

ARTICLE

Received 25 Jul 2016 | Accepted 16 May 2017 | Published 23 Jun 2017

DOI: 10.1038/ncomms15927

OPEN

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 Chen¹, Vagheesh Narasimhan¹, Arthur Gilly¹, Qasim Ayub¹, Vincenza Colonna^{1,3}, Lorraine Southam^{1,4}, Christopher Finan¹, Andrea Massaia^{1,5}, Himanshu Chheda⁶, Priit Palta^{6,7}, Graham Ritchie^{1,8,9}, Jennifer Asimit¹, George Dedoussis¹⁰, Paolo Gasparini¹¹, Aarno Palotie^{1,6,12,13,14,15,16}, Samuli Ripatti^{1,6,17}, Nicole Soranzo^{1,18}, Daniela Toniolo¹⁹, James F. Wilson^{9,20}, Richard Durbin¹, Chris Tyler-Smith¹ & Eleftheria Zeggini¹

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 3,059 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, DV_{xy} and SV_{xy} . We also develop an isolation-index (I_{sx}) that predicts the overall level of such key genetic characteristics and can thus help guide population choice in future complex-trait association studies.

¹The Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SA, UK. ²Institute for Maternal and Child Health, IRCCS Burlo Garofolo, University of Trieste, 34137 Trieste, Italy. ³Consiglio Nazionale delle Ricerche, Istituto di Genetica e Biofisica 'Adriano Buzzati-Traverso', via Pietro Castellino 111, 80131 Napoli, Italy. ⁴Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK. ⁵National Heart and Lung Institute, Imperial College London, London SW7 2AZ, UK. ⁶Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Tukholmankatu 8, 00290 Helsinki, Finland. ⁷Estonian Genome Center, University of Tartu, 23B Riia Street, 51010 Tartu, Estonia. ⁸European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SD, UK. ⁹MRC Human Genetics Unit, MRC IGMM, University of Edinburgh, Western General Hospital, Crewe Road, Edinburgh EH4 2XU, UK. ¹⁰Department of Nutrition and Dietetics, Harokopio University Athens, Athens, Eleftheriou Venizelou 70, Kallithea 176 76, Greece. ¹¹Medical 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. ¹²Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. ¹³Program in Medical and Population Genetics, The Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02114, USA. ¹⁴The Stanley Center for Psychiatric Research, The Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02114, USA. ¹⁵Psychiatric & Neurodevelopmental Genetics Unit, Department of Psychiatry, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. ¹⁶Department of Neurology, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. ¹⁷Department of Public Health, University of Helsinki, Helsinki FI-00014, Finland. ¹⁸Department of Haematology, University of Cambridge, Cambridge CB2 0XY, UK. ¹⁹Division of Genetics and Cell Biology, San Raffaele Scientific Institute, via Olgettina 60, 20132 Milan, Italy. ²⁰Usher Institute of Population Health Sciences and Informatics, University of Edinburgh, Teviot Place, Edinburgh, EH8 9AG Scotland, UK. * These authors contributed equally to this work. Correspondence and requests for materials should be addressed to Y.X. (email: ylx@sanger.ac.uk) or to E.Z. (email: Eleftheria@sanger.ac.uk).

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 studies². We thus generate and analyse low-depth ($4 \times -10 \times$) 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 data set includes newly generated low-depth ($4 \times -10 \times$) WGS from eight cohorts drawn from isolated European populations: one each from Kuusamo in Finland (FIK) and Crete in Greece (GRM¹³), four from Friuli-Venezia Giulia in Italy (IF1, IF2, IF3 and IF4 (ref. 14)), and one each from Val Borbera in Italy (IVB¹⁵) and the Orkney Islands in the UK (UKO¹⁶); and the closest non-isolated (general) population: Finland (FIG⁹), Greece (GRG), together with publicly available data for Italy (ITG¹⁷) and UK (UKG¹⁸) (Fig. 1a and Supplementary Table 1). We generated a superset of variants called in these cohorts and all 26 population samples in the 1000 Genomes Project Phase 3 (ref. 17), and performed multi-sample genotype calling across all 9,375 samples (3,059 from the current study, 2,353 from the 1000 Genomes Project Phase 3 release and 3,781 from UK10K). Both individual population and amalgamated genotype call data, which have greater than 99% concordance with genotyping data (Supplementary Table 2), are available to the scientific community (Data availability).

General description of the variants in the isolates. We identified approximately 12.2 million variants with minor allele frequency (MAF) $\leq 2\%$ (rare), 5.5 million with $2 < \text{MAF} \leq 5\%$ (low-frequency) and 8.3 million variants with $\text{MAF} > 5\%$ (common) across the ten populations newly sequenced here (eight isolates, GRG and FIG). Of these, 10.5, 0.7 and 0.3%, respectively, are novel (Table 1 and Supplementary Table 3). As expected, most of the isolates have lower numbers of variant sites per genome than their closest general population (Supplementary Fig. 1, Supplementary Table 5). We find $\sim 188,000 - \sim 513,000$ variants that are common with $\text{MAF} > 5.6\%$ in each isolate but with $\text{MAF} \leq 1.4\%$ in its closest general population (Table 1); $\sim 30,000 - 122,000$ of these per isolate have frequency $\leq 1.4\%$ in all the general samples studied, among which $\sim 150 - \sim 700$ in coding regions and $\sim 500 - \sim 2,800$ genome-wide are deleterious (Supplementary Table 4). These common and low-frequency variants are thus useful markers for whole-genome association studies in these populations and some of them (if absent from the general population) could potentially lead to novel association signals. They include known examples such as rs76353203 (R19X) in *APOC3* in GRM, which is associated with high-density lipoprotein and triglyceride levels⁶.

