



## Supplementary Materials for

### **Evolution and Functional Impact of Rare Coding Variation from Deep Sequencing of Human Exomes**

Jacob A. Tennessen, Abigail W. Bigham, Timothy D. O'Connor, Wenqing Fu, Eimear E. Kenny, Simon Gravel, Sean McGee, Ron Do, Xiaoming Liu, Goo Jun, Hyun Min Kang, Daniel Jordan, Suzanne M. Leal, Stacey Gabriel, Mark J. Rieder, Goncalo Abecasis, David Altshuler, Deborah A. Nickerson, Eric Boerwinkle, Shamil Sunyaev, Carlos D. Bustamante, Michael J. Bamshad,\* Joshua M. Akey,\* Broad GO, Seattle GO, on behalf of the NHLBI Exome Sequencing Project

\*To whom correspondence should be addressed. E-mail: akeyj@uw.edu (J.M.A.);  
mbamshad@u.washington.edu (M.J.B.)

Published 17 May 2012 on *Science Express*  
DOI: 10.1126/science.1219240

#### **This PDF file includes:**

Materials and Methods  
Supplementary Text  
Figs. S1 to S19  
Tables S1 to S7  
References

# Supporting Online Material

## Table of Contents

Study Sample .....	2
Cohort descriptions .....	3
Informed Consent .....	6
Exome sequencing .....	6
Variant calling .....	9
Data filtering.....	10
Identifying related individuals.....	11
Assignment of ancestry labels .....	11
Data annotation .....	11
Prediction of functional variation .....	12
Analysis of batch and cohort effects .....	12
Validation .....	12
Power of variant discovery.....	14
Comparison to 1000 Genomes Project.....	14
Demographic inference.....	15
Analysis of natural selection .....	16
Rare variant association power analysis .....	18
Further acknowledgements .....	20
Supplementary tables .....	28
Supplementary figures.....	36

## Study sample

The 2,440 exomes used in the analysis were generated from samples ascertained from 15 different cohorts.

Phenotypes associated with samples and cohort contributions:

Phenotype	Cohort	Number of samples contributed
Early onset myocardial infarction (n=1094)	Atherosclerosis Risk in Communities	258
	Cleveland Clinic GeneBank	3
	Framingham Heart Study	177
	Heart Attack Risk in Puget Sound	13
	Jackson Heart Study	112
	Massachusetts General Hospital-Premature Coronary Artery Disease	147
	Penn-Cath	34
	Translational Research Investigating Underlying Disparities in Acute Myocardial Infarction Patients' Health Status	77
	Women's Health Initiative	273
Cystic Fibrosis (n=91)	Early Pseudomonas Infection Control cystic fibrosis cohort	74
	University of North Carolina cystic fibrosis cohort	17
Chronic obstructive pulmonary disease (n=89)	Lung Health Study	88
	COPD Gene	1
Pulmonary hypertension (n=79)	Pulmonary Arterial Hypertension	79
LDL (n=393)	Atherosclerosis Risk in Communities	292
	Cardiovascular Health Study	24
	Framingham Heart Study	25
	Jackson Heart Study	52
Body mass index and type 2 diabetes (n=450)	Women's Health Initiative	450
Early-onset stroke (n=244)	Women's Health Initiative	244

## Cohort descriptions

**Women's Health Initiative (WHI):** The WHI is a long-term national health study that has focused on strategies for preventing heart disease, breast and colorectal cancer, and osteoporotic fractures in postmenopausal women. The original WHI study included 161,808 postmenopausal women enrolled between 1993 and 1998. The Fred Hutchinson Cancer Research Center in Seattle, WA serves as the WHI Clinical Coordinating Center for data collection, management, and analysis of the WHI.

**Framingham Heart Study (FHS):** FHS is a study to identify the common factors or characteristics that contribute to cardiovascular disease (CVD) by following its development over a long period of time in a large group of participants who had not yet developed overt symptoms of CVD or suffered a heart attack or stroke. In 1948, researchers recruited 5,209 men and women between the ages of 30 and 62 from the town of Framingham, Massachusetts, and began the first round of extensive physical examinations and lifestyle interviews that they would later analyze for common patterns related to CVD development. Since 1948, the subjects have returned to the study every two years for an examination consisting of a detailed medical history, physical examination, and laboratory tests. In 1971, the study enrolled a second-generation cohort, 5,124 of the original participants' adult children and their spouses, to participate in similar examinations. The second examination of the offspring cohort occurred eight years after the first examination, and subsequent examinations have occurred approximately every four years thereafter. In April 2002, the Study entered a new phase: the enrollment of a third generation of participants, the grandchildren of the original cohort. The first examination of the Third Generation Study was completed in July 2005 and involved 4,095 participants. Thus, the FHS has evolved into a prospective, community-based, three-generation family study. The FHS is a joint project of the National Heart, Lung and Blood Institute and Boston University.

**Jackson Heart Study (JHS):** JHS is a large, community-based, observational study whose participants were recruited from urban and rural areas of the three counties (Hinds, Madison, and Rankin) that make up the Jackson, Mississippi, metropolitan statistical area (MSA). Participants were enrolled from each of four recruitment pools: random, 17%; volunteer, 30%; currently enrolled in the Atherosclerosis Risk in Communities (ARIC) Study, 22%; and secondary family members, 31%. Recruitment was limited to non-institutionalized adult African Americans 35-84 years old, except in the family cohort where those 21 to 34 years of age were eligible. The final cohort of 5,301 participants includes 6.59% of all African American Jackson MSA residents aged 35-84 (N=76,426, US Census 2000). Major components of each exam include medical history, physical examination, blood/urine analyses and interview questions on areas such as: physical activity; stress, coping and spirituality; racism and discrimination; socioeconomic position; and access to health care. At 12-month intervals after the baseline clinic visit (Exam 1), participants are contacted by telephone to update information, confirm vital statistics, document interim medical events, hospitalizations, and functional status, and obtain additional sociocultural information. Questions about medical events, symptoms of cardiovascular disease, and functional status are

repeated annually. Ongoing cohort surveillance includes abstraction of medical records and death certificates for relevant International Classification of Diseases (ICD) codes and adjudication of nonfatal events and deaths.

**Atherosclerosis Risk in Communities (ARIC):** The ARIC Study, sponsored by the National Heart, Lung and Blood Institute (NHLBI), is a prospective epidemiologic study conducted in four U.S. communities. The four communities are Forsyth County, NC; Jackson, MS; the northwest suburbs of Minneapolis, MN; and Washington County, MD. ARIC is designed to investigate the etiology and natural history of atherosclerosis, the etiology of clinical atherosclerotic diseases, and variation in cardiovascular risk factors, medical care and disease by race, gender, location, and date.

ARIC includes two parts: the Cohort Component and the Community Surveillance Component. The Cohort Component began in 1987, and each ARIC field center randomly selected and recruited a cohort sample of approximately 4,000 individuals aged 45-64 from a defined population in their community. A total of 15,792 participants received an extensive examination, including medical, social, and demographic data. These participants were reexamined every three years with the first screen (baseline) occurring in 1987-1989, the second in 1990-1992, the third in 1993-1995, and the fourth and last exam was in 1996-1998. Follow-up occurs yearly by telephone to maintain contact with participants and to assess health status of the cohort.

**Cardiovascular Health Study (CHS):** The CHS is a study of risk factors for development and progression of coronary heart disease (CHD) and stroke in people aged 65 years and older. The objectives of the Cardiovascular Health Study are to: 1) quantify associations of conventional and hypothesized risk factors with CHD and stroke; 2) assess the associations of non-invasive measures of subclinical disease with the incidence of CHD and stroke; 3) quantify the associations of risk factors with subclinical disease; 4) characterize the natural history of CHD and stroke, and identify factors associated with clinical course; and 5) describe the prevalence and distributions of risk factors for CHD and stroke, non-invasive measures of subclinical disease, and clinical endpoints of CHD and stroke. The study originated in 1988 from the recommendations of an NHLBI workshop on the management of CHD in the elderly. This is the most extensive study undertaken by the NHLBI to study CVD exclusively in an elderly population. Initially funded for six years, it was renewed for a second six-year period in 1994, and recently was renewed for continued morbidity and mortality follow-up. The 5,888 study participants were recruited from four U.S. communities and have undergone extensive clinic examinations for evaluation of markers of subclinical cardiovascular disease. The original cohort totaled 5,201 participants. A new cohort was recruited in 1992. The 687 participants in the new cohort are predominately African-American and were recruited at three of the four field centers.

**Cleveland Clinic GeneBank:** The Cleveland Clinic GeneBank is a prospective cohort based study that enrolled patients undergoing elective diagnostic coronary angiography.

**Heart Attack Risk in Puget Sound (HARPS):** HARPS is a population-based case-

control study that enrolled cases with incident myocardial infarction (MI) presenting to a network of hospitals in Washington State. In HARPS, eligible cases were men <50 and women < 60 years of age.

**Massachusetts General Hospital - Premature Coronary Artery Disease (MGH-PCAD):** MGH-PCAD is a hospital-based case-control study that enrolled cases who were hospitalized with MI at MGH. In MGH-PCAD, eligible cases were men <50 and women < 60 years of age.

**Penn-Cath:** Penn-Cath is a catheterization-lab based cohort study from the University of Pennsylvania Medical Center and enrolled subjects at the time of cardiac catheterization and coronary angiography.

**Translational Research Investigating Underlying Disparities in Acute Myocardial Infarction Patients' Health Status (TRIUMPH):** TRIUMPH is an observational, multi-center prospective registry that enrolled subjects presenting with MI at participating medical centers.

**Lung Health Study (LHS):** The LHS is a randomized multicenter clinical trial with 5,887 participants carried out from October 1986 to April 1994, designed to test the effectiveness of smoking cessation and bronchodilator administration in smokers aged 35 to 60 with mild lung function impairment. Participants were randomly assigned to one of three groups: (1) usual care, who received no intervention; (2) smoking intervention with the inhaled bronchodilator ipratropium bromide; or (3) smoking intervention with an inhaled placebo. The effect of intervention was evaluated by the rate of decline of forced expiratory volume in one second (FEV1).

**Pulmonary Arterial Hypertension (PAH):** The PAH cohort was comprised of idiopathic PAH (IPAH) and scleroderma-associated PAH (PAH-SSc) cases and healthy controls in the Hopkins SCCOR program. Pulmonary hypertension was defined in IPAH or PAH-SSc patients as a mean pulmonary artery pressure greater than 25 mm Hg proven by right heart catheterization. For patients with scleroderma, the presence of disease was defined as systemic sclerosis with diffuse or limited scleroderma meeting the American College of Rheumatology criteria.

**Cystic Fibrosis (CF) cohorts (2):** CF cases were ascertained from two separate cohorts. The first cohort consisted of 1,704 cases that represent participants in the Early *Pseudomonas* Infection Control (EPIC) Observational Study, a multicenter, longitudinal, prospective cohort of early lung disease in young CF patients. The second CF cohort consisted of 1,208  $\Delta F508$  homozygotes who are at the extremes of lung disease severity ("severe", worst 25<sup>th</sup> percentile of birth cohort vs. "mild", best 25<sup>th</sup> percentile) based on ~22 measures of lung function for each patient (> 5 years of age) developed at the University of North Carolina.

## Informed consent

All study participants in each of the component studies provided written informed consent for the use of their DNA in studies aimed at identifying genetic risk variants for disease and for broad data sharing. Institutional certification was obtained for each sample to allow deposition of phenotype and genotype data in dbGaP and BAM files in the short-read archive.

## Exome sequencing

### University of Washington exome sequencing / analysis

#### QC of sample DNA

Initial quality control (QC) performed on all samples included sample quantification (PicoGreen; Life Technologies, Grand Island, NY), confirmation of high-molecular weight DNA, test PCR amplification (four amplicons), and sex determination using a Taq-man assay (31). Prior to preparation for exome sequencing, all samples were genotyped (Illumina BeadXpress; Illumina, San Diego, CA) for 96 high frequency (30-50% minor allele frequency) exome-specific SNPs derived from the content found on large-scale (GWAS) genotyping chips from both Illumina and Affymetrix. Exome data at these variant sites were used to ensure sample tracking integrity through sample preparation and the sequencing pipeline. Samples failed QC if: the total mass, concentration, or integrity of DNA was low; low (<90%) genotype call rates were observed; sex-typing was inconsistent with the sample manifest. Following QC, all 3.5 ug of genomic DNA was reformatted into 96 well plates for library preparation and exome capture.

#### Library production and exome capture

All protocols for library construction and exome capture were automated on a Perkin-Elmer Janus II liquid handling robot or multi-channel pipettors, and performed in 96-well plate format. Samples were prepared by subjecting ~3.5 ug of genomic DNA to a series of shotgun library construction steps, including fragmentation through acoustic fragmentation (Covaris, Woburn, MA), end-polishing and A-tailing, ligation of sequencing adaptors, and PCR amplification. Sample shotgun libraries were captured for exome enrichment using one of three in-solution capture products: CCDS 2008 (~26 Mb), Roche/Nimblegen SeqCap EZ Human Exome Library v1.0 (~32 Mb; Roche Nimblegen EZ Cap v1), or EZ Cap v2 (~34 Mb). Briefly, 1 µg of shotgun library was hybridized to biotinylated capture probes for 72 hours and recovered via streptavidin beads. Unbound DNA was washed away, and the captured DNA PCR amplified. Following capture, washing, and PCR, libraries were assessed again on the Bioanalyzer (Agilent, Santa Clara, CA) for concentration, molecular weight distribution, and the presence of PCR artifacts. The fragment size distributions of the libraries were highly consistent (typically  $125 \pm 15$  bp).

