isolated populations (isolates) have generally experienced bottlenecks and strong genetic drift, so by chance some deleterious rare variants have increased in frequency while some neutral rare variation is lost, both helpful characteristics for the discovery of novel rare variant signals underpinning complex traits[1-3]. Studies to date have focused on individual isolates and have identified several disease-associated signals[4-12]. However, isolates differ in the time when they became isolated, their initial population size, the level of gene flow from outside and other historical demographic factors, and consequently also differ in their power for association studies[2]. We thus generate and analyze low-depth (4x-10x) whole-genome sequences (WGS) from eight cohorts drawn from isolated European populations and compare each isolate with the closest non-isolated (general) population, for which we also generate or access WGS data. We then investigate empirically how these historical differences influence the population-genetic properties of isolates, and frame these insights in terms of their consequences for study design in complex trait association studies.

## Results

**Samples, sequencing and QC.** The dataset includes newly-generated low-depth (4x-10x) WGS from eight cohorts drawn from isolated European populations: one each from Kuusamo in Finland (FIK) and Crete in Greece (GRM[13]), four from Friuli-Venezia Giulia in Italy (IF1, IF2, IF3 and IF4[14]), and one each from Val Borbera in

3

104    Italy (IVB[15]) and the Orkney Islands in the UK (UKO[16]); and the closest non-isolated
105    (general) population: Finland (FIG[9]), Greece (GRG), together with publicly available
106    data for Italy (ITG[17]) and UK (UKG[18]) (Fig. 1a and Supplementary Table 1). We
107    generated a superset of variants called in these cohorts and all 26 population
108    samples in the 1000 Genomes Project Phase 3[17], and performed multi-sample
109    genotype calling across all 9375 samples (3059 from the current study, 2353 from
110    the 1000 Genomes Project Phase 3 release, and 3781 from UK10K). Both individual
111    population and amalgamated genotype call data, which have greater than 99%
112    concordance with genotyping data (Supplementary Table 2), are available to the
113    scientific community (Data availability).
114
115    **General description of the variants in the isolates.** We identified approximately
116    12.2 million (M) variants with minor allele frequency (MAF) ≤2% (rare), 5.5M with
117    MAF >2-≤5% (low-frequency) and 8.3M variants with MAF >5% (common) across
118    the 10 populations newly sequenced here (eight isolates, GRG and FIG). Of these,
119    10.5%, 0.7% and 0.3%, respectively, are novel (Table 1 and Supplementary Table 3).
120    As expected, most of the isolates have lower numbers of variant sites per genome
121    than their closest general population (Supplementary Fig.1, Supplementary Table 5).
122    We find ~188,000-~513,000 variants that are common with MAF >5.5% in each
123    isolate but with MAF ≤1.4% in its closest general population (Table 1); ~30,000-
124    122,000 of these per isolate have frequency ≤1.4% in all the general samples
125    studied, among which ~150-~700 in coding regions and ~500-~2800 genome-wide
126    are deleterious (Supplementary Table 4).  These common and low-frequency
127    variants are thus useful markers for whole-genome association studies in these
128    populations and some of them (if absent from the general population) could
129    potentially lead to novel association signals. They include known examples such as
130    rs76353203 (R19X) in *APOC3* in GRM, which is associated with high-density
131    lipoprotein and triglyceride levels[6].
132
133    **Population-genetic analyses in the isolates.** Previous population-genetic studies
134    of isolates have, with some exceptions[11,19], been based on common variants found
135    on genotyping arrays, and have illustrated general characteristics such as low
136    genetic diversity and longer shared haplotypes[9,13-15,19,20]. Rare variants discovered
137    from sequencing are on average more recent in origin than common variants[21] and
138    therefore more powerful for distinguishing closely-related populations and more
139    informative about recent demographic history. We find that isolates are, as expected,
140    genetically close to their matched general population in principal component
141    analyses (PCA), ADMIXTURE[22] and TreeMix[23] using common variants (Fig. 1b,
142    Supplementary Figs. 2-5 and Supplementary Table 6), but PCA using rare and low-
143    frequency variants, as found previously[24], distinguishes them more clearly from the
144    general population and also from other isolates, particularly among the Italian
145    samples (Fig. 1c, Supplementary Fig. 2). The majority of sharing of variants present
146    just twice across all samples of 36 individuals from each population ($f_2$ variants[21])
147    takes place within the same population, and the isolates generally share more with
148    their closest general population than with other populations. This latter trend,
149    however, is not apparent for IF1-IF4, who show little sharing with any other