Population-genetic analyses in the isolates. Previous population-genetic studies of isolates have, with some exceptions^{11,19}, been based on common variants found on genotyping arrays, and have illustrated general characteristics such as low genetic diversity and longer shared haplotypes^{9,13–15,19,20}. Rare variants discovered from sequencing are on average more recent in origin than common variants²¹ and therefore more powerful for distinguishing closely related populations and more informative about recent demographic history. We find that isolates are, as expected, genetically close to their matched general population in principal component analyses (PCA), ADMIXTURE²² and TreeMix²³ using common variants (Fig. 1b, Supplementary Figs 2–5 and Supplementary Table 6), but PCA using rare and low-frequency variants, as found previously²⁴, distinguishes them more clearly from the general population and also from other isolates, particularly among the Italian samples (Fig. 1c, Supplementary Fig. 2). The majority of sharing of variants present just twice across all samples of 36 individuals from each population (f_2 variants²¹) takes place within the same population, and the isolates generally share more with their closest general population than with other populations. This latter trend, however, is not apparent for IF1–IF4, who show little sharing with any other population, pointing to a greater level of isolation and lower level of gene flow with their general population (Fig. 1d, upper triangle and Supplementary Fig. 7), which is confirmed by f_3 -statistics²⁵ comparing with a worldwide population panel of HGDP-CEPH samples using common SNPs (Supplementary Fig. 6). $f_3 - f_{10}$ variant sharing demonstrates sharing by ITG and IVB with both Greek and UK populations (Fig. 1d, lower triangle and Supplementary Fig. 7), potentially indicative of their more ancient heritage.

Population demographic history. All populations studied here, both isolates and general, appear to have shared a comparable effective population size (N_e) history before 20 thousand years ago (KYA) based on the multiple sequentially Markovian coalescent method²⁶ (Supplementary Fig. 9). The isolates diverged from their general populations within the last $\sim 5,000$ years based on LD estimations²⁷ (Supplementary Table 7 and Supplementary Fig. 8) and yet had sharp decreases in their population sizes in more recent times as estimated using inferred long segments of identity by descent (IBD)²⁸ (Fig. 1e,f and Supplementary Fig. 10). Different isolates also split from their respective general populations at different times. For example, IF1–IF4 split from ITG $\sim 4 - 5$ KYA, while most other isolates split from their general populations within the last $\sim 1,000$ years (Supplementary Table 7).

The different demographic histories of different isolates should lead to different genetic characteristics. To summarize these features in a single quantitative measure that can be calculated from genotype data, as well as sequence data, we developed an isolation index (Isx) which combines information on the divergence time from the general population (T_{dg}), N_e and migration rate (M), such that early-divergence-time isolates with small N_e and low M have a high Isx value (Fig. 2a and Supplementary Fig. 11). The different isolates show different Isx values: IF1, IF2, IF3 and IF4 have the highest, while IVB has the lowest (Supplementary Table 8). Isx values are highly correlated with other population-genetic characteristics (for example, Fig. 2b,c, Supplementary Table 11), such as genome-wide pairwise F_{ST} between isolates and their matching general population (reflecting the genetic drift of the isolates) (Supplementary Fig. 12), the total length and number of runs of homozygosity (ROH) (Supplementary Fig. 13), inbreeding coefficient (F) (Supplementary Fig. 14) and length of LD