#### Clustering and sequencing

Library concentration and flow-cell loading cluster densities were determined using a standardized qPCR protocol (Kapa Biosystems, Woburn, MA). Using the automated Illumina cBot cluster station, non-multiplexed samples were processed in batches of eight (one for each lane of the flow-cell), diluted, and denatured to their final effective loading concentrations. Hybridization was followed by cluster generation via bridge PCR as per standard protocols (Illumina). Enriched libraries were sequenced on an Illumina GAIIx or HiSeq 2000 using either paired-end 76 base or 50 base runs, respectively.

### **Read mapping and variant analysis**

Samples were processed from real-time base-calls (RTA 1.7 software [Bustard], converted to qseq.txt files, and aligned to a human reference (hg19) using BWA (Burrows-Wheeler Aligner; 32). Read-pairs not mapping within  $\pm 2$  standard deviations of the average library size ( $\sim 125 \pm 15$  bp for exomes) were removed. Data was processed using the Genome Analysis ToolKit (GATK refv1.2905; 33). All aligned read data was subjected to “duplicate removal”, indel realignment (GATK IndelRealigner), and base quality recalibration (GATK TableRecalibration). Variant detection and genotyping were performed using the UnifiedGenotyper (UG) tool from GATK, and only performed on the targeted exome regions. Variant data for each sample was formatted (variant call format [VCF]) as “raw” calls for all samples, and lower quality/false positives sites were flagged using the filtration walker (GATK) (low quality scores ( $\leq 50$ ), allelic imbalance ( $\geq 0.75$ ), long homopolymer runs ( $> 3$ ), and/or low quality by depth  $< 5$ ). Samples were considered complete when exome targeted read coverage was  $> 8x$  over  $> 90\%$  of the exome target.

### **Data analysis QC**

All individual exome sequencing data were evaluated against QC metrics including assessment of: (1) total reads: a minimum of 30M PE reads; (2) library complexity: the ratio of unique reads to total reads mapped to target; (3) capture efficiency: the ratio of reads mapped to target versus the reads mapped to human; (4) coverage distribution: 90% at  $\geq 8x$  required for completion; (5) capture uniformity; (6) raw error rates; (7) Ti/Tv ratio (3.2 for known sites and 2.9 for novel sites); (8) distribution of known and novel variants relative to dbSNP; (9) fingerprint concordance  $> 99\%$ ; (10) homozygosity; (11) heterozygosity. All QC metrics for both single-lane and merged data were reviewed to identify data deviations from known or historical norms. Lanes/samples that failed QC were re-queued for library prep or further sequencing.

## **Broad Institute exome sequencing / analysis**

### **Receipt/QC of sample DNA**

Samples were shipped to the Biological Samples Platform laboratory at the Broad Institute of MIT and Harvard. DNA concentration was determined by the Picogreen assay before storage in 2D-barcoded 0.75 mL Matrix tubes at  $-20^{\circ}\text{C}$  in the SmarTStore™ (RTS, Manchester, UK) automated sample handling system. We performed initial QC on all samples involving sample quantification (PicoGreen), confirmation of high-molecular weight DNA and fingerprint genotyping, and sex determination (Illumina iSelect). Samples failed if the total mass, concentration, integrity of DNA or quality of preliminary genotyping data was too low.



### **Library construction and in-solution hybrid selection**

Starting with 3µg of genomic DNA, library construction and in-solution hybrid selection were performed as described previously (34). A subset of samples were, however, prepared using the Fisher et al. protocol (34) with slight modifications. Initial genomic DNA input into shearing was reduced from 3µg to 100ng in 50µL of solution. In addition, for adapter ligation, Illumina paired end adapters were replaced with palindromic forked adapters with unique eight base index sequences embedded within the adapter.

### **Preparation of libraries for cluster amplification and sequencing**

After in-solution hybrid selection, libraries were quantified using qPCR (KAPA Biosystems) with probes specific to the ends of the adapters. This assay was automated using Agilent's Bravo liquid handling platform. Based on qPCR quantification, libraries were normalized to 2nM and then denatured using 0.1 N NaOH using Perkin-Elmer's MultiProbe liquid handling platform. A subset of the samples prepared using forked, indexed adapters was quantified using qPCR, normalized to 2nM using Perkin-Elmer's Mini-Janus liquid handling platform, and pooled by equal volume using the Agilent Bravo. Pools were then denatured using 0.1 N NaOH. Denatured samples were diluted into strip tubes using the Perkin-Elmer MultiProbe.

### **Cluster amplification and sequencing**

Cluster amplification of denatured templates was performed according to the manufacturer's protocol (Illumina) using either Genome Analyzer v3, Genome Analyzer v4, or HiSeq 2000 v2 cluster chemistry and flowcells. After cluster amplification, SYBR Green dye was added to all flowcell lanes, and a portion of each lane visualized using a light microscope, in order to confirm target cluster density. Flowcells were sequenced either on Genome AnalyzerII using v3 and v4 Sequencing-by-Synthesis Kits, then analyzed using RTA v1.7.48, or on HiSeq 2000 using HiSeq 2000 v2 Sequencing-by-Synthesis Kits, then analyzed using RTA v1.10.15. All samples were run on 76 cycle, paired end runs. For samples prepared using forked, indexed adapters, Illumina's Multiplexing Sequencing Primer Kit was also used.

### **Read mapping and variant analysis**

Samples were processed from real-time base-calls (RTA 1.7 software [Bustard], converted to qseq.txt files, and aligned to a human reference (hg19) using BWA (32). Aligned reads duplicating the start position of another read were flagged as duplicates and not analyzed ("duplicate removal"). Data were processed using the Genome Analysis ToolKit (GATK v1.1.3; (33)). Reads were locally realigned (GATK IndelRealigner) and their base qualities were recalibrated (GATK TableRecalibration). Variant detection and genotyping were performed on both exomes and flanking 50 bp of intronic sequence using the UnifiedGenotyper (UG) tool from the GATK. Variant data for each sample was formatted (variant call format [VCF]) as "raw" calls for all samples. SNP and Indel sites were flagged using the Variant Filtration walker (GATK) to mark sites of low quality that were likely false positives. SNPs were marked as potential errors if they exhibited strong strand bias (SB  $\geq 0.10$ ), low average quality (QD  $< 5.0$ ), or fell in a homopolymer run (HRun  $> 4$ ). Indels were marked as potential errors for low quality (QUAL  $< 30.0$ ), low average quality (QD  $< 2.0$ ), or if the site exhibited strong

strand bias ( $SB > -1.0$ ). Samples were considered complete when exome targeted read coverage was  $\geq 20\times$  over  $\geq 80\%$  of the exome target.

### **Data analysis QC**

Sequence data that was processed were required to match known fingerprint genotypes for their respective samples, and to achieve the sequence coverage noted above ( $20\times$  over 80% target). Variant calls were evaluated on both bulk and per-sample properties: novel and known variant counts, Ti/Tv ratio, Het/Hom ratio, and Deletion/Insertion ratio. Both bulk and sample metrics were compared to historical values for exome sequencing projects at the Broad. No significant deviation of the ESP calls or ESP samples from historical values were noted.

## **Variant calling**

### **Input data**

Variant discovery and genotyping was performed simultaneously across the initial set of 2,520 exomes sequenced by the NHLBI Exome Sequencing Project. These consisted of 1,384 samples sequenced at the University of Washington and 1,136 samples sequenced at the Broad Institute. Most sequenced samples were of self-identified European or African American Ancestry, although self-identified ancestry was not available for a subset of samples. Sequencing centers carried out initial alignment and pre-processing for each sample. Briefly, read data were aligned to NCBI build 37 reference using BWA (32) and indel-realigned and recalibrated using the Genome Analysis Toolkit (GATK; 33). The initial pre-processing and alignment steps carried out at each center were broadly similar.

As indicated above, three different target solutions were used for exome capture across the 2,520 samples: Broad used Agilent capture solution (Agilent SureSelect Human All Exon Kit v2 [31 Mb]) and UW used two different targets of Nimblegen-designed solutions (custom RefSeq/CCDS design [28 Mb]; SeqCap EZ v1 [32 Mb]). In a preliminary analysis of chromosome 20, we found variants called within 50 bp of target regions had similar transition-transversion ratios, variant density, and overlap with dbSNP compared to variants called on-target; thus we attempted to call variants within 50bp of target regions in each sample. When X chromosome variant calls were generated, sex information was available only for a subset of samples; hence we used normalized read depths across the sex chromosomes to infer sex for the remaining samples, identifying 774 samples as male and 1,746 samples as female.

### **Calling pipeline overview**

Briefly, we took a two-step approach for discovering variants and genotyping candidate sites. First, genotype likelihood files (GLFs) were generated using SAMtools pileup on individual BAM files. Next, we used glfMultiples - a multi-sample variant caller - to generate initial SNV calls. This two-step approach is very scalable and is especially beneficial for exome sequence data, because the size of GLF files largely depends on the number of sites targeted for analysis ( $\sim 52.6$  M sites are within 50 bp of one of the

three ESP capture reagents). This process allowed us to perform multisample calling of variants simultaneously across thousands of samples to avoid batch effects.

For the initial variant calls, we assumed that the prior probability that a site was polymorphic was  $Prior(SNP) = \theta \sum_{i=1}^{2n} \frac{1}{i} = 0.0091$ , corresponding to an estimate of the prior for a segregating site in a simple population genetics model and an estimated per sample per base pair heterozygosity of  $\theta = 0.001$ . We assumed transition and transversion mutations were equally likely – while this assumption is not optimal, it makes it easier to use transition-transversion ( $T_r/T_v$ ) ratio as a diagnostic of SNP call quality. Details of the likelihood model implemented in glfMultiples are given in (35) in the section entitled “Identifying Potential Polymorphic Sites”. Briefly, the distribution of observed bases and quality scores at each location (conditional on the true genotype) is modeled using the MAQ model (36). Using maximum likelihood, we then estimated an allele frequency for each site under the assumption that genotypes were segregating in Hardy-Weinberg proportions. Genotypes for each site and corresponding posterior probabilities were then calculated via Bayes theorem, using the Hardy-Weinberg proportions as priors and the MAQ error model to describe the conditional probability of bases and quality scores given the true genotype.

### Initial variant filtering

After these initial SNV calls were generated, we re-examined the BAM files to collect additional information about each variant site. Filters considered the total read depth, the number of individuals with coverage at the site, the fraction of variant reads in each heterozygote, the ratio of forward and reverse strand reads for reads carrying reference and variant alleles and the average position of variant alleles along a read.

The final SNV call set included variants that were called with posterior probability >99% (glfMultiples SNP quality >20), were >5bp away from an indel detected in the 1000 Genomes Pilot Project (37), were covered by at least one read in 85% of samples, and had total depth across samples of between 2,500 and 2,500,000 reads (~1-100 reads per sample). Sites having >65% of reads as heterozygotes carrying the variant allele or where the absolute squared correlation between allele (variant or reference) and strand (forward or reverse) was >0.15 were excluded. Additional quality metrics were examined, including the fraction of unexpected alleles (those that did not match either the reference or alternate allele) and the average location of variant calls along each read. These metrics were not used to filter this initial dataset. All processing steps were automated using GNU make, which facilitates parallelization of data processing steps, dependency checking, and repeats of failed runs.

## Data filtering

The lightly filtered data set described above was further subjected to more stringent filtering in order to obtain high quality genotypes suitable for population genetics analyses. Because three distinct capture targets were used in sequencing, we restricted all analyses to the intersection of the capture targets; that is, sites not appearing on all

three targets were excluded. We removed individual genotypes if quality (GQ) was under 30 and/or filtered depth (DP or GD) was under 10. After such filtering, if more than 10% of individuals had a missing genotype, we excluded the entire site from analysis. We identified clusters of African Americans and European Americans based on ancestry-informative markers and mean heterozygosity per individual (see below). We used these clusters to remove sites violating Hardy-Weinberg equilibrium (chi-square p-value  $< 10^{-6}$ ) with heterozygote or homozygote excess within either cluster. Similarly, we removed sites with overall heterozygote excess in the full dataset, but not sites with overall homozygote excess as these could represent population structure. After filtering Ti/Tv ratios for synonymous, missense, nonsense, and splice variants were 5.60, 2.31, 2.13, and 1.69, respectively.

## Identifying related individuals

We used KING (38) to identify duplicate samples and kin with a third degree (e.g. first cousin) or closer relationship (Figure S8). Five exomes were removed due to low coverage and 15 exomes showed high and presumably spurious relatedness to a large number of other samples, likely due to low quality or contamination; we conservatively removed these exomes from subsequent analyses as well. For all duplicate (N=22) and kinship pairs (N=48), the individual with lower mean filtered depth was removed. We removed 60 samples as some individuals have multiple relationships (Figure S8).

## Assignment of ancestry labels

Designated ancestry was determined using heterozygosity and ancestry-informative markers (AIMs), which we defined as autosomal SNVs estimated from preexisting datasets to exhibit large allele frequency differences between individuals of European and African ancestry (Figure S19). We chose 29 exomic AIMs varying between Europeans and Africans in the HapMap dataset, and 17 exomic AIMs varying between those continents and Native Americans in the HGDP dataset. The one self-identified Native American in our dataset was a clear outlier, being the only individual matching over 40% of Native American AIMs and fewer than 65% European AIMs, and was designated as Native American. The remaining samples formed two very distinct clusters: 97% of samples had either >65% European AIMs and 13-15% of SNVs heterozygous, or else < 50% European AIMs and 16.5-18.5% of SNVs heterozygous. In order to assign the few ambiguous individuals falling between clusters, we designated all individuals with over 50% European AIMs and <16% of SNVs heterozygous as European American, and all other individuals as African American. This approach resulted in less ambiguous ancestry designation compared to a simple PC analysis.