4

150   population, pointing to a greater level of isolation and lower level of gene flow with
151   their general population (Fig. 1d, lower triangle and Supplementary Fig. 7), which is
152   confirmed by f3-statistics[25] comparing with a worldwide population panel of HGDP-
153   CEPH samples using common SNPs (Supplementary Fig. 6). $f_3$-$f_{10}$ variant sharing
154   demonstrates sharing by ITG and IVB with both Greek and UK populations (Fig. 1d,
155   upper triangle and Supplementary Fig. 7), potentially indicative of their more
156   ancient heritage.
157
158   **Population demographic history.** All populations studied here, both isolates and
159   general, appear to have shared a comparable effective population size (*Ne*)history
160   before 20 thousand years ago (KYA) based on the multiple sequentially Markovian
161   coalescent (MSMC) method[26] (Supplementary Fig. 9). The isolates diverged from
162   their general populations within the last ~5000 years based on LD estimations[27]
163   (Supplementary Table 7 and Supplementary Fig. 8) and yet had sharp decreases in
164   their population sizes in more recent times as estimated using inferred long
165   segments of identity by descent (IBD)[28] (Fig. 1e, f and Supplementary Fig.10).
166   Different isolates also split from their respective general populations at different
167   times. For example, IF1-IF4 split from ITG ~4-5 KYA, while most other isolates split
168   from their general populations within the last ~1,000 years (Supplementary Table
169   7).
170
171   The different demographic histories of different isolates should lead to different
172   genetic characteristics. To summarize these features in a single quantitative
173   measure that can be calculated from genotype data, as well as sequence data, we
174   developed an isolation index (*Isx*) which combines information on the divergence
175   time from the general population (*Tdg*), *Ne* and migration rate (*M*), such that early-
176   divergence-time isolates with small *Ne* and low *M* have a high *Isx* value (Fig. 2a and
177   Supplementary Fig. 11). The different isolates show different *Isx* values: IF1, IF2, IF3
178   and IF4 have the highest, while IVB has the lowest (Supplementary Table 8). *Isx*
179   values are highly correlated with other population-genetic characteristics (e.g. Fig.
180   2b, c, Supplementary Table 11), such as genome-wide pairwise $F_{ST}$ between isolates
181   and their matching general population (reflecting the genetic drift of the isolates)
182   (Supplementary Fig. 12), the total length and number of runs of homozygosity
183   (ROH) (Supplementary Fig. 13), inbreeding coefficient (F) (Supplementary Fig. 14)
184   and length of LD (Supplementary Figs. 15-16 and Supplementary Table 9, 10). All
185   these characteristics are correlated, but the pairwise correlation coefficients show
186   that *Isx* is a slightly better overall predictor of the other measures than any single
187   existing measure (Fig. 2c, Supplementary Fig. 17 and Supplementary Table 11);
188   moreover, it is potentially more robust to confounding factors as it is calculated
189   from three demographic parameters, while the others are all based on single
190   measurements.
191
192   **Purifying selection analyses.** Several lines of evidence suggest relaxed purifying
193   selection in the isolates due to their reduced *Ne*, although as expected we do not
194   detect substantially increased genetic load per genome using the *Rxy* statistic[29]
195   based on all of the variants in the genomes (Fig. 3a and Supplementary Table 12).

196    First, we see different levels of enrichment of low-frequency functional variants in
197    isolates (Fig. 3b and c, Supplementary Tables 13 and14, Supplementary Figs. 18a)
198    quantified by a new statistic, *DVxy-coding*, developed here (DV: drifted variants).
199    *DVxy-coding* measures the ratio of functional coding variants (missense plus loss-of-
200    function (LoF)) in isolates compared to the closest general population (and vice-
201    versa), adjusted for the corresponding ratios of intergenic variants in order to
202    correct for the effect of genetic drift. We applied this only to a subclass of DVs,
203    defined as low-frequency (2-5%, the best choice according to the sample size we
204    have) in any isolate, yet at least three-fold higher than in the closest general
205    population (and vice versa). We find that *DVxy-coding* is >1 in all isolates and <1 in
206    all general populations (Fig. 3c, Supplementary Fig. 18a and Supplementary Table
207    13). We also calculated a similar *DVxy-wg* statistic by stratifying whole-genome
208    variants according to their combined annotation dependent depletion (CADD) score
209    (0-5, neutral variants; 5-10, mildly deleterious; 10-20, deleterious; and >20, highly
210    deleterious; these cut-off choices balance the number of variants in each bin to allow
211    us comparable statistical power among all bins, although the conclusions are robust
212    to the particular cut-off values chosen and different bins (Supplementary Figs. 18b
213    and Supplementary Fig. 19)). The *DVxy-wg* values are differentiated for variants
214    with CADD score of 10-20 and significantly so (assessed using the jack-knife
215    bootstrap method) for ones with CADD scores >20, with *DVxy-wg* values >1 in all
216    isolates and <1 in all general populations (Fig. 3b, Supplementary Fig. 18b and
217    Supplementary Table 14). This demonstrates enrichment of low-frequency
218    functional variants, both coding and genome-wide with CADD score >10, in the
219    isolated populations. Moreover, both *DVxy-coding* and *DVxy-wg* values are
220    correlated with *Isx*, suggesting that different isolation characteristics lead to
221    different levels of enrichment of functional variants.
222
223    We also investigated the relaxation of purifying selection by assessing functional
224    (missense) singleton variants (SV) pooled for all of the genes that have at least one
225    singleton missense or synonymous variant in a pair of populations (one isolate and
226    its general population), correcting with pooled synonymous variants (*SVxy* statistic,).
227    We find a substantial deviation from 1 for functional singletons in all of the isolates
228    (Fig. 3d and Supplementary Table 15), with *SVxy* values positively correlating with
229    *Isx* (Fig. 2c and Supplementary Fig. 20). We also find that the proportion of relaxed
230    essential genes[30] with *SVxy* >1 in isolates is significantly higher than in the general
231    population (Supplementary Table 15). Such rare and low-frequency drifted
232    functional variants, measured by both *SVxy* and *DVxy*, are particularly relevant for
233    boosting the power of association studies[6].
234
235    **Positive selection analyses.** We do not find convincing evidence for positive
236    selection in any isolate using deltaDAF[31], PCAdapt[32] or SDS[33], although we do
237    identify some highly differentiated variants (Supplementary Fig. 21 and
238    Supplementary Tables 16,17), including in the protein-coding genes *ALK*, *SPNS2*,
239    *SLC39A11* and *ACSS2*, which can nevertheless be accounted for by drift.
240    Interestingly, we also find six highly-differentiated variants shared between
241    different isolates from Italy, IF2, IF3 and IF4, but interpret them as likely to result