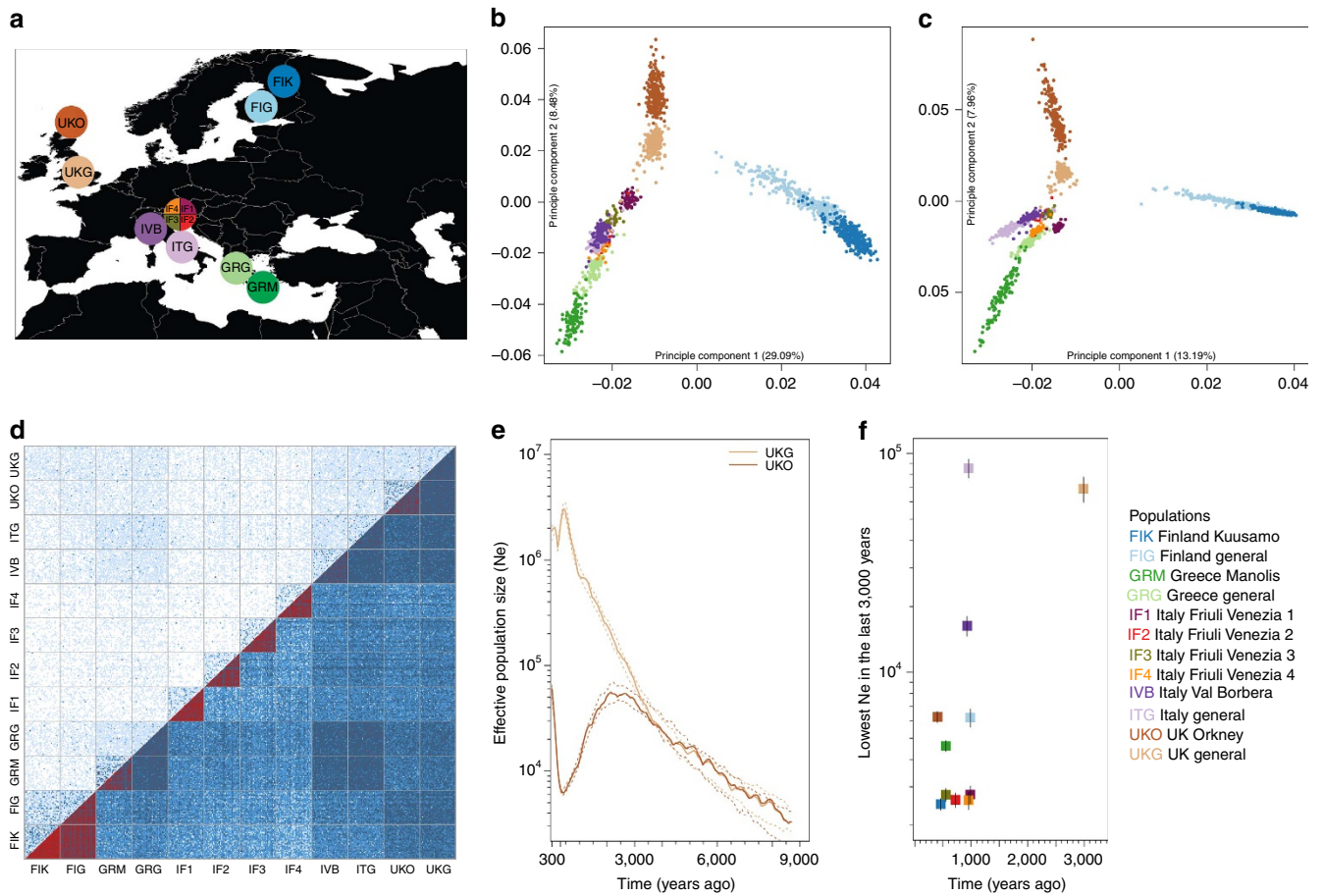


Figure 1 | General characteristics and demographic history of isolated and matched general populations. (a) Geographical locations of samples. The base map was plotted in R using the mapdata package and circles were added using Photoshop. (b) PCA using common variants. (c) PCA using low-frequency variants. (d) Sharing of rare variants within and between populations. Upper left triangle: f_2 variants; lower right triangle f_3 – f_{10} variants. (e) Effective population size (N_e) inferred from IBDNe for UKO and UKG during the past nine KY. (f) The lowest N_e inferred by IBDNe for all populations for the past three KY, plotted as a function of the time at which it occurred.

Table 1 | Summary of variants discovered in this study.

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	1,564	6x	6,548,833	11.80	1,540,915	0.80	6,053,704	0.70	na	na
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	na	1,321,955	na	5,842,537	na	na	na
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	3,096		12,218,797	10.50	5,503,179	0.70	8,301,524	0.30		

^{*}Novel[†] variants are those not found in 1000 Genomes Project Phase 3 or UK10K project.

^{*}Variants that are common (minor allele frequency, MAF ≥ 5.6%, alternative allele count ≥ 4) in an isolated population but not common (MAF < 1.4%, alternative allele count ≤ 1) in its closest general population.

[†]Variants that are common (MAF ≥ 5.6%, alternative allele count ≥ 4) in an isolated population but not (MAF < 1.4%, alternative allele count ≤ 1) in any of the general populations.

[‡]Different variant calling procedure in this population.

(Supplementary Figs 15 and 16 and Supplementary Tables 9 and 10). All these characteristics are correlated, but the pairwise correlation coefficients show that I_{sx} is a slightly better overall predictor of the other measures than any single existing measure

(Fig. 2c, Supplementary Fig. 17 and Supplementary Table 11); moreover, it is potentially more robust to confounding factors as it is calculated from three demographic parameters, while the others are all based on single measurements.