## Data annotation

We used SeattleSeq Annotation 131 to annotate all SNVs, with respect to site type, presence in existing databases, and other factors.

## Prediction of functional variation

To identify putatively functional variation, we used PolyPhen2 (39), SIFT (40), a likelihood ratio test (41), Mutation Taster (42), GERP (43), PhyloP (44), and a novel population genetics approach that combines conservation information with the SFS that we designate SFS-Del (see below). Thresholds in determining whether a given metric predicted a SNV to be functional were as follows: PolyPhen2 “Probably Damaging”, SIFT “Damaging”, likelihood ratio test “Deleterious”, MutationTaster “disease causing automatic” or “disease causing”,  $GERP \geq 5$ ,  $PhyloP \geq 3$ , SFS-Del probability  $\geq 0.90$ .

## Analyses of batch and cohort effects

Exome sequences were assessed for batch and cohort effects using several strategies. First, among individual exomes all samples showed similar properties of Ti/Tv, heterozygosity and frequency of missing genotypes (Figure S1). Second, no substantial differences were observed among sets of exomes sequences defined by target definition, target group, and sequencing center vs. heterozygote / homozygous non-reference (Figure S2) and Ti/Tv (Figure S3). Third, we assessed variant QC metrics including depth of coverage, frequency of missing genotypes, and minor allele frequency (MAF). Ti/Tv was relatively constant among a broad range of coverage, frequencies of missing data  $< 10\%$ , and MAF (Figure S4). Last, to confirm that patterns of population structure and other population parameters were not significantly biased by sequencing strategy, we relabeled the PCA analyses of population structure by: exome capture target (Figure S5), sequencing center (Figure S6), and cohort/phenotype information (Figure S7) for each sample. Plots are based on variants with a MAF  $> 5\%$  and pruned by LD.

## Validation

### Selection of variants for validation

Single nucleotide variants (SNVs) were selected for validation if they were novel – i.e. not present in dbSNP 131, a single base change (i.e. non-indel), biallelic, and functionally annotated as a missense or nonsense variant. This resulted in a total set of 247,906 variants over a frequency range from 1 observation (i.e., a singleton; 1/4880 chromosomes) to more common variants at a frequency  $>10\%$ . From this set, we selected novel variants for validation including: 400 singletons, 768 variants  $<10\%$ , and all 52 variants identified with  $>10\%$  MAF. SNVs selected for validation were not allowed to be within 100,000bp of each other, nor within 100bp of any other known SNV. These variants were evenly split between the sequencing centers for testing.

### University of Washington genotyping/validation

Directed PCR and Sanger sequencing were used to confirm all singleton variants (n=200). Primer design, optimization, and testing resulted in 146/200 passing amplicons. For each variant position, we selected the specific genotype positive singleton sample, along with six Coriell negative control samples. All chromatograms were generated on an ABI 3700XL, and aligned to a genomic reference sequence. Polyphred was used to identify putative variant positions within an amplicon and the specific position of interest was directly examined and validated by a data analyst. All samples were genotyped using PolyPhred or manually. Comparison of exome genotypes against the Sanger data resulted in 2/146 (1.3% discordant) genotypes. Of the singletons that failed confirmation, one variant was annotated by the filtration walker (GATK) as failing allelic balance (GATK filter status was not originally used as a selection criteria). The second singleton that could not be confirmed was supported by manual review of the exome data, however multiple PCR amplicons and sequencing did not confirm this sample genotype.

Singleton variants were validated by Sanger sequencing. Of 145 singleton variants that were assessed, 143 (99%) were validated. Non-singleton variants were validated using the Illumina BeadXpress (n=384 plex) assay. All 820 variants (768 < 10% MAF + 52 > 10% MAF) were submitted for design scores and the highest ranking variants used in a 384 plex genotyping pool (372 < 10% MAF + 12 > 10% MAF). All BeadXpress sample pools (10 pools x 96 samples) were scanned and initially clustered/genotyped using the BeadStudio software, with subsequent manual review. We selected 960 samples from the genotype positive sample set to be assayed across all 384 variants. 934/960 (98%) samples passed QC (>95% of variants called) and were used for validation. We found 332/384 (86%) variants passing QC (>90% of samples reporting a genotype) of which 316/332 (95%) of the variants were confirmed in the genotype positive samples. Of 332 non-singleton variants tested, 323 had a MAF < 10%. We validated 316/323 (98%) non-singleton variants with 1191/1255 (95%) confirmed genotypes. The observed confirmation failures (5%) appeared to be due to genotype clustering, with unconfirmed samples falling at or near cluster boundaries. Nine SNVs with MAF>10% were tested and all nine were validated with 1914/1925 (99%) confirmed genotypes.

### **Broad genotyping/validation**

Singleton and non-singleton variants were validated using iplex technology (Sequenom; John Hopkins Court, San Diego, CA). In all, 97 putative singletons identified in 94 DNA samples were genotyped in 8 SNP pools. Each SNP was genotyped on the specific genotype positive singleton samples and the other 93 DNA samples from which the initial pool of singletons were selected. Four of the sites demonstrated high missingness (a no-call rate of >5%) and were excluded from further analysis. Of the remaining 93 singletons, 91 (98%) were validated. The two putative singletons that could not be detected were determined to be monomorphic and this was confirmed by visual analysis of the raw intensity plots. For validation of non-singletons, 402 mis- or non-sense SNPs with  $MAF \leq 10\%$  were selected for testing. All 402 SNPs were sequenced in 1024 samples from the genotype positive sample set. 352 of the 402 SNPs (88%) tested passed initial QC with 331 confirmed in the genotype positive samples. Of 391 non-singleton variants tested, 370 had a MAF < 10%. We validated 362/370 (98%) non-

singleton variants. Twenty-one SNVs with MAF>10% were tested and 20/21 (95%) were validated.

## Power of variant discovery

We used two simple approaches for calculating the power of variant discovery. First, the probability that a variant with a population frequency  $f$  is observed in a sample of  $n$  individuals is  $1-(1-f)^{2n}$  (45). Second, we used the hypergeometric distribution to calculate the probability of sampling one or more alleles of a variant drawn from a population of size  $N$  that consists of  $m$  minor alleles,  $N-m$  major alleles (i.e., the more common allele), in a sample size of  $k$  chromosomes is  $1 - \text{hyper}(0, m, N-m, k)$ . We evaluated the probabilities of variant discovery for a range of allele frequencies and population sizes,  $N$ . Note, both of these approaches implicitly assume a randomly mating population and therefore should be regarded as rough approximations on the power of variant discovery.

## Comparison of variants discovered in the ESP with the 1000 Genomes Project

SNVs were compared between the ESP and 1000 Genomes Project (1KGP) for individuals of European ancestry (1000 Genomes Project Consortium; 7). In total, 1,351 ESP EA individuals were compared to 381 1KGP European individuals. The 1KGP European populations were: 87 Utah residents with Northern and Western European ancestry (CEU), 93 Finnish in Finland (FIN), 89 British from England and Scotland (GBR), 14 Iberian populations in Spain (IBS), and 98 Toscani in Italia (TSI). Data was downloaded from the 1KGP-Interim Phase I data (June 2011 release [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim\\_phase1\\_release/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20101123/interim_phase1_release/)).

The comparison focused on 140,323 autosomal regions (~23.92 Mb), with sequence coverage rates > 90% in the ESP. Among these regions, 261,251 SNVs were identified in the 1,351 EA ESP individuals and 90,961 SNVs were identified in the 381 European individuals from 1KGP (Table S1). Among these variants, 59,759 were identified in both the ESP and in 1KGP. For SNVs with a MAF  $\leq$  5%, 201,379 variants were uniquely observed in the ESP, compared with 25,641 variants uniquely observed in 1KGP. For SNVs with a MAF > 5%, 113 SNPs were uniquely observed in ESP, compared with 5,561 variants uniquely observed in 1KGP.

Among the 5,561 common SNVs that were not observed in ESP, 4,893 (87.98%) sites were filtered out in the ESP because >10% of individuals were called as missing (GQ<30 and/or DP<10), denoted as C90; 322 (5.79%) sites were filtered out in ESP because of departure from HWE in EAs (chi-square > 24), denoted as HWD; 151 (2.72%) sites were invariant (i.e., genotypes of all individuals are reference homozygous) in ESP, noted as INV; 120 (2.16%) sites were not called or were not targeted, noted as NoC; and 75 (1.35%) sites were filtered out for other reasons in ESP, such as allele balance, nearby indels, etc, denoted as OTH.



Figure S10 illustrates characteristics, such as reference allele frequency and HWE, of common SNVs specific to the 1KGP data. SNVs filtered out of the ESP data because of HWD or NoC (Figure S10B and D) also exhibit extreme departures from HWE in 1KGP. This departure from HWE in 1KGP is not explained by population structure compared with Figure S9A. SNPs filtered out of the ESP data because they were invariant (INV) show a higher frequency of the reference allele in 1KGP (Figure S9C), some of which were only observed in specific European populations, such as FIN, GBR and IBS (Table S2).

Thus, of the 5,561 common SNVs in 1KGP that were not observed in ESP, all but 151 variants can be accounted for by being filtered out of the ESP data or not present in the capture target. The remaining 151 SNPs can largely be accounted for by European population stratification (Table S2), as the 1KGP included a broader sampling of European populations.

## Demographic Inference

Demographic inference based on the site frequency spectrum was performed using the package *dadi version 1.5.2* (46), using a demographic model similar to a previously described three-population model (46,47), but with three added parameters to allow for a recent growth, namely a time for growth change and two recent growth rates (one in the European and one in the African population).

Because we want to resolve the site frequency spectrum for population frequencies of the order of 0.0005, and *dadi* calculates the site-frequency spectrum by numerically solving diffusion equations, we require a very fine discretization of frequency space, which is computationally demanding and makes parameter inference challenging. As we expect the extra information provided by the current data to be particularly informative of recent migration, we fixed the parameters of ancient demographic events to those obtained in (47), allowing only the recent European and African demography to change, namely the two recent growth parameters in Europeans and Africans, the time for the accelerated growth, and the growth rate in Europeans since the split with the ancestors to the Asian population. Even with these constraints, the requirements for simulating a 3-population model at this level of accuracy are high. Since the current data does not include individuals of reported Asian ancestry, we did not simulate the Asian population in this model after the split from the European population. Because the original model includes migration between all populations, we verified that this omission did not overly affect the joint African-American/European site-frequency spectrum--the maximum difference between two-population and three-population models in any of the bins was within 21.2% of the expected Poisson sampling noise in the data, when the simulation was carried with 100 grid points, and this number decreased to 16% with 500 grid points (we expect this number to decrease further when convergence is reached).

For the optimization process, we used grids of 500, 750, and 1000 points per population and extrapolated to infinite sample size using the method described in (46). Once we



obtained optimal parameters, we used an even finer grid (1000, 1500, and 2000) to generate plots of frequency spectra and perform comparisons with the real data.

To obtain confidence intervals for the recent growth parameters, we performed a bootstrap analysis over exons. More specifically, we first defined 8239 target regions separated by more than 50kb, by merging together target regions separated by a smaller distance. Then, we generated 24 bootstrapped samples by selecting 8239 genomic regions with replacement among the initial regions. For each bootstrap sample, we calculated the allele frequency spectrum, and optimized the likelihood as a function of the four recent demographic parameters, as for the initial point estimates. Convergence of the optimization was ensured by performing six optimizations from different starting positions and verifying that the highest-likelihood values were attained for at least two initial conditions. For computational efficiency, we used coarser grids than in the point estimate (extrapolating from results 100, 200, and 400 grid points per population), as we expect the width of the confidence intervals to be less affected than point estimates by finite grid effects. We use the obtained standard deviation to estimate confidence interval for the point estimate reported in the text.

## Analysis of natural selection

### Inference of positive selection

All positive selection analyses were restricted to the autosomes, so that all compared genes shared the same demographic history. To find genes showing evidence of positive selection in humans, we first identified human-chimp fixed differences by comparison to panTro3, excluding sites that were polymorphic in our dataset. We retained only human-specific substitutions by excluding sites where the human allele was identical to macaque (rheMac2); sites where neither the chimp nor the human allele matched macaque were assigned a 50% probability of being a human-specific substitution. We used these sites to identify genes with high ratios of human-specific divergence to within-human polymorphism, as well as genes with high ratios of nonsynonymous/synonymous divergence relative to their ratio of nonsynonymous/synonymous polymorphism. For all measures of within-human polymorphism, we used  $\pi$  rather than counts of segregating sites in order to minimize the signal from rare and likely deleterious variants. To find genes showing evidence of population-specific adaptation, we identified candidates with unusually high or low ratios of AA polymorphism to EA polymorphism. We also identified genes with at least one SNV showing  $F_{ST}$  at least 0.3 and either AA polymorphism or EA polymorphism in the lowest 20% for that sample.