242  from drift or positive selection for the ancestral allele in the ITG (Supplementary
243  Table 17). We find that the SDS method has little power in our samples because of
244  their small size, and failed to detect selection even at the lactose tolerance SNP in
245  the UKO, a known strong signal of recent selection (Supplementary Fig. 22).
246
247

## Discussion

249

250  Isolated populations have special characteristics that can be leveraged to increase
251  the power of association studies, as several previous studies have shown[19,34].
252  Nevertheless, only a small proportion of functional variants have increased in
253  frequency in any one isolate, so multiple isolates must be investigated to reveal the
254  full diversity of associated variants. Here, we probed an extended allele frequency
255  spectrum of variants potentially underpinning human complex disease through the
256  analysis of whole-genome sequence data in multiple isolates matched to nearby
257  non-isolated populations, capturing common, low-frequency and rare variants. We
258  quantified different levels of isolation resulting from different demographic
259  histories and have demonstrated that the *Isx* statistic, calculated even from SNP-chip
260  data, reliably captures these relevant features. This study provides a systematic
261  evaluation of the genetic characteristics of multiple European isolates and for the
262  first time empirically demonstrates enrichment of rare functional variants across
263  multiple isolates. With the advent of large-scale whole-genome sequencing, studies
264  in isolates are poised to continue as major contributors to our understanding of
265  complex disease etiology.
266
267

## Methods

269

270  **Dataset and variant calling:** The dataset includes 3059 whole-genome low-depth
271  sequences generated at The Wellcome Trust Sanger Institute using the Illumina
272  Genome Analyzer II and Illumina HiSeq 2000 platforms, as well as 100 high-depth
273  sequences from the Illumina HiSeq X Ten (Fig. 1a and Supplementary Table 1).
274  Informed consent was obtained from all subjects and the study was approved by the
275  HMDMC (Human Materials and Data Management Committee) of the Welcome Trust
276  Sanger Institute. The multi-sample genotype calling across all of the low-coverage
277  sequencing data from the current study, as well as 2353 from the 1000 Genomes
278  Project Phase 3 release, and 3781 from UK10K (a total of 9375) was performed with
279  the defined site selection criteria (Supplementary Note). Genotype likelihoods were
280  calculated with samtools/bcftools (0.2.0-rc9) and then genotypes were called and
281  phased using Beagle v4 (r1274)[35]. We assessed the performance of the genotype
282  calling from the low coverage data using the available genotype chip data for a
283  subset of the cohorts consisting of 4665 individuals, and calculated the discordance
284  rates on chromosome 20 separately for the categories REF-REF, REF-ALT and ALT-
285  ALT.
286

287   The sample sizes are very different across these collections, and we used three
288   different standard-sized subsets of the samples for different analyses: (1) the whole
289   dataset; (2) the sample-size-matched dataset, obtained either by randomly selecting
290   samples from general population to match the isolated population (for example, we
291   randomly select 377 from FIG to match FIK), or by randomly selecting a subset of
292   the isolated population to match the general population (for example, we randomly
293   select 108 IVB to match the general population ITG); (3) the minimum-sample-size
294   dataset of 36 individuals per population. By doing this, we maximize the use of the
295   data for different analyses, and we specify which dataset is used for each analysis.
296   The sequencing depth is also different across different populations, within a 2.5-fold
297   range (apart from GRG, in which variants were called differently, details in
298   Supplementary Notes), and we allowed for these differences when interpreting the
299   results.
300
301   **Variant counts:** We first re-annotated all variants using the Variant Effect Predictor
302   (VEP) annotation from Ensembl 76 with the "- pick" option, which gives one
303   annotation per variant. We then performed variant counting at both the population
304   and individual level, stratifying by functional categories and frequency bins. These
305   counts were either plotted in figures or summarized as median values in tables. We
306   carried out these analyses using both the sample-size-matched dataset and the
307   minimum-sample-size dataset.
308
309   **Population-genetic analyses:** We used the whole dataset for the analyses in this
310   section, unless otherwise specified. Principal component analyses (PCAs) were
311   performed separately with common variants or rare variants using EIGENSTRAT
312   v.501[36]. Shared ancestry between the populations studied here was evaluated using
313   ADMIXTURE v1.22[22]. The relationships between the populations studied here,
314   combined with worldwide populations from the HGDP-CEPH panel[37], were also
315   examined using ancestry graph analyses implemented in TreeMix v.1.12[23]. We also
316   used formal test of f3-satitisitcs[25] to investigate population mixture in the history of
317   the populations studied here, as well as worldwide populations from the HGDP-
318   CEPH panel. Rare $f_2$ variants (with only two copies of the alternative allele in the
319   minimum-sample-size dataset) and moderately rare $f_{3-10}$ variants (3-10 copies of the
320   alternative allele in the same dataset) are particularly informative for investigating
321   recent human history[21]. We investigated the sharing pattern of these two types of
322   variant by summing all $f_2$ variants or any random two alleles of the $f_{3-10}$ variants
323   shared by pairs of individuals. We plotted the results as a heat map using the image[1]
324   function from the base R package (https://stat.ethz.ch/R-manual/R-
325   devel/library/graphics/html/image.html). Variants were aggregated by pair of
326   individuals using the 'count' function of the plyr package, then arranged in matrix
327   form and colorized using 'colorRampPalette' from the colorspace package
328   (https://cran.r-project.org/web/packages/colorspace/index.html). Runs of
329   homozygosity (ROH), inbreeding coefficient (F) as well as the length of LD-blocks
330   were calculated in PLINK, and finally genome-wide $F_{ST}$ values between isolates and
331   their general populations were calculated with the software 4P[38] using the
332   minimum-sample-size dataset.