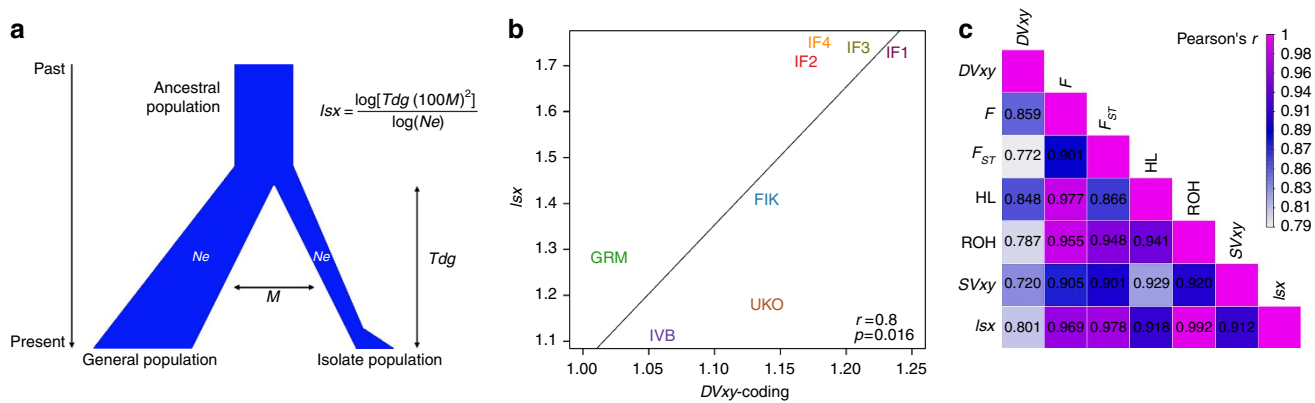


Figure 2 | Isolation index (Isx) and its correlation with other genetic measures. (a) Information summarized in Isx . (b) Example of the correlation between Isx and other statistics, here DV_{xy} -coding. (c) Summary of the correlations between Isx and other population-genetic statistics. All the correlation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

Purifying selection analyses. Several lines of evidence suggest relaxed purifying selection in the isolates due to their reduced Ne , although as expected we do not detect substantially increased genetic load per genome using the Rxy statistic²⁹ based on all of the variants in the genomes (Fig. 3a and Supplementary Table 12). First, we see different levels of enrichment of low-frequency functional variants in isolates (Fig. 3b,c, Supplementary Tables 13 and 14, Supplementary Fig. 18a) quantified by a new statistic, DV_{xy} -coding, developed here (DV: drifted variants). DV_{xy} -coding measures the ratio of functional coding variants (missense plus loss-of-function (LoF)) in isolates compared to the closest general population (and vice-versa), adjusted for the corresponding ratios of intergenic variants in order to correct for the effect of genetic drift. We applied this only to a subclass of DVs, defined as low-frequency (2–5%, the best choice according to the sample size we have) in any isolate, yet at least three-fold higher than in the closest general population (and vice versa). We find that DV_{xy} -coding is >1 in all isolates and <1 in all general populations (Fig. 3c, Supplementary Fig. 18a and Supplementary Table 13). We also calculated a similar DV_{xy} -wg statistic by stratifying whole-genome variants according to their combined annotation dependent depletion (CADD) score (0–5, neutral variants; 5–10, mildly deleterious; 10–20, deleterious; and >20 , highly deleterious; these cut-off choices balance the number of variants in each bin to allow us comparable statistical power among all bins, although the conclusions are robust to the particular cut-off values chosen and different bins (Supplementary Figs 18b and 19)). The DV_{xy} -wg values are differentiated for variants with CADD score of 10–20 and significantly so (assessed using the jack-knife bootstrap method) for ones with CADD scores >20 , with DV_{xy} -wg values >1 in all isolates and <1 in all general populations (Fig. 3b, Supplementary Fig. 18b and Supplementary Table 14). This demonstrates enrichment of low-frequency functional variants, both coding and genome-wide with CADD score >10 , in the isolated populations. Moreover, both DV_{xy} -coding and DV_{xy} -wg values are correlated with Isx , suggesting that different isolation characteristics lead to different levels of enrichment of functional variants.

We also investigated the relaxation of purifying selection by assessing functional (missense) singleton variants (SV) pooled for all of the genes that have at least one singleton missense or synonymous variant in a pair of populations (one isolate and its general population), correcting with pooled synonymous variants (SV_{xy} statistic). We find a substantial deviation from 1 for functional singletons in all of the isolates (Fig. 3d and

Supplementary Table 15), with SV_{xy} values positively correlating with Isx (Fig. 2c and Supplementary Fig. 20). We also find that the proportion of relaxed essential genes³⁰ with $SV_{xy} >1$ in isolates is significantly higher than in the general population (Supplementary Table 15). Such rare and low-frequency drifted functional variants, measured by both SV_{xy} and DV_{xy} , are particularly relevant for boosting the power of association studies⁶.

Positive selection analyses. We do not find convincing evidence for positive selection in any isolate using ΔDAF ³¹, $PCAdapt$ ³² or singleton density score (SDS)³³, although we do identify some highly differentiated variants (Supplementary Fig. 21 and Supplementary Tables 16 and 17), including in the protein-coding genes *ALK*, *SPNS2*, *SLC39A11* and *ACSS2*, which can nevertheless be accounted for by drift. Interestingly, we also find six highly differentiated variants shared between different isolates from Italy, IF2, IF3 and IF4, but interpret them as likely to result from drift or positive selection for the ancestral allele in the ITG (Supplementary Table 17). We find that the SDS method has little power in our samples because of their small size, and failed to detect selection even at the lactose tolerance SNP in the UKO, a known strong signal of recent selection (Supplementary Fig. 22).