### Inferences of purifying selection based on the SFS

We inferred purifying selection by examining the joint SFS (in EA and AA) of different site classes. We classified sites according to four basic type categories (synonymous, missense, nonsense/splice, and other) and GERP score rounded to the nearest integer. We assumed that synonymous sites with a GERP score of -5 or less (site class S) represented a neutral class of sites with a joint SFS reflecting only demographic history. We assumed that the effects of positive and balancing selection on segregating

polymorphism are negligible for the exome as a whole. We compared joint SFS for sites S against all non-S site classes, assuming that any differences are due to the presence of deleterious variants (D) in the non-S class, in addition to the neutral non-S variants (N) which behave just like S. We assumed that all sites which reach a frequency of 50% in both EA and AA are effectively neutral, so all non-S variants are N. For a given joint frequency rarer than 50% ( $F_{EA}$  and  $F_{AA}$ ) the ratio of non-S deleterious SNVs (D+N) to selectively neutral SNVs (S) is can be approximated by the ratio at 50% (N/S) multiplied by several constants which scale with the logs of the frequencies and their interaction term:

$$\ln((D + N)/S) = \ln(N/S) + \beta_1(\ln(0.5/F_{EA})) + \beta_2(\ln(0.5/F_{AA})) + \beta_3((\ln(0.5/F_{EA})) * (\ln(0.5/F_{AA})))$$

We used multiple regression to solve for the  $\beta$  values for all site classes, combining some similar site classes to ensure a large sample size for all regression calculations. Thus, for a site in a given site class at a given frequency, we could estimate the probability that it was D rather than N. By summing these probabilities across the exome, once for every heterozygous or derived hemizygous variant and twice for every derived homozygous variant, we estimated the number of deleterious variants per individual. In effect, we estimated the number of variants that would be prevented by selection from reaching a 50% frequency in both populations if the populations continued to experience the same selective pressures and demographic processes as they have in the past. The advantage of this approach is that we have to make very few *a priori* assumptions about which site types are likely to be deleterious; rather we can empirically estimate a probability for each site class based on how disproportionately prevalent that site class is at increasingly rarer frequencies.

### **Purifying selection and protein structure**

We analyzed six features that are potentially important for protein structural stability: normalized solvent-accessible surface area, change in amino acid side chain volume, change in amino acid hydrophobic potential, average number of hydrogen bonds formed by the residue, and secondary structure element. We also analyzed three features that are potentially important for protein biochemical function: number of active site contacts, number of heteroatom contacts, and number of side chain contacts. Variants with normalized accessible surface area above the median were classified as surface variants; variants with normalized accessible surface area below the median were classified as buried variants. Variants with side chain volume change above the top quintile (large positive volume change) were classified as overpacking variants; variants with side chain volume change below the bottom quintile (large negative volume change) were classified as cavity-forming variants. Variants with hydrophobic potential change above the top quintile were classified as hydrophobic potential-changing variants. Secondary structure elements were grouped into helices, strands, and loops. Variants in sites that formed any hydrogen bonds were classified as hydrogen-bonding; variants in sites that made any active site contacts were classified as active site variants; variants in sites that made any heteroatom contacts were

classified as ligand-binding variants; and variants in sites that made any interchain contacts were classified as variants contacting another protein chain.

Nonsynonymous variants were mapped to PDB structures using PolyPhen2, using a minimum sequence identity of 50% and a minimum alignment length of 100 amino acids. Each variant that mapped to a structure was then scored for the nine features described above. For each of the categories described above, we performed a random permutation test to determine whether the category contains a statistically significant enrichment of rare variants. For each category, we chose 100,000 random sets of variants that mapped to PDB structures and counted the number of variants with a MAF < 0.005.

### **Purifying selection and synonymous variation**

The  $w$  scores of all possible human codons were taken from (48). The change in  $w$  score was calculated as the difference between that of the codon with the derived allele and that of the codon with the ancestral allele. The hypothesis that the correlation between the derived allele frequency (DAF) and the functional or  $w$  scores is equal to zero was tested using Fisher's transformation and a  $z$  score test.

### **Coordinated purifying selection on synonymous and nonsynonymous variation within genes**

For each gene we estimated Spearman's rank correlation coefficients (RCCs) between the average  $w$  score changes of synonymous variants and the average functional prediction scores (SIFT, Polyphen2, LRT and MutationTaster) of the nonsynonymous variants. We observed a weak but statistically significant negative correlation between  $w$  score changes and functional prediction scores from SIFT, Polyphen2 and MutationTaster (Table S7), which suggests potential coordinated purifying selection acting on both the synonymous variants and the non-synonymous variants within a gene.

## **Rare variant association power analysis**

The power of association studies on individual rare variants is extremely low, and numerous methods have been developed to pool, or aggregate, rare variants at a gene to improve power (28). Here, we used a simple Fisher's Exact Test (FET) to explore power on a gene by gene basis. In order to simulate phenotypes we used a logistic model of the form:

$$\text{logit}(p) = \beta \cdot (x_1 + x_2 \dots x_n) + \varepsilon$$

where the  $x$ 's are each individual variant in the genic region, and  $\beta$  was set to either 0.405 or 1.609, which leads to an OR of 1.5 and 5, respectively.  $\varepsilon$  is the random error term drawn from a standard normal distribution. The logistic model yields a continuous probability of being a case or control, which is subsequently determined using a uniformly distributed random number compared to the probability. The data were

preprocessed to include only 1000 EA and AA, which will be sub-selected for 400 cases and 400 controls; sites with a MAF < 1%, calculated at the 1000 individual level for each sample (ignoring those greater than 1% from random sampling); and sites with a functional designation of missense, nonsense, splice-site, or GERP  $\geq 5$ .

In an effort to remove spurious associations due to population demography and structure, the 1000 individuals for each population were selected in a PCA aware manner. For both samples, we use a global PCA of sites with a MAF > 5% and pruned for LD (plink --indep-pairwise 50 5 0.5) (see Figure S17). For the AA sample, we selected the 1,000 individuals with the lowest PC1 values, which is indicative of level of European admixture. With each analysis of AA, PC1 will be included as a covariate. For the EA sample, we selected the 1,000 individuals with the lowest PC2 values, which is indicative of a North/South European cline. This removes all the individuals that cluster with Eastern European Jewish populations and the most extreme individuals that cluster with Tuscans from HapMap3 (see below). The power results per gene for an OR = 5 are shown in the main text (Figure 4A) and for an OR = 1.5 in Figure S18.

## Further acknowledgements

### HeartGO:

**Atherosclerosis Risk in Communities (ARIC):** NHLBI (N01 HC-55015, N01 HC-55016, N01HC-55017, N01 HC-55018, N01 HC-55019, N01 HC-55020, N01 HC-55021); **Cardiovascular Health Study (CHS):** NHLBI (N01-HC-85239, N01-HC-85079 through N01-HC-85086, N01-HC-35129, N01 HC-15103, N01 HC-55222, N01-HC-75150, N01-HC-45133, and grant HL080295), with additional support from NINDS and from NIA (AG-023629, AG-15928, AG-20098, and AG-027058); **Coronary Artery Risk Development in Young Adults (CARDIA):** NHLBI (N01-HC95095 & N01-HC48047, N01-HC48048, N01-HC48049, and N01-HC48050); **Framingham Heart Study (FHS):** NHLBI (N01-HC-25195 and grant R01 NS17950) with additional support from NIA (AG08122 and AG033193); **Jackson Heart Study (JHS):** NHLBI and the National Institute on Minority Health and Health Disparities (N01 HC-95170, N01 HC-95171 and N01 HC-95172); **Multi-Ethnic Study of Atherosclerosis (MESA):** NHLBI (N01-HC-95159 through N01-HC-95169 and RR-024156).

### Lung GO:

**Cystic Fibrosis (CF):** Cystic Fibrosis Foundation (GIBSON07K0, KNOWLE00A0, OBSERV04K0, RDP R026), the NHLBI (R01 HL-068890, R02 HL-095396), NIH National Center for Research Resources (UL1 RR-025014), and the National Human Genome Research Institute (NHGRI) (5R00 HG-004316). **Chronic Obstructive Pulmonary Disease (COPDGene):** NHLBI (U01 HL-089897, U01 HL-089856), and the COPD Foundation through contributions made to an Industry Advisory Board comprised of AstraZeneca, Boehringer Ingelheim, Novartis, Pfizer, and Sunovion. The COPDGene clinical centers and investigators are available at [www.copdgene.org](http://www.copdgene.org). **Acute Lung Injury (ALI):** NHLBI (RC2 HL-101779). **Lung Health Study (LHS):** NHLBI (RC2 HL-066583), the NHGRI (HG-004738), and the NHLBI Division of Lung Diseases (HR-46002). **Pulmonary Arterial Hypertension (PAH):** NIH (P50 HL-084946, K23 AR-52742), and the NHLBI (F32 HL-083714). **Asthma:** NHLBI (RC2 HL-101651), and the NIH (HL-077916, HL-69197, HL-76285, M01 RR-07122).

### SWISS and ISGS:

Siblings with Ischemic Stroke Study (SWISS): National Institute of Neurological Disorders and Stroke (NINDS) (R01 NS039987); Ischemic Stroke Genetics Study (ISGS): NINDS (R01 NS042733)

### WHISP:

**Women's Health Initiative (WHI):** The WHI Sequencing Project is funded by the NHLBI (HL-102924) as well as the National Institutes of Health (NIH), U.S. Department of Health and Human Services through contracts N01WH22110, 24152, 32100-2, 32105-6, 32108-9, 32111-13, 32115, 32118-32119, 32122, 42107-26, 42129-32, and 44221. The authors thank the WHI investigators and

staff for their dedication, and the study participants for making the program possible. A full listing of WHI investigators can be found at: [http://www.whiscience.org/publications/WHI\\_investigators\\_shortlist.pdf](http://www.whiscience.org/publications/WHI_investigators_shortlist.pdf).

## **NHLBI GO Exome Sequencing Project**

### **BroadGO**

Stacey B. Gabriel (Broad Institute)<sup>4, 5, 11, 16, 17</sup>, David M. Altshuler (Broad Institute, Harvard Medical School, Massachusetts General Hospital)<sup>1, 5, 7, 17</sup>, Gonçalo R. Abecasis (University of Michigan)<sup>3, 5, 9, 13, 15, 17</sup>, Hooman Allayee (University of Southern California)<sup>5</sup>, Sharon Cresci (Washington University School of Medicine)<sup>5</sup>, Mark J. Daly (Broad Institute, Massachusetts General Hospital), Paul I. W. de Bakker (Broad Institute, Harvard Medical School, University Medical Center Utrecht)<sup>3, 15</sup>, Mark A. Depristo (Broad Institute)<sup>4, 13, 15, 16</sup>, Ron Do (Broad Institute)<sup>5, 9, 13, 15</sup>, Peter Donnelly (University of Oxford)<sup>5</sup>, Deborah N. Farlow (Broad Institute)<sup>3, 4, 5, 14, 12, 16, 17</sup>, Tim Fennell (Broad Institute), Kiran Garimella (University of Oxford)<sup>4, 16</sup>, Stanley L. Hazen (Cleveland Clinic)<sup>5</sup>, Youna Hu (University of Michigan)<sup>3, 9, 15</sup>, Daniel M. Jordan (Harvard Medical School, Harvard University)<sup>13</sup>, Goo Jun (University of Michigan), Sekar Kathiresan (Broad Institute, Harvard Medical School, Massachusetts General Hospital)<sup>5, 8, 9, 14, 12, 15, 17, 20</sup>, Steven Kawut (University of Pennsylvania)<sup>10</sup>, Adam Kiezun (Massachusetts Institute of Technology)<sup>5, 13, 15</sup>, Grigoriy Kryukov (Broad Institute), Guillaume Lettre (Broad Institute, Montreal Heart Institute, Université de Montréal)<sup>1, 2, 13, 15</sup>, Bingshan Li (University of Michigan)<sup>3</sup>, Mingyao Li (University of Pennsylvania)<sup>5</sup>, Christopher H. Newton-Cheh (Broad Institute, Massachusetts General Hospital, Harvard Medical School)<sup>3, 8, 15</sup>, Sandosh Padmanabhan (University of Glasgow School of Medicine)<sup>3, 12, 15</sup>, Sara Pulit (Broad Institute)<sup>3, 15</sup>, Daniel J. Rader (University of Pennsylvania)<sup>5</sup>, David Reich (Broad Institute, Harvard Medical School)<sup>15</sup>, Muredach P. Reilly (University of Pennsylvania)<sup>5</sup>, Manuel A. Rivas (Broad Institute, Massachusetts General Hospital)<sup>5</sup>, Steve Schwartz (Fred Hutchinson Cancer Research Center)<sup>5, 12</sup>, Laura Scott (University of Michigan)<sup>1</sup>, Johnathon A. Spertus (University of Missouri Kansas City)<sup>5</sup>, Nathaniel O. Stitzel (Brigham and Women's Hospital)<sup>5, 15</sup>, Nina Stoletzki (Brigham and Women's Hospital, Broad Institute, Harvard Medical School)<sup>13</sup>, Shamil R. Sunyaev (Brigham and Women's Hospital, Broad Institute, Harvard Medical School)<sup>1, 3, 5, 13, 15</sup>, Benjamin F. Voight (Broad Institute, Massachusetts General Hospital), Cristen J. Willer (University of Michigan)<sup>1, 9, 13, 15</sup>