333
334 **Demographic inferences:** LD-based[39-41] demographic inference was performed in
335 the NeON R package[27] using the minimum-sample-size dataset; the median and
336 confidence interval were estimated using the 50th, 5th and 95th percentiles of the
337 distribution of long-term *Ne* in each time interval. We used the multiple sequentially
338 Markovian coalescent (MSMC) method[26] to infer demographic changes before
339 20,000 years ago using four individual sequences from each population. In order to
340 account for some loss of heterozygous sites in the low-depth data, we used a slow
341 mutation rate of 0.8 x 10$^{-8}$ mutations per nucleotide per generation and a longer
342 generation time of 33 years. We then estimated more recent demographic changes
343 (from the present to ~9,000 years ago) using IBDNe[28] with the minimum-sample-
344 size dataset. We used IBDseq[42] to detect IBD segments in sequence data from
345 chromosome 2 in all populations. We then used IBDNe with the default parameters
346 and a minimum IBD segment length of 2 centiMorgan (cM) units. We assumed a
347 generation time of 29 years.
348
349 **Isolation index:** In order to quantify the different isolation levels of different
350 isolates, we developed an index that combines three demographic parameters: (a)
351 *Tdg*, (b) *Ne*, and (c) the level of private isolate ancestry (*M*). We call this estimate the
352 Isolation index (*Isx*). It is defined as:
353

$$Isx = \frac{\log(Tdg(100 * M)^2)}{\log(Ne)}$$

354
355 Both *Tdg* and *Ne* were inferred from the LD-based method using the NeON R
356 package[27]. *M* is difficult to estimate directly from SNP genotype data, so here we
357 estimated the difference of shared ancestral components between an isolate and its
358 general population from ADMIXTURE analysis. We ran ADMIXTURE with only one
359 isolate and it closest general population using K=2. We then estimated the difference
360 in the means of ancestry between the isolate and its general population. The *M*
361 parameter was defined as Delta Ancestry.
362
363 ***Rxy* analysis:** *Rxy* statistics[29] between each pair of populations (an isolate and its
364 closest general population) for different functional categories were calculated using
365 the matched-sample-size data for missense and LoF variants, including stop gain,
366 splice donor and acceptor variants, using synonymous variants as controls (we did
367 not use intragenic variants as control because of the ascertainment in the ITG which
368 has high-depth exome sequences and low depth for the rest of the genome). We also
369 calculated *Rxy* statistics for variants with CADD scores[43] greater than 10 and 20,
370 using variants with CADD scores less than 5 as controls. The mean and standard
371 deviation for each *Rxy* value were obtained from 100 bootstraps.
372
373 ***DVxy* analysis:** A new statistic, *DVxy*, was developed to quantify the enrichment of
374 low-frequency functional variants in the isolates using both the matched-sample-
375 size and minimum-sample-size datasets. It calculates the proportion of functional

376  variants in each isolate compared with its general population, correcting for genetic
377  drift at the same time. We calculated *DVxy* specifically for the subset of variants with
378  DAF 2-5% in the isolate, and at least three times lower in its closest general
379  population, or vice-versa. We called these variants "drifted variants" (DV). *DVxy* was
380  calculated for both coding regions and whole genomes.
381
382  For coding variants, we defined missense or missense plus LoF variants as
383  functional variants.  We counted the number of functional DVs and neutral
384  (intergenic) DVs in each isolate (population *x*) and the corresponding general
385  population (population *y*). The ratio between the fraction of DV variants from the
386  isolated population (corrected by the count of intergenic variants) and the
387  corresponding fraction of DV variants from its general population was defined as
388  the *DVxy* statistic. If *DVxy* is equal to 1, there is no enrichment for the functional DVs
389  in the isolate; less than 1 indicates depletion, and greater than 1 indicates
390  enrichment.
391

$$DVxy\_coding = \frac{\dfrac{\%DVx\ missense}{\%DVx\ intergenic}}{\dfrac{\%DVy\ missense}{\%DVy\ intergenic}}$$

392
393  For the whole genome, we used different CADD score cut-offs and bins. We
394  calculated a DV statistic by stratifying the variants according to their CADD scores
395  (0-5, neutral variants; 5-10, mildly deleterious; 10-20, deleterious; and greater than
396  20, highly deleterious) for each isolate and its closest general population. We finally
397  calculated a ratio of the fraction of DV variants (from each class) between the isolate
398  and its general population, and vice-versa. The following formula shows the *DVxy-*
399  *wg* calculation for variants with CADD score between *i* (isolate) and *j* (general
400  population).
401

$$DVxy_{CADD(ij)} = \frac{\%\ DVx\ (CADD\ i-j)}{\%DVy(CADD\ i-j)}$$

402
403  The 95% confidence interval for each calculation was obtained by randomly
404  sampling data from 20 chromosomes 100 times.
405
406  ***SVxy* analysis:** We further investigated the relaxation of purifying selection in the
407  isolated populations using singleton variants. Here, we also used the minimum-
408  sample-size dataset. Another new statistic, *SVxy*, was developed to measure the ratio
409  of missense vs synonymous singletons per gene in each population, as well as the
410  ratio of the sum of singletons in all genes which have at least one singleton in the
411  pair of the populations (one isolate and one general population). We counted the
412  number of missense singletons and synonymous singletons per gene in each
413  population, and *SVgene* was calculated as:
414