Discussion

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

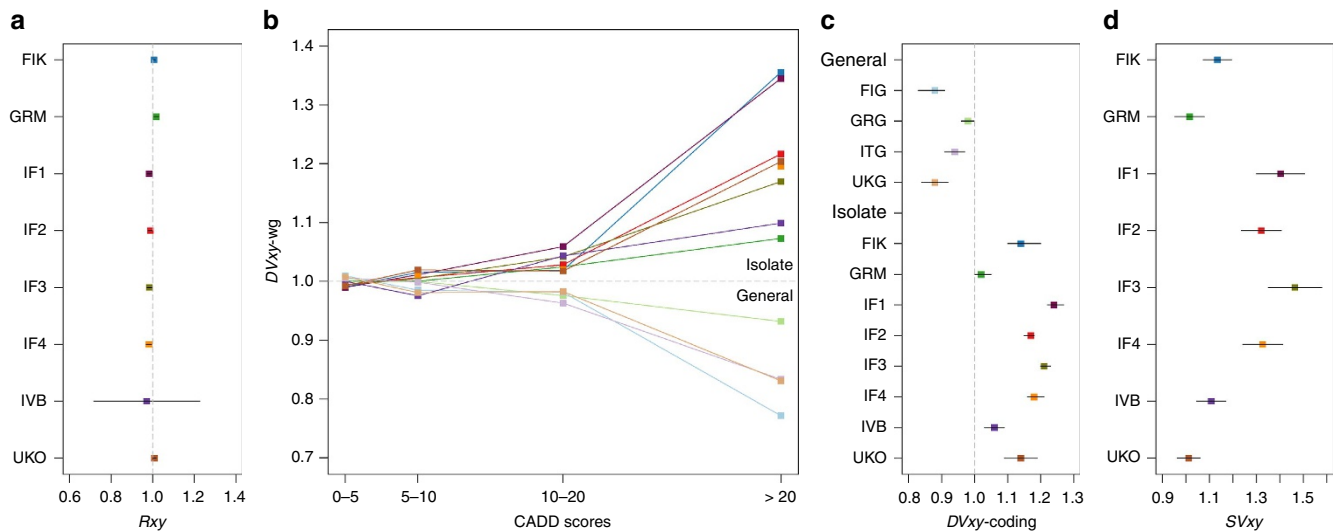


Figure 3 | Purifying selection in the isolates and general populations. (a) R_{xy} -missense statistic in each isolate, showing no evidence for increased genetic load in the isolates. The mean and s.d. for each R_{xy} value from 100 bootstraps are shown. (b) DV_{xy-wg} (DV_{xy} -whole genome) statistic in isolates and general populations, stratified by CADD score, showing enrichment of highly functional low-frequency variants. (c) $DV_{xy-coding}$ statistic in isolates and general populations, showing enrichment of low-frequency missense variants in isolates. (d) SV_{xy} -missense statistic in each isolate, showing relaxation of purifying selection in isolates in singletons. The s.e.'s for both DV_{xy} and SV_{xy} were calculated by randomly sampling data from 20 chromosomes 100 times. All of these analyses are based on the minimum-sample-size data set (36 individuals from each population).

Methods

Data set and variant calling. The data set includes 3,059 whole-genome low-depth sequences generated at The Wellcome Trust Sanger Institute using the Illumina Genome Analyzer II and Illumina HiSeq 2000 platforms, as well as 100 high-depth sequences from the Illumina HiSeq X Ten (Fig. 1a and Supplementary Table 1). Informed consent was obtained from all subjects and the study was approved by the HMDMC (Human Materials and Data Management Committee) of the Wellcome Trust Sanger Institute. The multi-sample genotype calling across all of the low-coverage sequencing data from the current study, as well as 2,353 from the 1000 Genomes Project Phase 3 release, and 3,781 from UK10K (a total of 9,375) was performed with the defined site selection criteria (Supplementary Note). Genotype likelihoods were calculated with samtools/bcftools (0.2.0-rc9) and then genotypes were called and phased using Beagle v4 (r1274) (ref. 35). We assessed the performance of the genotype calling from the low coverage data using the available genotype chip data for a subset of the cohorts consisting of 4,665 individuals, and calculated the discordance rates on chromosome 20 separately for the categories REF-REF, REF-ALT and ALT-ALT.

The sample sizes are very different across these collections, and we used three different standard-sized subsets of the samples for different analyses: (1) the whole data set; (2) the sample-size-matched data set, obtained either by randomly selecting samples from general population to match the isolated population (for example, we randomly select 377 from FIG to match FIK), or by randomly selecting a subset of the isolated population to match the general population (for example, we randomly select 108 IVB to match the general population ITG); (3) the minimum-sample-size data set of 36 individuals per population. By doing this, we maximize the use of the data for different analyses, and we specify which data set is used for each analysis. The sequencing depth is also different across different populations, within a 2.5-fold range (apart from GRG, in which variants were called differently, details in Supplementary Notes), and we allowed for these differences when interpreting the results.