### **HeartGO**

Stephen S. Rich (University of Virginia)<sup>2, 4, 7, 8, 9, 14, 11, 15, 17, 18, 31</sup>, Ermeg Akylbekova (Jackson State University, University of Mississippi Medical Center)<sup>29</sup>, Larry D. Atwood (Boston University)<sup>1, 11, 28</sup>, Christie M. Ballantyne (Baylor College of Medicine, Methodist DeBakey Heart Center)<sup>9, 22</sup>, Maja Barbalic (University of Texas Health Science Center Houston)<sup>9, 14, 15, 17, 22</sup>, R. Graham Barr (Columbia University Medical Center)<sup>10, 31</sup>, Emelia J. Benjamin (Boston University)<sup>14, 20, 28</sup>, Joshua Bis (University of Washington)<sup>15, 23</sup>, Chris Bizon (University of North Carolina Chapel Hill)<sup>3, 9, 13, 15, 23, 34</sup>, Eric Boerwinkle (University of Texas Health Science Center Houston)<sup>3, 5, 9, 13, 15, 17, 22</sup>, Donald W. Bowden (Wake Forest University)<sup>1, 31</sup>, Jennifer Brody (University of Washington)<sup>3, 5, 15, 23</sup>, Matthew Budoff (Harbor-UCLA Medical Center)<sup>31</sup>, Greg Burke (Wake Forest

University)<sup>5, 31</sup>, Sarah Buxbaum (Jackson State University)<sup>3, 13, 15, 29</sup>, Jeff Carr (Wake Forest University)<sup>25, 29, 31</sup>, Donna T. Chen (University of Virginia)<sup>6, 11</sup>, Ida Y. Chen (Cedars-Sinai Medical Center)<sup>1, 31</sup>, Wei-Min Chen (University of Virginia)<sup>13, 15, 18</sup>, Pat Concannon (University of Virginia)<sup>11</sup>, Jacy Crosby (University of Texas Health Science Center Houston)<sup>22</sup>, L. Adrienne Cupples (Boston University)<sup>1, 3, 5, 9, 13, 15, 18, 28</sup>, Ralph D'Agostino (Boston University)<sup>28</sup>, Anita L. DeStefano (Boston University)<sup>13, 18, 28</sup>, Albert Dreisbach (University of Mississippi Medical Center)<sup>3, 29</sup>, Josée Dupuis (Boston University)<sup>1, 28</sup>, J. Peter Durda (University of Vermont)<sup>15, 23</sup>, Jaclyn Ellis (University of North Carolina Chapel Hill)<sup>1</sup>, Aaron R. Folsom (University of Minnesota)<sup>5, 22</sup>, Myriam Fornage (University of Texas Health Science Center Houston)<sup>3, 18, 25</sup>, Caroline S. Fox (National Heart, Lung, and Blood Institute)<sup>1, 28</sup>, Ervin Fox (University of Mississippi Medical Center)<sup>3, 9, 29</sup>, Vincent Funari (Cedars-Sinai Medical Center)<sup>1, 11, 31</sup>, Santhi K. Ganesh (Johns Hopkins University, National Human Genome Research Institute, National Heart, Lung, and Blood Institute, University of Michigan)<sup>2, 22, 28</sup>, Julius Gardin (Hackensack University Medical Center)<sup>25</sup>, David Goff (Wake Forest University)<sup>25</sup>, Ora Gordon (Cedars-Sinai Medical Center)<sup>11, 31</sup>, Wayne Grody (University of California Los Angeles)<sup>11, 31</sup>, Myron Gross (University of Minnesota)<sup>1, 5, 14, 25</sup>, Xiuqing Guo (Cedars-Sinai Medical Center)<sup>3, 15, 31</sup>, Ira M. Hall (University of Virginia), Nancy L. Heard-Costa (Boston University)<sup>1, 11, 28</sup>, Susan R. Heckbert (University of Washington)<sup>10, 14, 15, 23</sup>, Nicholas Heintz (University of Vermont), David M. Herrington (Wake Forest University)<sup>5, 31</sup>, DeMarc Hickson (Jackson State University, University of Mississippi Medical Center)<sup>29</sup>, Jie Huang (National Heart, Lung, and Blood Institute)<sup>5, 28</sup>, Shih-Jen Hwang (Boston University, National Heart, Lung, and Blood Institute)<sup>3, 28</sup>, David R. Jacobs (University of Minnesota)<sup>25</sup>, Nancy S. Jenny (University of Vermont)<sup>1, 2, 23</sup>, Andrew D. Johnson (National Heart, Lung, and Blood Institute)<sup>2, 5, 11, 28</sup>, Craig W. Johnson (University of Washington)<sup>15, 31</sup>, Richard Kronmal (University of Washington)<sup>31</sup>, Raluca Kurz (Cedars-Sinai Medical Center)<sup>11, 31</sup>, Ethan M. Lange (University of North Carolina Chapel Hill)<sup>3, 5, 9, 13, 34</sup>, Leslie A. Lange (University of North Carolina Chapel Hill)<sup>1, 2, 3, 5, 9, 12, 13, 15, 17, 18, 20, 25, 34</sup>, Martin G. Larson (Boston University)<sup>3, 15, 28</sup>, Mark Lawson (University of Virginia), Daniel Levy (National Heart, Lung, and Blood Institute)<sup>3, 15, 17, 28</sup>, Dalin Li (Cedars-Sinai Medical Center)<sup>11, 15, 31</sup>, Honghuang Lin (Boston University)<sup>20, 28</sup>, Chunyu Liu (National Heart, Lung, and Blood Institute)<sup>3, 28</sup>, Jiankang Liu (University of Mississippi Medical Center)<sup>1, 29</sup>, Kiang Liu (Northwestern University)<sup>25</sup>, Xiaoming Liu (University of Texas Health Science Center Houston)<sup>15, 22</sup>, Yongmei Liu (Wake Forest University)<sup>2, 5, 31</sup>, William T. Longstreth (University of Washington)<sup>18, 23</sup>, Cay Loria (National Heart, Lung, and Blood Institute)<sup>25</sup>, Thomas Lumley (University of Auckland)<sup>9, 23</sup>, Kathryn Lunetta (Boston University)<sup>28</sup>, Aaron J. Mackey (University of Virginia)<sup>16, 18</sup>, Rachel Mackey (University of Pittsburgh)<sup>1, 23, 31</sup>, Ani Manichaikul (University of Virginia)<sup>8, 15, 18</sup>, Taylor Maxwell (University of Texas Health Science Center Houston)<sup>22</sup>, Barbara McKnight (University of Washington)<sup>15, 23</sup>, James B. Meigs (Brigham and Women's Hospital, Harvard Medical School, Massachusetts General Hospital)<sup>1, 28</sup>, Alanna C. Morrison (University of Texas Health Science Center Houston)<sup>3, 15, 17</sup>, Solomon K. Musani (University of Mississippi Medical Center)<sup>3, 29</sup>, Josyf C. Mychaleckyj (University of Virginia)<sup>13, 15, 31</sup>, Jennifer A. Nettleton (University of Texas Health Science Center Houston)<sup>9, 22</sup>, Kari North (University of North Carolina Chapel Hill)<sup>1, 3, 9, 10, 13, 15, 17, 34</sup>, Christopher J. O'Donnell (Massachusetts General Hospital, National Heart, Lung, and

Blood Institute)<sup>2, 5, 9, 14, 11, 12, 15, 17, 20, 28</sup>, Daniel O'Leary (Tufts University School of Medicine)<sup>25, 31</sup>, Frank Ong (Cedars-Sinai Medical Center)<sup>3, 11, 31</sup>, Walter Palmas (Columbia University)<sup>3, 15, 31</sup>, James S. Pankow (University of Minnesota)<sup>1, 22</sup>, Nathan D. Pankratz (Indiana University School of Medicine)<sup>15, 25</sup>, Shom Paul (University of Virginia), Marco Perez (Stanford University School of Medicine), Sharina D. Person (University of Alabama at Birmingham, University of Alabama at Tuscaloosa)<sup>25</sup>, Joseph Polak (Tufts University School of Medicine)<sup>31</sup>, Wendy S. Post (Johns Hopkins University)<sup>3, 9, 14, 11, 20, 31</sup>, Bruce M. Psaty (Group Health Research Institute, University of Washington)<sup>3, 5, 9, 14, 11, 15, 23</sup>, Aaron R. Quinlan (University of Virginia)<sup>18, 19</sup>, Leslie J. Raffel (Cedars-Sinai Medical Center)<sup>6, 11, 31</sup>, Vasani S. Ramachandran (Boston University)<sup>3, 28</sup>, Alexander P. Reiner (Fred Hutchinson Cancer Research Center, University of Washington)<sup>1, 2, 3, 5, 9, 11, 12, 13, 14, 15, 20, 25, 34</sup>, Kenneth Rice (University of Washington)<sup>15, 23</sup>, Jerome I. Rotter (Cedars-Sinai Medical Center)<sup>1, 3, 6, 8, 11, 15, 31</sup>, Jill P. Sanders (University of Vermont)<sup>23</sup>, Pamela Schreiner (University of Minnesota)<sup>25</sup>, Sudha Seshadri (Boston University)<sup>18, 28</sup>, Steve Shea (Brigham and Women's Hospital, Harvard University)<sup>28</sup>, Stephen Sidney (Kaiser Permanente Division of Research, Oakland, CA)<sup>25</sup>, Kevin Silverstein (University of Minnesota)<sup>25</sup>, David S. Siscovick (University of Washington)<sup>5, 1, 25</sup>, Nicholas L. Smith (University of Washington)<sup>2, 15, 20, 23</sup>, Nona Sotoodehnia (University of Washington)<sup>3, 15, 23</sup>, Asoke Srinivasan (Tougaloo College)<sup>29</sup>, Herman A. Taylor (Jackson State University, Tougaloo College, University of Mississippi Medical Center)<sup>5, 29</sup>, Kent Taylor (Cedars-Sinai Medical Center)<sup>31</sup>, Fridtjof Thomas (University of Texas Health Science Center Houston)<sup>3, 22</sup>, Russell P. Tracy (University of Vermont)<sup>5, 9, 14, 11, 12, 15, 17, 20, 23</sup>, Michael Y. Tsai (University of Minnesota)<sup>9, 31</sup>, Kelly A. Volcik (University of Texas Health Science Center Houston)<sup>22</sup>, Christina L. Wassel (University of California San Diego)<sup>9, 15, 31</sup>, Karol Watson (University of California Los Angeles), Gina Wei (National Heart, Lung, and Blood Institute)<sup>25</sup>, Wendy White (Tougaloo College)<sup>29</sup>, Kerri L. Wiggins (University of Vermont)<sup>23</sup>, Jemma B. Wilk (Boston University)<sup>10, 28</sup>, O. Dale Williams (Florida International University)<sup>25</sup>, James G. Wilson (University of Mississippi Medical Center)<sup>1, 2, 5, 8, 9, 14, 11, 12, 17, 20, 29</sup>, Phillip Wolf (Boston University)<sup>28</sup>, Neil A. Zakai (University of Vermont)<sup>2, 23</sup>

## **ISGS and SWISS**

John Hardy (Reta Lila Weston Research Laboratories, Institute of Neurology, University College London)<sup>18</sup>, James F. Meschia (Mayo Clinic)<sup>18</sup>, Michael Nalls (National Institute on Aging)<sup>2, 18</sup>, Stephen S. Rich (University of Virginia)<sup>2, 4, 7, 8, 9, 14, 11, 15, 17, 18, 31</sup>, Andrew Singleton (National Institute on Aging)<sup>18</sup>, Brad Worrall (University of Virginia)<sup>18</sup>

## **LungGO**

Michael J. Bamshad (Seattle Children's Hospital, University of Washington)<sup>4, 6, 7, 8, 10, 11, 13, 15, 17, 27</sup>, Kathleen C. Barnes (Johns Hopkins University)<sup>2, 10, 14, 12, 15, 17, 20, 24, 30, 32</sup>, Ibrahim Abdulhamid (Children's Hospital of Michigan)<sup>27</sup>, Frank Accurso (University of Colorado)<sup>27</sup>, Ran Anbar (Upstate Medical University)<sup>27</sup>, Terri Beaty (Johns Hopkins University)<sup>24, 30</sup>, Abigail Bigham (University of Washington)<sup>13, 15</sup>, Phillip Black (Children's Mercy Hospital)<sup>27</sup>, Eugene Bleecker (Wake Forest University)<sup>33</sup>, Kati Buckingham (University of Washington)<sup>27</sup>, Anne Marie Cairns (Maine Medical Center)<sup>27</sup>, Daniel Caplan (Emory University)<sup>27</sup>, Barbara Chatfield (University of Utah)<sup>27</sup>, Aaron Chidekel