415   *SVgene* = (SV missense count +1)/ (SV synonymous count +1)
416
417   *SVgene* >1 indicates relaxation of purifying selection; *SVgene* = 1 indicates
418   neutrality; and *SVgene* <1 indicates purifying selection.
419
420   We then divided the gene list into essential genes[30] and non-essential genes (the
421   rest), and calculated a statistic, $G_{SV}$, for each population, defined as:
422
423   $G_{SV}$ = percentage of essential genes with *SVgene* >1/percentage of non-essential
424   genes with *SVgene* >1
425
426   We finally calculated a statistic, *SVxy*, which is the ratio of *SVpop* of each isolate to
427   *SVpop* of its general population. *SVpop* for each isolate and its general population
428   was calculated using all genes which have at least one singleton in the pair of the
429   populations and defined as *SVpop* = Σ (SV missense counts)/Σ(SV synonymous
430   counts).
431
432   We used the same annotation as in the variant counts. We calculated a confidence
433   interval for each estimate using bootstrapping of 80% of the genes 100 times.
434
435   **Correlation analyses**: We calculated pair-wise correlation coefficients between the
436   *Isx* values, population-genetic measurements ROH, F, $F_{ST}$, and number and length of
437   LD blocks, as well as the newly-developed statistics *DVxy* and *SVxy* using the
438   Pearson correlation in R.
439
440   **Positive selection analyses:** We calculated genome-wide pairwise derived allele
441   frequency differences (deltaDAF) for each pair of populations (an isolate and its
442   general population) as described previously[31] using the matched-sample-size
443   dataset. We also carried out PCAdapt analyses[32] for each pair of populations using
444   the whole dataset. Both analyses look for high derived allele frequency variants in
445   the isolates, and will not be affected by sample size. Finally, we ran the singleton
446   density score (SDS) method[33] using the whole UKO and UKG datasets, which have
447   the largest sample sizes for both isolate and its general population, and thus the
448   greatest power for this method.
449
450   **Data availability:**
451
452   Amalgamated genotype calls across all populations studied are available through
453   the European Genome/Phenome Archive (EGAD00001002014) with Data Access
454   Agreement described in the Supplementary Information.
455
456
457   **References**
458

459    1.    Zeggini, E. Using genetically isolated populations to understand the genomic
460          basis of disease. *Genome Med* **6**, 83 10.1186/s13073-014-0083-5 (2014).
461    2.    Hatzikotoulas, K., Gilly, A. & Zeggini, E. Using population isolates in genetic
462          association studies. *Brief Funct Genomics* **13**, 371-377 10.1093/bfgp/elu022
463          (2014).
464    3.    Zuk, O. *et al.* Searching for missing heritability: designing rare variant
465          association studies. *Proc Natl Acad Sci U S A* **111**, E455-464
466          10.1073/pnas.1322563111 (2014).
467    4.    Pollin, T.I. *et al.* A null mutation in human APOC3 confers a favorable plasma
468          lipid profile and apparent cardioprotection. *Science* **322**, 1702-1705
469          10.1126/science.1161524 (2008).
470    5.    Gudmundsson, J. *et al.* A study based on whole-genome sequencing yields a
471          rare variant at 8q24 associated with prostate cancer. *Nat Genet* **44**, 1326-
472          1329 10.1038/ng.2437 (2012).
473    6.    Tachmazidou, I. *et al.* A rare functional cardioprotective *APOC3* variant has
474          risen in frequency in distinct population isolates. *Nat Commun* **4**, 2872
475          10.1038/ncomms3872 (2013).
476    7.    Huyghe, J.R. *et al.* Exome array analysis identifies new loci and low-frequency
477          variants influencing insulin processing and secretion. *Nat Genet* **45**, 197-201
478          10.1038/ng.2507 (2013).
479    8.    Li, A.H. *et al.* Analysis of loss-of-function variants and 20 risk factor
480          phenotypes in 8,554 individuals identifies loci influencing chronic disease.
481          *Nat Genet* **47**, 640-642 10.1038/ng.3270 (2015).
482    9.    Lim, E.T. *et al.* Distribution and medical impact of loss-of-function variants in
483          the Finnish founder population. *PLoS Genet* **10**, e1004494
484          10.1371/journal.pgen.1004494 (2014).
485    10.   Moltke, I. *et al.* A common Greenlandic TBC1D4 variant confers muscle
486          insulin resistance and type 2 diabetes. *Nature* **512**, 190-193
487          10.1038/nature13425 (2014).
488    11.   Sidore, C. *et al.* Genome sequencing elucidates Sardinian genetic architecture
489          and augments association analyses for lipid and blood inflammatory markers.
490          *Nat Genet* **47**, 1272-1281 10.1038/ng.3368 (2015).
491    12.   Steinthorsdottir, V. *et al.* Identification of low-frequency and rare sequence
492          variants associated with elevated or reduced risk of type 2 diabetes. *Nat
493          Genet* **46**, 294-298 10.1038/ng.2882 (2014).
494    13.   Panoutsopoulou, K. *et al.* Genetic characterization of Greek population
495          isolates reveals strong genetic drift at missense and trait-associated variants.
496          *Nat Commun* **5**, 5345 10.1038/ncomms6345 (2014).
497    14.   Esko, T. *et al.* Genetic characterization of northeastern Italian population
498          isolates in the context of broader European genetic diversity. *Eur J Hum Genet*
499          **21**, 659-665 10.1038/ejhg.2012.229 (2013).
500    15.   Colonna, V. *et al.* Small effective population size and genetic homogeneity in
501          the Val Borbera isolate. *Eur J Hum Genet* **21**, 89-94 10.1038/ejhg.2012.113
502          (2013).