Variant counts. We first re-annotated all variants using the Variant Effect Predictor annotation from Ensembl 76 with the '-pick' option, which gives one annotation per variant. We then performed variant counting at both the population and individual level, stratifying by functional categories and frequency bins. These counts were either plotted in figures or summarized as median values in tables. We carried out these analyses using both the sample-size-matched data set and the minimum-sample-size data set.

Population-genetic analyses. We used the whole data set for the analyses in this section, unless otherwise specified. PCAs were performed separately with common variants or rare variants using EIGENSTRAT v.501 (ref. 36). Shared ancestry between the populations studied here was evaluated using ADMIXTURE v1.22 (ref. 22). The relationships between the populations studied here, combined with worldwide populations from the HGDP-CEPH panel³⁷, were also examined using

ancestry graph analyses implemented in TreeMix v.1.12 (ref. 23). We also used formal test of f_3 -statistics²⁵ to investigate population mixture in the history of the populations studied here, as well as worldwide populations from the HGDP-CEPH panel. Rare f_2 variants (with only two copies of the alternative allele in the minimum-sample-size data set) and moderately rare f_{3-10} variants (3–10 copies of the alternative allele in the same data set) are particularly informative for investigating recent human history²¹. We investigated the sharing pattern of these two types of variant by summing all f_2 variants or any random two alleles of the f_{3-10} variants shared by pairs of individuals. We plotted the results as a heat map using the image¹ function from the base R package (<https://stat.ethz.ch/R-manual/R-devel/library/graphics/html/image.html>). Variants were aggregated by pair of individuals using the 'count' function of the plyr package, then arranged in matrix form and colorized using 'colorRampPalette' from the colorspace package (<https://cran.r-project.org/web/packages/colorspace/index.html>). ROH, inbreeding coefficient (F) as well as the length of LD-blocks were calculated in PLINK, and finally genome-wide F_{ST} values between isolates and their general populations were calculated with the software 4P (ref. 38) using the minimum-sample-size data set.

Demographic inferences. LD-based^{39–41} demographic inference was performed in the NeON R package²⁷ using the minimum-sample-size data set; the median and confidence interval were estimated using the 50th, 5th and 95th percentiles of the distribution of long-term N_e in each time interval. We used the multiple sequentially Markovian coalescent method²⁶ to infer demographic changes before 20,000 years ago using four individual sequences from each population. In order to account for some loss of heterozygous sites in the low-depth data, we used a slow mutation rate of 0.8×10^{-8} mutations per nucleotide per generation and a longer generation time of 33 years. We then estimated more recent demographic changes (from the present to ~9,000 years ago) using IBDNe²⁸ with the minimum-sample-size data set. We used IBDseq⁴² to detect IBD segments in sequence data from chromosome 2 in all populations. We then used IBDNe with the default parameters and a minimum IBD segment length of 2 centiMorgan (cM) units. We assumed a generation time of 29 years.

Isolation index. In order to quantify the different isolation levels of different isolates, we developed an index that combines three demographic parameters: (a) T_{dg} , (b) N_e and (c) the level of private isolate ancestry (M). We call this estimate the Isolation index (I_{sx}). It is defined as:

$$I_{sx} = \frac{\log(T_{dg}(100 \times M)^2)}{\log(N_e)}$$

Both T_{dg} and N_e were inferred from the LD-based method using the NeON R package²⁷. M is difficult to estimate directly from SNP genotype data, so here we estimated the difference of shared ancestral components between an isolate and its general population from ADMIXTURE analysis. We ran ADMIXTURE with only one isolate and it closest general population using $K=2$. We then estimated the

difference in the means of ancestry between the isolate and its general population. The M parameter was defined as Delta Ancestry.

Rxy analysis. Rxy statistics²⁹ between each pair of populations (an isolate and its closest general population) for different functional categories were calculated using the matched-sample-size data for missense and LoF variants, including stop gain, splice donor and acceptor variants, using synonymous variants as controls (we did not use intragenic variants as control because of the ascertainment in the ITG which has high-depth exome sequences and low depth for the rest of the genome). We also calculated Rxy statistics for variants with CADD scores⁴³ greater than 10 and 20, using variants with CADD scores less than 5 as controls. The mean and s.d. for each Rxy value were obtained from 100 bootstraps.

DVxy analysis. A new statistic, DV_{xy} , was developed to quantify the enrichment of low-frequency functional variants in the isolates using both the matched-sample-size and minimum-sample-size data sets. It calculates the proportion of functional variants in each isolate compared with its general population, correcting for genetic drift at the same time. We calculated DV_{xy} specifically for the subset of variants with DAF 2–5% in the isolate, and at least three times lower in its closest general population, or vice-versa. We called these variants ‘drifted variants’ (DV). DV_{xy} was calculated for both coding regions and whole genomes.