(A.I. Dupont Institute Medical Center)<sup>27</sup>, Michael Cho (Brigham and Women's Hospital, Harvard Medical School)<sup>13, 15, 24</sup>, David C. Christiani (Massachusetts General Hospital)<sup>21</sup>, James D. Crapo (National Jewish Health)<sup>24, 30</sup>, Julia Crouch (Seattle Children's Hospital)<sup>6</sup>, Denise Daley (University of British Columbia)<sup>30</sup>, Anthony Dang (University of North Carolina Chapel Hill)<sup>26</sup>, Hong Dang (University of North Carolina Chapel Hill)<sup>26</sup>, Alicia De Paula (Ochsner Health System)<sup>27</sup>, Joan DeCelie-Germana (Schneider Children's Hospital)<sup>27</sup>, Allen Dozor (New York Medical College, Westchester Medical Center)<sup>27</sup>, Mitch Drumm (University of North Carolina Chapel Hill)<sup>26</sup>, Maynard Dyson (Cook Children's Med. Center)<sup>27</sup>, Julia Emerson (Seattle Children's Hospital, University of Washington)<sup>27</sup>, Mary J. Emond (University of Washington)<sup>10, 13, 15, 17, 27</sup>, Thomas Ferkol (St. Louis Children's Hospital, Washington University School of Medicine)<sup>27</sup>, Robert Fink (Children's Medical Center of Dayton)<sup>27</sup>, Cassandra Foster (Johns Hopkins University)<sup>30</sup>, Deborah Froh (University of Virginia)<sup>27</sup>, Li Gao (Johns Hopkins University)<sup>24, 30, 32</sup>, William Gershan (Children's Hospital of Wisconsin)<sup>27</sup>, Ronald L. Gibson (Seattle Children's Hospital, University of Washington)<sup>10, 27</sup>, Elizabeth Godwin (University of North Carolina Chapel Hill)<sup>26</sup>, Magdalen Gondor (All Children's Hospital Cystic Fibrosis Center)<sup>27</sup>, Hector Gutierrez (University of Alabama at Birmingham)<sup>27</sup>, Nadia N. Hansel (Johns Hopkins University, Johns Hopkins University School of Public Health)<sup>10, 15, 30</sup>, Paul M. Hassoun (Johns Hopkins University)<sup>10, 14, 32</sup>, Peter Hiatt (Texas Children's Hospital)<sup>27</sup>, John E. Hokanson (University of Colorado)<sup>24</sup>, Michelle Howenstine (Indiana University, Riley Hospital for Children)<sup>27</sup>, Laura K. Hummer (Johns Hopkins University)<sup>32</sup>, Jamshed Kanga (University of Kentucky)<sup>27</sup>, Yoonhee Kim (National Human Genome Research Institute)<sup>24, 32</sup>, Michael R. Knowles (University of North Carolina Chapel Hill)<sup>10, 26</sup>, Michael Konstan (Rainbow Babies & Children's Hospital)<sup>27</sup>, Thomas Lahiri (Vermont Children's Hospital at Fletcher Allen Health Care)<sup>27</sup>, Nan Laird (Harvard School of Public Health)<sup>24</sup>, Christoph Lange (Harvard School of Public Health)<sup>24</sup>, Lin Lin (Harvard Medical School)<sup>21</sup>, Tin L. Louie (University of Washington)<sup>13, 15, 27</sup>, David Lynch (National Jewish Health)<sup>24</sup>, Barry Make (National Jewish Health)<sup>24</sup>, Thomas R. Martin (University of Washington, VA Puget Sound Medical Center)<sup>10, 21</sup>, Steve C. Mathai (Johns Hopkins University)<sup>32</sup>, Rasika A. Mathias (Johns Hopkins University)<sup>10, 13, 15, 30, 32</sup>, John McNamara (Children's Hospitals and Clinics of Minnesota)<sup>27</sup>, Sharon McNamara (Seattle Children's Hospital)<sup>27</sup>, Deborah Meyers (Wake Forest University)<sup>33</sup>, Susan Millard (DeVos Children's Butterworth Hospital, Spectrum Health Systems)<sup>27</sup>, Peter Mogayzel (Johns Hopkins University)<sup>27</sup>, Richard Moss (Stanford University)<sup>27</sup>, Tanda Murray (Johns Hopkins University)<sup>30</sup>, Dennis Nielson (University of California at San Francisco)<sup>27</sup>, Blakeslee Noyes (Cardinal Glennon Children's Hospital)<sup>27</sup>, Wanda O'Neal (University of North Carolina Chapel Hill)<sup>26</sup>, David Orenstein (Children's Hospital of Pittsburgh)<sup>27</sup>, Brian O'Sullivan (University of Massachusetts Memorial Health Care)<sup>27</sup>, Rhonda Pace (University of North Carolina Chapel Hill)<sup>26</sup>, Peter Pare (St. Paul's Hospital, University of Washington)<sup>30</sup>, H. Worth Parker (Dartmouth-Hitchcock Medical Center, New Hampshire Cystic Fibrosis Center)<sup>27</sup>, Mary Ann Passero (Rhode Island Hospital)<sup>27</sup>, Elizabeth Perkett (Vanderbilt University)<sup>27</sup>, Adrienne Prestridge (Children's Memorial Hospital)<sup>27</sup>, Nicholas M. Rafaels (Johns Hopkins University)<sup>30</sup>, Bonnie Ramsey (Seattle Children's Hospital, University of Washington)<sup>27</sup>, Elizabeth Regan (National Jewish Health)<sup>24</sup>, Clement Ren (University of Rochester)<sup>27</sup>, George Retsch-Bogart (University of North Carolina Chapel Hill)<sup>27</sup>,

Michael Rock (University of Wisconsin Hospital and Clinics)<sup>27</sup>, Antony Rosen (Johns Hopkins University)<sup>32</sup>, Margaret Rosenfeld (Seattle Children's Hospital, University of Washington)<sup>27</sup>, Ingo Ruczinski (Johns Hopkins University School of Public Health)<sup>13, 15, 30</sup>, Andrew Sanford (University of British Columbia)<sup>30</sup>, David Schaeffer (Nemours Children's Clinic)<sup>27</sup>, Cindy Sell (University of North Carolina Chapel Hill)<sup>26</sup>, Daniel Sheehan (Children's Hospital of Buffalo)<sup>27</sup>, Edwin K. Silverman (Brigham and Women's Hospital, Harvard Medical School)<sup>24, 30</sup>, Don Sin (Children's Medical Center of Dayton)<sup>30</sup>, Terry Spencer (Children's Hospital Boston)<sup>27</sup>, Jackie Stonebraker (University of North Carolina Chapel Hill)<sup>26</sup>, Holly K. Tabor (Seattle Children's Hospital, University of Washington)<sup>6, 10, 11, 17, 27</sup>, Laurie Varlotta (St. Christopher's Hospital for Children)<sup>27</sup>, Candelaria I. Vergara (Johns Hopkins University)<sup>30</sup>, Robert Weiss<sup>30</sup>, Fred Wigley (Johns Hopkins University)<sup>32</sup>, Robert A. Wise (Johns Hopkins University)<sup>30</sup>, Fred A. Wright (University of North Carolina Chapel Hill)<sup>26</sup>, Mark M. Wurfel (Harvard School of Public Health, Massachusetts General Hospital, University of Washington)<sup>10, 14, 21</sup>, Robert Zanni (Monmouth Medical Center)<sup>27</sup>, Fei Zou (University of North Carolina Chapel Hill)<sup>26</sup>

## SeattleGO

Deborah A. Nickerson (University of Washington)<sup>3, 4, 5, 7, 8, 9, 11, 15, 17, 18, 19</sup>, Mark J. Rieder (University of Washington)<sup>4, 11, 13, 15, 16, 17, 19</sup>, Phil Green (University of Washington), Jay Shendure (University of Washington)<sup>1, 8, 14, 16, 17</sup>, Bryan Paeper (University of Washington), Joshua M. Akey (University of Washington)<sup>14, 13, 15</sup>, Michael J. Bamshad (Seattle Children's Hospital, University of Washington)<sup>4, 6, 7, 8, 10, 11, 13, 15, 17, 27</sup>, Carlos D. Bustamante (Stanford University School of Medicine)<sup>3, 13, 15</sup>, David R. Crosslin (University of Washington)<sup>2, 9</sup>, Evan E. Eichler (University of Washington)<sup>19</sup>, P. Keolu Fox<sup>2</sup>, Adam Gordon (University of Washington)<sup>11</sup>, Simon Gravel (Stanford University School of Medicine)<sup>13, 15</sup>, Gail P. Jarvik (University of Washington)<sup>9, 15</sup>, Jill M. Johnsen (Puget Sound Blood Center, University of Washington)<sup>2</sup>, Eimear E. Kenny (Stanford University School of Medicine)<sup>3, 13, 15</sup>, Jeffrey M. Kidd (Stanford University School of Medicine)<sup>13, 15</sup>, Fremiet Lara-Garduno (Baylor College of Medicine)<sup>15</sup>, Suzanne M. Leal (Baylor College of Medicine)<sup>1, 13, 15, 16, 17, 19, 20</sup>, Dajiang J. Liu (Baylor College of Medicine)<sup>13, 15</sup>, Sean McGee (University of Washington)<sup>13, 15, 19</sup>, Peggy D. Robertson (University of Washington)<sup>4</sup>, Joshua D. Smith (University of Washington)<sup>4, 16, 19</sup>, Jeffrey C. Staples (University of Washington), Emily H. Turner (University of Washington), Gao Wang (Baylor College of Medicine)

## WHISP

Rebecca Jackson (Ohio State University)<sup>1, 2, 4, 5, 8, 12, 14, 15, 17, 18, 20, 34</sup>, Kari North (University of North Carolina Chapel Hill)<sup>1, 3, 9, 10, 13, 15, 17, 34</sup>, Ulrike Peters (Fred Hutchinson Cancer Research Center)<sup>1, 3, 11, 12, 13, 15, 17, 18, 34</sup>, Christopher S. Carlson (Fred Hutchinson Cancer Research Center, University of Washington)<sup>1, 2, 3, 5, 14, 12, 13, 15, 16, 17, 18, 19, 34</sup>, Garnet Anderson (Fred Hutchinson Cancer Research Center)<sup>34</sup>, Hoda Anton-Culver (University of California at Irvine)<sup>34</sup>, Themistocles L. Assimes (Stanford University School of Medicine)<sup>5, 9, 11, 34</sup>, Paul L. Auer (Fred Hutchinson Cancer Research Center)<sup>1, 2, 3, 5, 11, 12, 13, 15, 16, 18, 34</sup>, Shirley Beresford (Fred Hutchinson Cancer Research Center)<sup>34</sup>, Chris Bizon (University of North Carolina Chapel Hill)<sup>3, 9, 13, 15, 23, 34</sup>, Henry

Black (Rush Medical Center)<sup>34</sup>, Robert Brunner (University of Nevada)<sup>34</sup>, Robert Brzyski (University of Texas Health Science Center San Antonio)<sup>34</sup>, Dale Burwen (National Heart, Lung, and Blood Institute WHI Project Office)<sup>34</sup>, Bette Caan (Kaiser Permanente Division of Research, Oakland, CA)<sup>34</sup>, Cara L. Carty (Fred Hutchinson Cancer Research Center)<sup>18, 34</sup>, Rowan Chlebowski (Los Angeles Biomedical Research Institute)<sup>34</sup>, Steven Cummings (University of California at San Francisco)<sup>34</sup>, J. David Curb (University of Hawaii)<sup>9, 18, 34</sup>, Charles B. Eaton (Brown University, Memorial Hospital of Rhode Island)<sup>12, 34</sup>, Leslie Ford (National Heart, Lung, and Blood Institute WHI Project Office)<sup>34</sup>, Nora Franceschini (University of North Carolina Chapel Hill)<sup>2, 3, 9, 10, 15, 34</sup>, Stephanie M. Fullerton (University of Washington)<sup>6, 11, 34</sup>, Margery Gass (University of Cincinnati)<sup>34</sup>, Nancy Geller (National Heart, Lung, and Blood Institute WHI Project Office)<sup>34</sup>, Gerardo Heiss (University of North Carolina Chapel Hill)<sup>5, 34</sup>, Barbara V. Howard (Howard University, MedStar Research Institute)<sup>34</sup>, Li Hsu (Fred Hutchinson Cancer Research Center)<sup>1, 13, 15, 18, 34</sup>, Carolyn M. Hutter (Fred Hutchinson Cancer Research Center)<sup>13, 15, 18, 34</sup>, John Ioannidis (Stanford University School of Medicine)<sup>11, 34</sup>, Shuo Jiao (Fred Hutchinson Cancer Research Center)<sup>34</sup>, Karen C. Johnson (University of Tennessee Health Science Center)<sup>3, 34</sup>, Emond Kabagambe (University of Alabama at Birmingham)<sup>34</sup>, Charles Kooperberg (Fred Hutchinson Cancer Research Center)<sup>1, 5, 9, 14, 13, 15, 17, 18, 34</sup>, Lewis Kuller (University of Pittsburgh)<sup>34</sup>, Andrea LaCroix (Fred Hutchinson Cancer Research Center)<sup>34</sup>, Kamakshi Lakshminarayan (University of Minnesota)<sup>18, 34</sup>, Dorothy Lane (State University of New York at Stony Brook)<sup>34</sup>, Ethan M. Lange (University of North Carolina Chapel Hill)<sup>3, 5, 9, 13, 34</sup>, Leslie A. Lange (University of North Carolina Chapel Hill)<sup>1, 2, 3, 5, 9, 12, 13, 15, 17, 18, 20, 25, 34</sup>, Norman Lasser (University of Medicine and Dentistry of New Jersey)<sup>34</sup>, Erin LeBlanc (Kaiser Permanente Center for Health Research, Portland, OR)<sup>34</sup>, Cora E. Lewis (University of Alabama at Birmingham)<sup>34</sup>, Marian Limacher (University of Florida)<sup>34</sup>, Danyu Lin (University of North Carolina Chapel Hill)<sup>1, 3, 9, 13, 15, 34</sup>, Benjamin A. Logsdon (Fred Hutchinson Cancer Research Center)<sup>2, 34</sup>, Shari Ludlam (National Heart, Lung, and Blood Institute WHI Project Office)<sup>34</sup>, JoAnn E. Manson (Brigham and Women's Hospital, Harvard School of Public Health)<sup>34</sup>, Karen Margolis (University of Minnesota)<sup>34</sup>, Lisa Martin (George Washington University Medical Center)<sup>9, 34</sup>, Joan McGowan (National Heart, Lung, and Blood Institute WHI Project Office)<sup>34</sup>, Keri L. Monda (Amgen, Inc.)<sup>1, 15, 34</sup>, Jane Morley Kotchen (Medical College of Wisconsin)<sup>34</sup>, Lauren Nathan (University of California Los Angeles)<sup>34</sup>, Judith Ockene (Fallon Clinic, University of Massachusetts)<sup>34</sup>, Mary Jo O'Sullivan (University of Miami)<sup>34</sup>, Lawrence S. Phillips (Emory University)<sup>34</sup>, Ross L. Prentice (Fred Hutchinson Cancer Research Center)<sup>34</sup>, Alexander P. Reiner (Fred Hutchinson Cancer Research Center, University of Washington)<sup>1, 2, 3, 5, 9, 11, 12, 13, 14, 15, 20, 25, 34</sup>, John Robbins (University of California at Davis)<sup>34</sup>, Jennifer G. Robinson (University of Iowa)<sup>9, 11, 18, 34</sup>, Jacques E. Rossouw (National Heart, Lung, and Blood Institute, National Heart, Lung, and Blood Institute WHI Project Office)<sup>5, 14, 17, 20, 34</sup>, Haleh Sangi-Haghepeykar (Baylor College of Medicine)<sup>34</sup>, Gloria E. Sarto (University of Wisconsin)<sup>34</sup>, Sally Shumaker (Wake Forest University)<sup>34</sup>, Michael S. Simon (Wayne State University)<sup>34</sup>, Marcia L. Stefanick (Stanford University School of Medicine)<sup>34</sup>, Evan Stein (Medical Research Labs)<sup>34</sup>, Hua Tang (Stanford University)<sup>2, 34</sup>, Kira C. Taylor (University of Louisville)<sup>1, 3, 13, 15, 20, 34</sup>, Cynthia A. Thomson (University of Arizona)<sup>34</sup>, Timothy A. Thornton (University of Washington)<sup>13, 15, 18, 34</sup>,