503    16.    Vitart, V. *et al.* SLC2A9 is a newly identified urate transporter influencing
504           serum urate concentration, urate excretion and gout. *Nat Genet* **40**, 437-442
505           10.1038/ng.106 (2008).
506    17.    The 1000 Genomes Project Consortium. A global reference for human genetic
507           variation. *Nature* **526**, 68-74 10.1038/nature15393 (2015).
508    18.    The UK10K Consortium. The UK10K project identifies rare variants in health
509           and disease. *Nature* **526**, 82-90 10.1038/nature14962 (2015).
510    19.    Gudbjartsson, D.F. *et al.* Large-scale whole-genome sequencing of the
511           Icelandic population. *Nat Genet* **47**, 435-444 10.1038/ng.3247 (2015).
512    20.    McQuillan, R. *et al.* Runs of homozygosity in European populations. *Am J Hum*
513           *Genet* **83**, 359-372 10.1016/j.ajhg.2008.08.007 (2008).
514    21.    Mathieson, I. & McVean, G. Demography and the age of rare variants. *PLoS*
515           *Genet* **10**, e1004528 10.1371/journal.pgen.1004528 (2014).
516    22.    Alexander, D.H., Novembre, J. & Lange, K. Fast model-based estimation of
517           ancestry in unrelated individuals. *Genome Res* **19**, 1655-1664
518           10.1101/gr.094052.109 (2009).
519    23.    Pickrell, J.K. & Pritchard, J.K. Inference of population splits and mixtures from
520           genome-wide allele frequency data. *PLoS Genet* **8**, e1002967
521           10.1371/journal.pgen.1002967 (2012).
522    24.    O'Connor, T.D. *et al.* Rare variation facilitates inferences of fine-scale
523           population structure in humans. *Mol Biol Evol* **32**, 653-660
524           10.1093/molbev/msu326 (2015).
525    25.    Reich, D., Thangaraj, K., Patterson, N., Price, A.L. & Singh, L. Reconstructing
526           Indian population history. *Nature* **461**, 489-494 10.1038/nature08365
527           (2009).
528    26.    Schiffels, S. & Durbin, R. Inferring human population size and separation
529           history from multiple genome sequences. *Nat Genet* **46**, 919-925
530           10.1038/ng.3015 (2014).
531    27.    Mezzavilla, M. & Ghirotto, S. *Neon*: An R package to estimate human effective
532           population size and divergence time from patterns of linkage disequilibrium
533           between SNPs. *J Comput Sci Syst Biol* **8**, 37-44 10.4172/jcsb.1000168 (2015).
534    28.    Browning, S.R. & Browning, B.L. Accurate non-parametric estimation of
535           recent effective population size from segments of identity by descent. *Am J*
536           *Hum Genet* **97**, 404-418 10.1016/j.ajhg.2015.07.012 (2015).
537    29.    Do, R. *et al.* No evidence that selection has been less effective at removing
538           deleterious mutations in Europeans than in Africans. *Nat Genet* **47**, 126-131
539           10.1038/ng.3186 (2015).
540    30.    Blomen, V.A. *et al.* Gene essentiality and synthetic lethality in haploid human
541           cells. *Science* **350**, 1092-1096 10.1126/science.aac7557 (2015).
542    31.    Colonna, V. *et al.* Human genomic regions with exceptionally high levels of
543           population differentiation identified from 911 whole-genome sequences.
544           *Genome Biol* **15**, R88 10.1186/gb-2014-15-6-r88 (2014).
545    32.    Duforet-Frebourg, N., Bazin, E. & Blum, M.B.G. Genome scans for detecting
546           local adaptation using a Bayesian factor model. *Mol Biol Evol* **31**, 2483-2495
547           10.1093/molbev/msu182 (2014).

548    33.    Field, Y. *et al.* Detection of human adaptation during the past 2000 years.
549            *Science* **354**, 760-764 10.1126/science.aag0776 (2016).
550    34.    Zoledziewska, M., Sidore, C. & Chiang, C.W. Height-reducing variants and
551            selection for short stature in Sardinia. *Nat Genet* **47**, 1352-1356
552            10.1038/ng.3403 (2015).
553    35.    Browning, S.R. & Browning, B.L. Rapid and accurate haplotype phasing and
554            missing-data inference for whole-genome association studies by use of
555            localized haplotype clustering. *Am J Hum Genet* **81**, 1084-1097
556            10.1086/521987 (2007).
557    36.    Price, A.L. *et al.* Principal components analysis corrects for stratification in
558            genome-wide association studies. *Nat Genet* **38**, 904-909 10.1038/ng1847
559            (2006).
560    37.    Li, J.Z. *et al.* Worldwide human relationships inferred from genome-wide
561            patterns of variation. *Science* **319**, 1100-1104 10.1126/science.1153717
562            (2008).
563    38.    Benazzo, A., Panziera, A. & Bertorelle, G. 4P: fast computing of population
564            genetics statistics from large DNA polymorphism panels. *Ecol Evol* **5**, 172-
565            175 10.1002/ece3.1261 (2014).
566    39.    Hill, W.G. Estimation of effective population size from data on linkage
567            disequilibrium. *Genetical Res.* **38**, 209-216 (1981).
568    40.    Hayes, B.J., Visscher, P.M., McPartlan, H.C. & Goddard, M.E. Novel multilocus
569            measure of linkage disequilibrium to estimate past effective population size.
570            *Genome Res* **13**, 635-643 10.1101/gr.387103 (2003).
571    41.    Tenesa, A. *et al.* Recent human effective population size estimated from
572            linkage disequilibrium. *Genome Res* **17**, 520-526 10.1101/gr.6023607 (2007).
573    42.    Browning, B.L. & Browning, S.R. Detecting identity by descent and estimating
574            genotype error rates in sequence data. *Am J Hum Genet* **93**, 840-851
575            10.1016/j.ajhg.2013.09.014 (2013).
576    43.    Kircher, M. *et al.* A general framework for estimating the relative
577            pathogenicity of human genetic variants. *Nat Genet* **46**, 310-315
578            10.1038/ng.2892 (2014).
579
580