For coding variants, we defined missense or missense plus LoF variants as functional variants. We counted the number of functional DVs and neutral (intergenic) DVs in each isolate (population x) and the corresponding general population (population y). The ratio between the fraction of DV variants from the isolated population (corrected by the count of intergenic variants) and the corresponding fraction of DV variants from its general population was defined as the DV_{xy} statistic. If DV_{xy} is equal to 1, there is no enrichment for the functional DVs in the isolate; less than 1 indicates depletion, and greater than 1 indicates enrichment.

$$DV_{xy_coding} = \frac{\% DVx\ missense}{\% DVx\ intergenic} / \frac{\% DVy\ missense}{\% DVy\ intergenic}$$

For the whole genome, we used different CADD score cut-offs and bins. We calculated a DV statistic by stratifying the variants according to their CADD scores (0–5, neutral variants; 5–10, mildly deleterious; 10–20, deleterious; and greater than 20, highly deleterious) for each isolate and its closest general population. We finally calculated a ratio of the fraction of DV variants (from each class) between the isolate and its general population, and vice-versa. The following formula shows the DV_{xy-wg} calculation for variants with CADD score between i and j in an isolate and its general population.

$$DV_{xyCADD(i,j)} = \frac{\% DVx(CADDi - j)}{\% DVy(CADDi - j)}$$

The 95% confidence interval for each calculation was obtained by randomly sampling data from 20 chromosomes 100 times.

SVxy analysis. We further investigated the relaxation of purifying selection in the isolated populations using SVs. Here, we also used the minimum-sample-size data set. Another new statistic, SV_{xy} , was developed to measure the ratio of missense versus synonymous singletons per gene in each population, as well as the ratio of the sum of singletons in all genes which have at least one singleton in the pair of the populations (one isolate and one general population). We counted the number of missense singletons and synonymous singletons per gene in each population, and SV_{gene} was calculated as:

$$SV_{gene} = \frac{(SV\ missense\ count + 1)}{(SV\ synonymous\ count + 1)}$$

$SV_{gene} > 1$ indicates relaxation of purifying selection; $SV_{gene} = 1$ indicates neutrality; and $SV_{gene} < 1$ indicates purifying selection.

We then divided the gene list into essential genes³⁰ and non-essential genes (the rest), and calculated a statistic, G_{SV} , for each population, defined as:

G_{SV} = percentage of essential genes with $SV_{gene} > 1$ / percentage of non-essential genes with $SV_{gene} > 1$

We finally calculated a statistic, SV_{xy} , which is the ratio of SV_{pop} of each isolate to SV_{pop} of its general population. SV_{pop} for each isolate and its general population was calculated using all genes which have at least one singleton in the pair of the populations and defined as $SV_{pop} = \Sigma (SV\ missense\ counts) / \Sigma (SV\ synonymous\ counts)$.

We used the same annotation as in the variant counts. We calculated a confidence interval for each estimate using bootstrapping of 80% of the genes 100 times.

Correlation analyses. We calculated pair-wise correlation coefficients between the I_{sx} values, population-genetic measurements ROH, F_{ST} , and number and length of LD blocks, as well as the newly developed statistics DV_{xy} and SV_{xy} using the Pearson correlation in R.

Positive selection analyses. We calculated genome-wide pairwise derived allele frequency differences (deltaDAF) for each pair of populations (an isolate and its general population) as described previously³¹ using the matched-sample-size data set. We also carried out PCAdapt analyses³² for each pair of populations using the whole data set. Both analyses look for high derived allele frequency variants in the isolates, and will not be affected by sample size. Finally, we ran the SDS method³³ using the whole UKO and UKG data sets, which have the largest sample sizes for both isolate and its general population, and thus the greatest power for this method.

Data availability. Amalgamated genotype calls across all populations studied are available through the European Genome/Phenome Archive (EGAD00001002014) with Data Access Agreement described in Supplementary Information.