Linda Van Horn (Northwestern University)<sup>34</sup>, Mara Vitolins (Wake Forest University)<sup>34</sup>, Jean Wactawski-Wende (University of Buffalo)<sup>34</sup>, Robert Wallace (University of Iowa)<sup>2, 34</sup>, Sylvia Wassertheil-Smoller (Boston University)<sup>18, 34</sup>

### **NHLBI**

Deborah Applebaum-Bowden (National Heart, Lung, and Blood Institute)<sup>4, 7, 12, 17</sup>, Michael Feolo (National Center for Biotechnology Information)<sup>12</sup>, Weiniu Gan (National Heart, Lung, and Blood Institute)<sup>7, 8, 16, 17</sup>, W. Keith Hoots (National Heart, Lung, and Blood Institute)<sup>17</sup>, James Kiley (National Heart, Lung, and Blood Institute)<sup>17</sup>, Michael Lauer (National Heart, Lung, and Blood Institute)<sup>17</sup>, Hilary Leeds (National Heart, Lung, and Blood Institute), Alan Michelson (National Heart, Lung, and Blood Institute)<sup>17</sup>, Dina N. Paltoo (National Heart, Lung, and Blood Institute)<sup>4, 6, 11, 17</sup>, Phyliss Sholinsky (National Heart, Lung, and Blood Institute)<sup>4, 17</sup>, Sonia Skarlatos (National Heart, Lung, and Blood Institute)<sup>17</sup>, Anne Sturcke (National Center for Biotechnology Information)<sup>12</sup>

### **ESP Groups**

<sup>1</sup>Anthropometry Project Team, <sup>2</sup>Blood Count/Hematology Project Team, <sup>3</sup>Blood Pressure Project Team, <sup>4</sup>Data Flow Working Group, <sup>5</sup>Early MI Project Team, <sup>6</sup>ELSI Working Group, <sup>7</sup>Executive Committee, <sup>8</sup>Family Study Project Team, <sup>9</sup>Lipids Project Team, <sup>10</sup>Lung Project Team, <sup>11</sup>Personal Genomics Project Team, <sup>12</sup>Phenotype and Harmonization Working Group, <sup>13</sup>Population Genetics and Statistical Analysis Working Group, <sup>14</sup>Publications and Presentations Working Group, <sup>15</sup>Quantitative Analysis Ad Hoc Task Group, <sup>16</sup>Sequencing and Genotyping Working Group, <sup>17</sup>Steering Committee, <sup>18</sup>Stroke Project Team, <sup>19</sup>Structural Variation Working Group, <sup>20</sup>Subclinical/Quantitative Project Team

### **ESP Cohorts**

<sup>21</sup>Acute Lung Injury (ALI), <sup>22</sup>Atherosclerosis Risk in Communities (ARIC), <sup>23</sup>Cardiovascular Health Study (CHS), <sup>24</sup>Chronic Obstructive Pulmonary Disease (COPD)Gene), <sup>25</sup>Coronary Artery Risk Development in Young Adults (CARDIA), <sup>26</sup>Cystic Fibrosis (CF), <sup>27</sup>Early Pseudomonas Infection Control (EPIC), <sup>28</sup>Framingham Heart Study (FHS), <sup>29</sup>Jackson Heart Study (JHS), <sup>30</sup>Lung Health Study (LHS), <sup>31</sup>Multi-Ethnic Study of Atherosclerosis (MESA), <sup>32</sup>Pulmonary Arterial Hypertension (PAH), <sup>33</sup>Severe Asthma Research Program (SARP), <sup>34</sup>Women's Health Initiative (WHI)

## Supplementary tables

**Table S1.** Comparison of the number of SNVs with different minor allele frequency (MAF) in individuals of European ancestry between ESP and 1KGP.

MAF	Exome Project	1000 Genomes Project
Singletons	161,906 (62.0%)	28,696 (31.5%)
≤0.001	186,602 (71.4%)	Not Available
≤0.005	220,086 (84.2%)	40,677 (44.7%)
≤0.01	227,364 (87.0%)	49,540 (54.5%)
≤0.05	239,824 (91.8%)	64,191 (70.6%)
>0.05	21,427 (8.2%)	26,770 (29.4%)
Total	261,251	90,961

**Table S2.** The distribution of MAF for the 151 common SNPs in different European ancestral populations of 1KGP, but invariant in ESP.

MAF (%)	CEU	FIN	GBR	IBS	TSI
Invariant	105 (69.5%)	0	0	16 (10.6%)	107 (70.9%)
≤0.01	3 (2.0%)	0	0	0	4 (2.6%)
≤0.05	17 (11.2%)	6 (4.0%)	5 (3.3%)	12 (7.9%)	23 (15.2%)
>0.05	29 (19.2%)	145 (96.0%)	146 (96.7%)	123 (81.4%)	21 (13.9%)

**Table S3.** Genes with unusually high proportions of rare variation.

Gene	Chrom	Pi	Tajima's D	S	CommonS	RareS	Length
<i>SAMD11</i>	1	1.08E-04	-2.42	49	7	42	743
<i>KLHL17</i>	1	3.84E-04	-2.35	67	6	61	1322
<i>PLEKHN1</i>	1	6.54E-04	-2.04	53	8	45	966
<i>SCNN1D</i>	1	2.31E-04	-2.45	67	7	60	1003
<i>TAS1R3</i>	1	8.04E-05	-2.51	62	3	59	1154
<i>ATAD3B</i>	1	1.02E-03	-2.09	83	15	68	1434
<i>PRAMEF2</i>	1	3.19E-03	-1.60	96	31	65	1324
<i>TMEM82</i>	1	6.01E-04	-2.22	47	7	40	670
<i>CROCC</i>	1	1.10E-04	-2.52	100	14	86	2111
<i>CELA3B</i>	1	1.52E-03	-2.01	52	12	40	788
<i>GJB4</i>	1	7.49E-04	-2.25	48	11	37	765
<i>HRNR</i>	1	9.62E-04	-2.17	311	53	258	4503
<i>FLG</i>	1	1.40E-03	-2.29	761	148	613	11166
<i>OR2T12</i>	1	1.44E-03	-1.94	50	11	39	820
<i>PRR21</i>	2	1.29E-03	-1.75	43	14	29	514

<i>ANO7</i>	2	6.00E-04	-2.27	98	12	86	2030
<i>TMEM175</i>	4	6.84E-04	-2.29	72	8	64	1424
<i>KIAA1530</i>	4	7.76E-04	-2.20	85	8	77	1791
<i>CRIPAK</i>	4	2.73E-03	-2.17	150	30	120	1299
<i>DOK7</i>	4	1.01E-04	-2.24	27	2	25	331
<i>WFS1</i>	4	1.51E-03	-2.21	175	25	150	2260
<i>SH3TC1</i>	4	4.74E-04	-2.42	145	21	124	2822
<i>GPR78</i>	4	6.49E-04	-2.00	35	5	30	570
<i>CPZ</i>	4	6.91E-04	-2.37	103	13	90	1537
<i>ZDHHC11</i>	5	6.12E-04	-2.30	62	8	54	1090
<i>SRA1</i>	5	4.39E-04	-2.44	58	11	47	499
<i>HIST1H3A</i>	6	3.82E-04	-2.23	28	3	25	378
<i>HIST1H4B</i>	6	1.61E-03	-2.03	30	2	28	315
<i>HIST1H2AB</i>	6	1.77E-03	-1.79	26	2	24	397
<i>HIST1H1C</i>	6	6.90E-04	-2.19	37	5	32	389
<i>HIST1H1T</i>	6	1.78E-03	-2.01	47	5	42	628
<i>HIST1H1D</i>	6	8.81E-04	-2.14	39	2	37	604
<i>MUC21</i>	6	1.95E-03	-2.10	126	40	86	1703
<i>HLA-DQA2</i>	6	2.85E-03	-1.70	56	15	41	781
<i>MICALL2</i>	7	8.04E-04	-2.04	83	19	64	1408
<i>AMZ1</i>	7	1.31E-04	-2.56	74	5	69	905
<i>SDK1</i>	7	7.71E-04	-2.33	275	41	234	6141
<i>SLC29A4</i>	7	1.49E-04	-2.47	58	2	56	1213
<i>CD36</i>	7	1.06E-04	-2.55	67	4	63	1464
<i>MUC17</i>	7	9.12E-04	-2.43	753	126	627	13448
<i>ERICH1</i>	8	6.94E-04	-2.28	66	8	58	1184
<i>MYOM2</i>	8	1.11E-03	-2.33	260	47	213	4529
<i>AMAC1L2</i>	8	2.03E-03	-1.91	66	13	53	1013
<i>CTSB</i>	8	8.78E-04	-2.19	55	7	48	1055
<i>TSNARE1</i>	8	1.04E-03	-2.05	66	8	58	1281
<i>IFNA14</i>	9	9.54E-05	-2.42	39	1	38	573
<i>FCN1</i>	9	1.36E-03	-2.00	54	5	49	1015
<i>FBXW5</i>	9	7.23E-05	-2.47	50	3	47	775
<i>SLC34A3</i>	9	3.30E-04	-2.29	52	3	49	952
<i>FAM166A</i>	9	3.14E-04	-2.39	49	11	38	784
<i>PPYR1</i>	10	3.16E-04	-2.47	66	7	59	1119
<i>MUPCDH</i>	11	8.35E-04	-2.12	90	14	76	1857
<i>OR4A5</i>	11	5.71E-04	-2.28	48	7	41	794
<i>OR4C46</i>	11	2.56E-03	-1.82	67	18	49	844
<i>OR4C16</i>	11	2.75E-03	-1.63	59	15	44	926
<i>KRT82</i>	12	5.10E-04	-2.34	68	6	62	1426
<i>FAM70B</i>	13	5.80E-04	-2.28	50	6	44	903
<i>DHRS2</i>	14	1.52E-04	-2.46	48	4	44	833
<i>DHRS4L2</i>	14	2.01E-03	-1.80	44	12	32	626
<i>OR4M2</i>	15	4.75E-04	-2.41	60	12	48	916
<i>C16ORF38</i>	16	3.12E-03	-1.05	63	19	44	961
<i>NUBP2</i>	16	1.50E-03	-1.91	45	6	39	734
<i>RPL3L</i>	16	2.80E-04	-2.46	62	6	56	1134
<i>PKD1</i>	16	1.32E-04	-2.45	146	15	131	2723