## Author contributions

Y.X., C.T.-S., R.D. and E.Z.: design and supervision of the project. G.D., P.G., A.P., S.R., N.S., D.T. and J.F.W.: population liaison, sampling and DNA provision. N.S., JF.W. and R.D: comments and approval of the manuscript on behalf of the population consortia. Y.X., M.M. and M.H.: statistical method development. M.M., M.H., S.M., V.N., A.G., Q.A., V.C., L.S., C.F., G.R., H.C. and P.R. and J.A.: population-genetic analyses, statistical analyses and data interpretation. Y.C. and A.M.: bioinformatics support. S.M., N.S. and R.D: data processing and QC. Y.X., M.M., M.H., S.M., V.C., C.T-.S. and E.Z.: manuscript drafting. All authors: approval of the final version of the manuscript.

## Competing financial interests

The authors declare no competing financial interests.

## Acknowledgements

**Figure legends**


Fig. 1. General characteristics and demographic history of isolated and matched general populations. a. Geographical locations of samples. The base map was plotted in R using the mapdata package and circles were added using Photoshop. b. PCA using common variants. c. PCA using low-frequency variants. d. Sharing of rare variants within and between populations. Upper left triangle: $f_2$ variants; lower right triangle $f_{3-f10}$ variants. e. Effective population size (Ne) inferred from IBDNe for UKO and UKG during the past 9 KY. f. The lowest Ne inferred by IBDNe for all populations for the past 3KY, plotted as a function of the time at which it occurred.

Fig. 2. Isolation index (*Isx*) and its correlation with other genetic measures. a. Information summarized in *Isx*. b. Example of the correlation between *Isx* and other statistics, here *DVxy-coding*. c. Summary of the correlations between *Isx* and other population-genetic statistics. All the correlation coefficients are high and statistically significant.

Fig. 3. Purifying selection in the isolates and general populations. a. *Rxy*-missense statistic in each isolate, showing no evidence for increased genetic load in the isolates. The mean and standard deviation for each Rxy value from 100 bootstraps are shown. b. *DVxy-wg* (*DVxy*-whole genome) statistic in isolates and general populations, stratified by CADD score, showing enrichment of highly-functional low-frequency variants. c. *DVxy-coding* statistic in isolates and general populations, showing enrichment of low-frequency missense variants in isolates. d. SVxy-missense statistic in each isolate, showing relaxation of purifying selection in isolates in singletons. The standard errors for both *DVxy* and *SVxy* were calculated by randomly sampling data from 20 chromosomes 100 times. All of these analyses are based on the minimum-sample-size dataset (36 individuals from each population).