References

- Zeggini, E. Using genetically isolated populations to understand the genomic basis of disease. *Genome Med.* **6**, 83 (2014).
- Hatzikotoulas, K., Gilly, A. & Zeggini, E. Using population isolates in genetic association studies. *Brief Funct. Genomics* **13**, 371–377 (2014).
- Zuk, O. *et al.* Searching for missing heritability: designing rare variant association studies. *Proc. Natl Acad. Sci. USA* **111**, E455–E464 (2014).
- Pollin, T. I. *et al.* A null mutation in human APOC3 confers a favorable plasma lipid profile and apparent cardioprotection. *Science* **322**, 1702–1705 (2008).
- Gudmundsson, J. *et al.* A study based on whole-genome sequencing yields a rare variant at 8q24 associated with prostate cancer. *Nat. Genet.* **44**, 1326–1329 (2012).
- Tachmazidou, I. *et al.* A rare functional cardioprotective APOC3 variant has risen in frequency in distinct population isolates. *Nat. Commun.* **4**, 2872 (2013).
- 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).
- 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**, 640–642 (2015).
- Lim, E. T. *et al.* Distribution and medical impact of loss-of-function variants in the Finnish founder population. *PLoS Genet.* **10**, e1004494 (2014).
- Moltke, I. *et al.* A common Greenlandic TBC1D4 variant confers muscle insulin resistance and type 2 diabetes. *Nature* **512**, 190–193 (2014).
- Sidore, C. *et al.* Genome sequencing elucidates Sardinian genetic architecture and augments association analyses for lipid and blood inflammatory markers. *Nat. Genet.* **47**, 1272–1281 (2015).
- 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**, 294–298 (2014).
- 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).
- Esco, T. *et al.* Genetic characterization of northeastern Italian population isolates in the context of broader European genetic diversity. *Eur. J. Hum. Genet.* **21**, 659–665 (2013).
- Colonna, V. *et al.* Small effective population size and genetic homogeneity in the Val Borbera isolate. *Eur. J. Hum. Genet.* **21**, 89–94 (2013).
- Vitart, V. *et al.* SLC2A9 is a newly identified urate transporter influencing serum urate concentration, urate excretion and gout. *Nat. Genet.* **40**, 437–442 (2008).
- The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
- The UK10K Consortium. The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82–90 (2015).
- Gudbjartsson, D. F. *et al.* Large-scale whole-genome sequencing of the Icelandic population. *Nat. Genet.* **47**, 435–444 (2015).
- McQuillan, R. *et al.* Runs of homozygosity in European populations. *Am. J. Hum. Genet.* **83**, 359–372 (2008).
- Mathieson, I. & McVean, G. Demography and the age of rare variants. *PLoS Genet.* **10**, e1004528 (2014).
- Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* **19**, 1655–1664 (2009).
- Pickrell, J. K. & Pritchard, J. K. Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS Genet.* **8**, e1002967 (2012).
- O’Connor, T. D. *et al.* Rare variation facilitates inferences of fine-scale population structure in humans. *Mol. Biol. Evol.* **32**, 653–660 (2015).
- Reich, D., Thangaraj, K., Patterson, N., Price, A. L. & Singh, L. Reconstructing Indian population history. *Nature* **461**, 489–494 (2009).
- Schiffels, S. & Durbin, R. Inferring human population size and separation history from multiple genome sequences. *Nat. Genet.* **46**, 919–925 (2014).
- Mezzavilla, M. & Ghirrotto, S. *Neon*: an R package to estimate human effective population size and divergence time from patterns of linkage disequilibrium between SNPs. *J. Comput. Sci. Syst. Biol.* **8**, 37–44 (2015).
- Browning, S. R. & Browning, B. L. Accurate non-parametric estimation of recent effective population size from segments of identity by descent. *Am. J. Hum. Genet.* **97**, 404–418 (2015).

29. Do, R. *et al.* No evidence that selection has been less effective at removing deleterious mutations in Europeans than in Africans. *Nat. Genet.* **47**, 126–131 (2015).
30. Blomen, V. A. *et al.* Gene essentiality and synthetic lethality in haploid human cells. *Science* **350**, 1092–1096 (2015).
31. Colonna, V. *et al.* Human genomic regions with exceptionally high levels of population differentiation identified from 911 whole-genome sequences. *Genome Biol.* **15**, R88 (2014).
32. Duforet-Frebourg, N., Bazin, E. & Blum, M. B. G. Genome scans for detecting local adaptation using a Bayesian factor model. *Mol. Biol. Evol.* **31**, 2483–2495 (2014).
33. Field, Y. *et al.* Detection of human adaptation during the past 2000 years. *Science* **354**, 760–764 (2016).
34. Zoledziewska, M., Sidore, C. & Chiang, C. W. Height-reducing variants and selection for short stature in Sardinia. *Nat. Genet.* **47**, 1352–1356 (2015).
35. Browning, S. R. & Browning, B. L. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am. J. Hum. Genet.* **81**, 1084–1097 (2007).
36. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909 (2006).
37. Li, J. Z. *et al.* Worldwide human relationships inferred from genome-wide patterns of variation. *Science* **319**, 1100–1104 (2008).
38. Benazzo, A., Panziera, A. & Bertorelle, G. 4P: fast computing of population genetics statistics from large DNA polymorphism panels. *Ecol. Evol.* **5**, 172–175 (2014).
39. Hill, W. G. Estimation of effective population size from data on linkage disequilibrium. *Genetical Res.* **38**, 209–216 (1981).
40. Hayes, B. J., Visscher, P. M., McPartlan, H. C. & Goddard, M. E. Novel multilocus measure of linkage disequilibrium to estimate past effective population size. *Genome Res.* **13**, 635–643 (2003).
41. Tenesa, A. *et al.* Recent human effective population size estimated from linkage disequilibrium. *Genome Res.* **17**, 520–526 (2007).
42. Browning, B. L. & Browning, S. R. Detecting identity by descent and estimating genotype error rates in sequence data. *Am. J. Hum. Genet.* **93**, 840–851 (2013).
43. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).

Acknowledgements

We thank all study participants for making this work possible. Our work was supported by the Wellcome Trust (098051). We also thank the UK10K Consortium and SISu Consortium for making their data available to this study; detailed acknowledgements for these contributions are included in Supplementary Information.

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., J.F.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.P. 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.

Additional information

Supplementary Information accompanies this paper at <http://www.nature.com/naturecommunications>

Competing interests: The authors declare no competing financial interests.

Reprints and permission information is available online at <http://npg.nature.com/reprintsandpermissions/>

How to cite this article: Xue, Y. *et al.* Enrichment of low-frequency functional variants revealed by whole-genome sequencing of multiple isolated European populations. *Nat. Commun.* **8**, 15927 doi: 10.1038/ncomms15927 (2017).

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

© The Author(s) 2017