637     **Tables:**

638

639     **Table 1. Summary of variants discovered in this study**

640

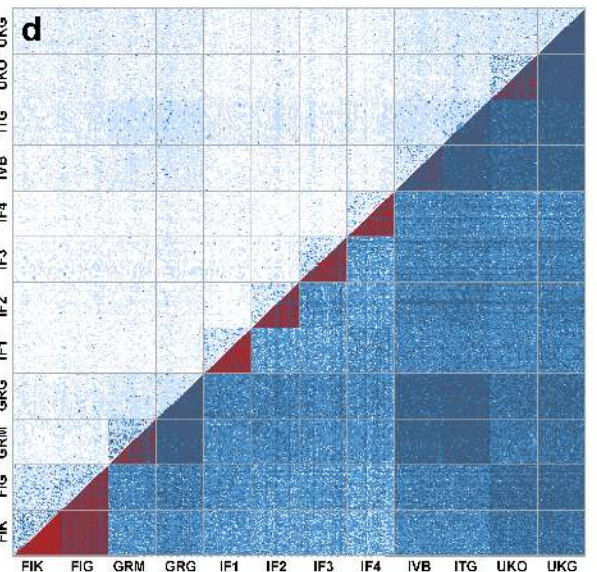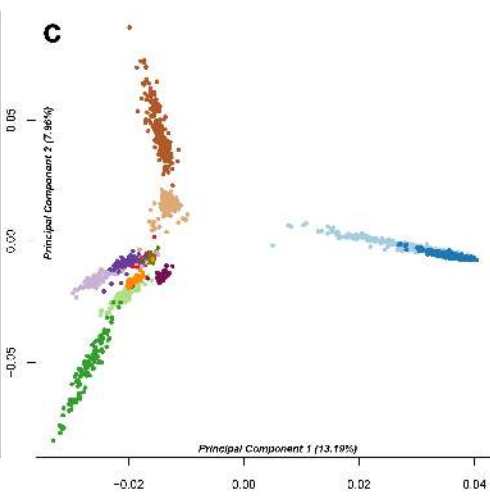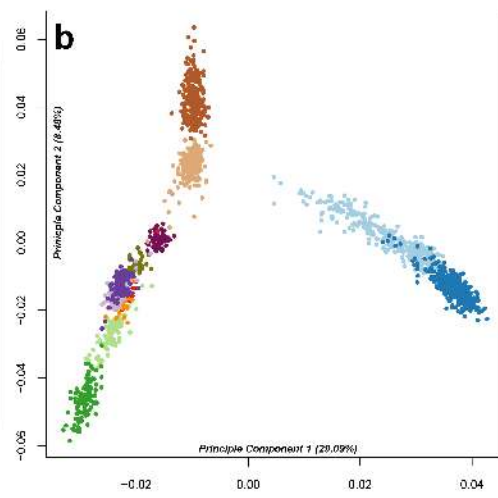| POP | n | average depth | MAF ≤2% | | MAF >2-≤5% | | MAF >5% | | Novel common SNPs in isolate* | Novel common SNPs in isolate** |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | total | novel % | total | novel % | total | novel % | | |
| **FIK** | 377 | 4x | 4,066,373 | 10.90 | 1,553,076 | 1.20 | 6,025,077 | 0.70 | 190,527 | 70,579 |
| **FIG** | 1564 | 6x | 6,548,833 | 11.80 | 1,540,915 | 0.80 | 6,053,704 | 0.70 | n.a. | n.a. |
| **GRM** | 249 | 4x | 5,129,513 | 7.20 | 1,447,981 | 1.10 | 6,111,923 | 0.80 | 513,272 | 49,884 |
| **GRG*** | 99 | 10-30x | 3,757,110 | n.a. | 1,321,955 | n.a. | 5,842,537 | n.a. | n.a. | n.a. |
| **IF1** | 60 | 4-10x | 1,456,881 | 1.30 | 1,420,929 | 1.30 | 5,890,714 | 0.80 | 320,191 | 119,157 |
| **IF2** | 45 | 4-10x | 1,063,098 | 1.30 | 1,554,145 | 1.00 | 6,001,568 | 0.80 | 273,694 | 94,496 |
| **IF3** | 47 | 4-10x | 961,059 | 1.30 | 1,455,284 | 1.10 | 6,068,304 | 0.80 | 299,603 | 107,281 |
| **IF4** | 36 | 4-10x | 1,030,673 | 1.30 | 1,124,789 | 1.10 | 6,001,625 | 0.80 | 308,356 | 122,254 |
| **IVB** | 222 | 6x | 4,857,767 | 1.60 | 1,396,799 | 0.80 | 6,112,476 | 0.80 | 188,972 | 30,284 |
| **UKO** | 397 | 4x | 5,963,416 | 11.70 | 1,471,782 | 0.80 | 6,047,383 | 0.80 | 193,300 | 36,512 |
| **Total** | 3096 | | 12,218,797 | 10.50 | 5,503,179 | 0.70 | 8,301,524 | 0.30 | | |

641

642     'Novel' variants are those not found in 1000 Genomes Project Phase 3 or UK10K

643     project. *Variants that are common (minor allele frequency, MAF ≥5.6%, alternative

644     allele count ≥ 4) in an isolated population but not common (MAF <1.4%, alternative

645     allele count ≤ 1) in its closest general population. **Variants that are common (MAF

646     ≥5.6%, alternative allele count ≥ 4) in an isolated population but not (MAF <1.4%,

647     alternative allele count ≤ 1) in *any* of the general populations. ***Different variant

648     calling procedure in this population.

649

650

651

652

a



b



c

d

e

f

**Populations**

FIK Finland Kuusamo
FIG Finland general
GRM Greece Manolis
GRG Greece general
IF1 Italy Friuli Venezia 1
IF2 Italy Friuli Venezia 2
IF3 Italy Friuli Venezia 3
IF4 Italy Friuli Venezia 4
IVB Italy Val Borbera
ITG Italy general
UKO UK Orkney
UKG UK general

**a**

Past

Ancestral population

$$Isx = \frac{\log[Tdg(100M^2)]}{\log(Ne)}$$

Ne $\quad$ Ne

Tdg

M

Present

General population $\qquad$ Isolate population

**b**

IF4 $\quad$ IF3

IF2 $\qquad$ IF1

Isx

1.7

1.6

1.5

1.4 $\qquad$ FIK

1.3

1.2 $\qquad$ GRM $\qquad\qquad$ UKO

1.1 $\qquad$ IVB

$r = 0.8$
$p = 0.016$

1.00 $\quad$ 1.05 $\quad$ 1.10 $\quad$ 1.15 $\quad$ 1.20 $\quad$ 1.25

*DVxy*-coding

**c**

| | DVxy | F | F_ST | HL | ROH | SVxy | Isx |
|---|---|---|---|---|---|---|---|
| DVxy | | | | | | | |
| F | 0.859 | | | | | | |
| F_ST | 0.772 | 0.901 | | | | | |
| HL | 0.848 | 0.977 | 0.866 | | | | |
| ROH | 0.787 | 0.955 | 0.948 | 0.941 | | | |
| SVxy | 0.720 | 0.905 | 0.901 | 0.929 | 0.920 | | |
| Isx | 0.801 | 0.969 | 0.978 | 0.918 | 0.992 | 0.912 | |

Pearson's *r*

1
0.98
0.96
0.94
0.91
0.89
0.87
0.86
0.83
0.81
0.79