<i>SRRM2</i>	16	3.27E-04	-2.58	313	22	291	7866
<i>DNASE1</i>	16	9.90E-04	-2.28	63	7	56	878
<i>TRAP1</i>	16	6.01E-04	-2.39	103	13	90	2038
<i>PPL</i>	16	7.45E-04	-2.29	198	24	174	4644
<i>FAM86A</i>	16	6.26E-04	-2.10	38	7	31	571
<i>TEKT5</i>	16	9.62E-04	-2.21	76	12	64	1471
<i>PDILT</i>	16	1.45E-03	-1.96	84	10	74	1799
<i>NECAB2</i>	16	1.25E-03	-2.13	70	9	61	978
<i>SLC38A8</i>	16	9.65E-04	-2.28	81	12	69	1344
<i>KIAA1609</i>	16	1.65E-03	-1.96	76	16	60	1399
<i>KLHDC4</i>	16	9.44E-04	-2.16	73	12	61	1415
<i>CDT1</i>	16	1.03E-03	-1.74	47	8	39	757
<i>KCNJ12</i>	17	2.50E-04	-2.49	66	9	57	1187
<i>ACTG1</i>	17	1.03E-03	-2.03	50	8	42	885
<i>THEG</i>	19	5.27E-04	-2.36	61	9	52	1171
<i>POLRMT</i>	19	4.30E-04	-2.18	79	13	66	1535
<i>PRSSL1</i>	19	7.21E-04	-1.78	24	3	21	309
<i>TJP3</i>	19	1.09E-03	-1.98	103	13	90	2072
<i>CAPS</i>	19	9.26E-04	-2.19	41	7	34	558
<i>PPAN-P2RY11</i>	19	3.46E-04	-2.40	83	6	77	1596
<i>P2RY11</i>	19	6.59E-04	-2.26	55	4	51	837
<i>CYP2A6</i>	19	1.13E-03	-2.24	92	22	70	1497
<i>CYP2B6</i>	19	6.48E-04	-2.40	88	6	82	1483
<i>PSG8</i>	19	2.33E-03	-1.96	96	14	82	1264
<i>PSG1</i>	19	1.94E-03	-2.01	88	18	70	1236
<i>PSG6</i>	19	1.11E-03	-2.38	108	18	90	1294
<i>PSG11</i>	19	2.28E-03	-1.76	61	12	49	1008
<i>PSG2</i>	19	2.03E-03	-1.92	67	14	53	1028
<i>PSG9</i>	19	1.39E-03	-2.14	81	20	61	1303
<i>LILRA2</i>	19	9.31E-04	-2.31	87	24	63	1415
<i>TRIB3</i>	20	1.00E-03	-2.06	50	6	44	946
<i>LAMA5</i>	20	4.05E-04	-2.38	270	38	232	5428
<i>TPTE</i>	21	3.00E-04	-2.52	87	13	74	1651
<i>TRPM2</i>	21	2.71E-04	-2.58	178	17	161	4007
<i>C21ORF29</i>	21	1.23E-04	-2.57	84	5	79	1776
<i>COL6A2</i>	21	8.75E-04	-2.20	136	18	118	2541
<i>IGLL1</i>	22	1.35E-03	-1.97	38	12	26	443
<i>IGLL3</i>	22	3.90E-03	-1.75	68	22	46	640
<i>CYP2D6</i>	22	7.50E-04	-2.18	58	6	52	810
<i>SERHL2</i>	22	5.54E-04	-2.16	39	5	34	537
<i>LOC553158</i>	22	2.10E-03	-2.12	117	18	99	1402
<i>ARHGAP8</i>	22	2.20E-03	-2.11	115	18	97	1340
<i>LMF2</i>	22	5.64E-04	-2.19	65	6	59	898
<i>PLCXD1</i>	X	5.18E-04	-2.19	58	5	53	995
<i>IL3RA</i>	X	2.10E-04	-2.38	57	3	54	1181
<i>VCX3A</i>	X	5.49E-04	-2.00	23	7	16	171

Column headers are Gene ID, chromosome, nucleotide diversity, Tajima's D, number of segregating sites (S), number of common variants (CommonS), number of rare variants (RareS), and observed gene length.

**Table S4.** Summary information of 114 genes with signatures of positive selection.

Gene	Chr	F <sub>ST</sub>	Pi (%)	AA Pi (%)	EA Pi (%)	Nsyn Pi (%)	Syn Pi (%)	Extreme Pattern
<i>C10RF63</i>	1	0.16	0.049	0.017	0.068	0.065	0.001	LowAA
<i>EIF2C1</i>	1	0.36	0.012	0.021	0.002	0	0.049	HighFst
<i>CDC14A</i>	1	0.34	0.033	0.062	0.003	0.003	0.125	LowEA
<i>BCL9</i>	1	0.22	0.016	0.007	0.021	0.016	0.016	LowAA
<i>OR6K6</i>	1	0.26	0.275	0.505	0.028	0.26	0.32	LowEA
<i>DARC</i>	1	0.25	0.065	0.024	0.086	0.085	0.005	LowAA
<i>KLHL20</i>	1	0.31	0.025	0.046	0.002	0.001	0.097	HighFst
<i>LAMC2</i>	1	0.21	0.069	0.085	0.049	0.025	0.203	MK
<i>PRG4</i>	1	0.17	0.046	0.05	0.04	0.029	0.098	LowPi
<i>LRRN2</i>	1	0.27	0.065	0.04	0.079	0.036	0.154	LowAA
<i>SLC41A1</i>	1	0.22	0.046	0.021	0.057	0.001	0.182	LowAA
<i>C10RF116</i>	1	0.36	0.079	0.147	0.004	0.099	0.018	LowEA
<i>C10RF101</i>	1	0.21	0.063	0.098	0.025	0.016	0.206	MK
<i>OR2W5</i>	1	0.17	0.164	0.196	0.131	0.045	0.519	MK
<i>OR2T10</i>	1	0.02	0.015	0.025	0.008	0.012	0.025	LowPi
<i>NT5C1B</i>	2	0.26	0.066	0.072	0.045	0.003	0.245	MK
<i>MAP4K3</i>	2	0.45	0.04	0.04	0.021	0.002	0.154	MK
<i>STK17B</i>	2	0.28	0.087	0.072	0.087	0.001	0.346	MK
<i>ARMC9</i>	2	0.08	0.077	0.078	0.074	0.006	0.289	MK
<i>DHFRL1</i>	3	0.11	0.091	0.036	0.128	0.12	0.006	LowAA
<i>UGT2B28</i>	4	0.03	0.015	0.023	0.008	0.009	0.032	LowPi
<i>UGT2A1</i>	4	0.08	0.07	0.04	0.092	0.081	0.037	LowAA
<i>BRD9</i>	5	0.20	0.072	0.133	0.013	0.012	0.252	LowEA
<i>UGT3A1</i>	5	0.04	0.008	0.016	0.001	0.003	0.021	LowPi
<i>C5ORF36</i>	5	0.36	0.03	0.052	0.001	0.04	0	HighFst
<i>C5ORF48</i>	5	0.40	0.07	0.116	0.002	0.003	0.272	HighFst
<i>LARS</i>	5	0.15	0.075	0.037	0.096	0.024	0.228	LowAA
<i>RARS</i>	5	0.21	0.074	0.029	0.102	0.016	0.248	LowAA
<i>FBXW11</i>	5	0.35	0.022	0.04	0.001	0	0.087	HighFst
<i>FAM153A</i>	5	0.08	0.074	0.15	0.01	0.099	0.001	LowPi
<i>HLA-G</i>	6	0.11	0.352	0.379	0.326	0.068	1.205	MK
<i>WDR46</i>	6	0.07	0.032	0.015	0.044	0.023	0.059	LowAA
<i>TMEM63B</i>	6	0.05	0.02	0.01	0.028	0.007	0.061	LowAA
<i>FUT9</i>	6	0.32	0.023	0.042	0.001	0.03	0.002	HighFst
<i>OR2A4</i>	6	0.00	0	0	0.001	0	0.001	LowPi
<i>C7ORF20</i>	7	0.11	0.115	0.142	0.089	0.001	0.455	MK
<i>MUC17</i>	7	0.17	0.092	0.121	0.063	0.089	0.101	LowPi
<i>TAS2R60</i>	7	0.04	0.059	0.064	0.054	0.011	0.204	MK
<i>SCARA3</i>	8	0.01	0.003	0.005	0.001	0.002	0.006	LowPi
<i>UNC5D</i>	8	0.40	0.04	0.065	0.006	0.002	0.152	LowEA
<i>PTK2</i>	8	0.42	0.012	0.02	0.002	0.001	0.046	HighFst
<i>C9ORF123</i>	9	0.03	0.17	0.16	0.174	0.002	0.672	MK
<i>CDC14B</i>	9	0.41	0.025	0.043	0.002	0.005	0.086	HighFst
<i>ZNF510</i>	9	0.21	0.067	0.131	0.007	0.058	0.095	LowEA
<i>ORM2</i>	9	0.42	0.173	0.312	0.008	0.109	0.364	LowEA
<i>OR1L3</i>	9	0.39	0.086	0.162	0.007	0.102	0.039	LowEA



<i>NMT2</i>	10	0.30	0.021	0.039	0.002	0.001	0.081	HighFst
<i>ZNF248</i>	10	0.05	0.025	0.01	0.035	0.001	0.095	LowAA
<i>CYP17A1</i>	10	0.01	0.065	0.062	0.066	0	0.258	MK
<i>PNLIPRP3</i>	10	0.21	0.081	0.14	0.023	0.07	0.113	LowEA
<i>OR51F2</i>	11	0.11	0.054	0.066	0.04	0.03	0.126	MK
<i>OR56A4</i>	11	0.23	0.111	0.092	0.116	0.043	0.316	MK
<i>KCNJ11</i>	11	0.23	0.106	0.052	0.133	0.094	0.142	LowAA
<i>PACIN3</i>	11	0.32	0.024	0.043	0.001	0.03	0.005	HighFst
<i>FAM111A</i>	11	0.05	0.034	0.042	0.026	0.013	0.099	MK
<i>VWCE</i>	11	0.21	0.021	0.043	0.002	0.005	0.07	MK
<i>GPR152</i>	11	0.47	0.139	0.227	0.005	0.047	0.414	LowEA
<i>MLL</i>	11	0.07	0.005	0.006	0.004	0.001	0.017	LowPi
<i>C11ORF63</i>	11	0.11	0.025	0.028	0.02	0.006	0.081	MK
<i>OR6T1</i>	11	0.14	0.146	0.126	0.146	0.032	0.486	MK
<i>TAS2R20</i>	12	0.16	0.291	0.154	0.36	0.383	0.016	LowAA
<i>SLCO1B1</i>	12	0.30	0.123	0.115	0.112	0.069	0.279	MK
<i>ZNF641</i>	12	0.20	0.054	0.023	0.069	0.036	0.108	LowAA
<i>AMHR2</i>	12	0.49	0.021	0.032	0.001	0.001	0.083	HighFst
<i>OR6C70</i>	12	0.16	0.118	0.089	0.131	0.107	0.149	LowPi
<i>OR6C4</i>	12	0.31	0.081	0.151	0.007	0.052	0.167	LowEA
<i>POLR3B</i>	12	0.13	0.022	0.012	0.029	0.004	0.035	LowAA
<i>HNF1A</i>	12	0.11	0.065	0.037	0.082	0.033	0.163	LowAA
<i>DIABLO</i>	12	0.32	0.035	0.062	0.002	0.002	0.133	HighFst
<i>IRS2</i>	13	0.00	0.001	0.002	0.001	0.002	0	LowPi
<i>PACS2</i>	14	0.42	0.049	0.081	0.004	0.021	0.133	LowEA
<i>PGBD4</i>	15	0.07	0.007	0.015	0.001	0.005	0.014	LowPi
<i>SLC24A5</i>	15	0.77	0.032	0.027	0.002	0.042	0.002	HighFst
<i>CD19</i>	16	0.16	0.062	0.031	0.079	0.042	0.121	LowAA
<i>KIF22</i>	16	0.28	0.034	0.065	0.002	0.009	0.109	LowEA
<i>ABCC12</i>	16	0.47	0.078	0.121	0.026	0.064	0.122	LowEA
<i>CES2</i>	16	0.13	0.007	0.013	0.002	0.001	0.024	LowPi
<i>LRRC36</i>	16	0.43	0.064	0.097	0.014	0.058	0.085	LowEA
<i>CDRT15</i>	17	0.06	0.058	0.048	0.063	0.062	0.046	LowPi
<i>TBC1D29</i>	17	0.11	0.062	0.075	0.049	0.003	0.17	LowPi
<i>WNK4</i>	17	0.24	0.029	0.054	0.003	0.018	0.061	LowEA
<i>CNTD1</i>	17	0.35	0.049	0.096	0.001	0.046	0.061	HighFst
<i>IMP5</i>	17	0.12	0.195	0.113	0.244	0.171	0.27	LowAA
<i>MAPT</i>	17	0.12	0.125	0.066	0.163	0.005	0.181	LowAA
<i>KIAA1267</i>	17	0.12	0.091	0.044	0.121	0.048	0.22	LowAA
<i>C17ORF64</i>	17	0.61	0.357	0.492	0.05	0.097	1.138	MK, LowEA
<i>APPBP2</i>	17	0.41	0.026	0.044	0.002	0.022	0.036	HighFst
<i>LRRC37A3</i>	17	0.17	0.014	0.029	0.002	0.017	0.006	LowPi
<i>BPTF</i>	17	0.36	0.011	0.02	0.002	0.004	0.029	HighFst; LowEA
<i>GPS1</i>	17	0.02	0.006	0.007	0.006	0.001	0.021	LowPi
<i>BRUNOL4</i>	18	0.06	0.038	0.03	0.044	0.001	0.151	MK
<i>FECH</i>	18	0.17	0.074	0.037	0.094	0.026	0.215	LowAA
<i>OR7G2</i>	19	0.10	0.138	0.115	0.146	0.134	0.15	LowPi
<i>ICAM5</i>	19	0.20	0.084	0.045	0.102	0.088	0.071	LowAA

<i>RGL3</i>	19	0.01	0.003	0.006	0.002	0.003	0.003	LowPi
<i>ZNF442</i>	19	0.20	0.056	0.103	0.01	0.065	0.028	LowEA
<i>USHBP1</i>	19	0.24	0.102	0.108	0.083	0.066	0.207	MK
<i>HIPK4</i>	19	0.03	0.014	0.016	0.013	0.016	0.01	LowPi
<i>LIPE</i>	19	0.23	0.03	0.058	0.005	0.033	0.016	LowEA
<i>VRK3</i>	19	0.19	0.114	0.217	0.012	0.093	0.176	LowEA
<i>ZNF473</i>	19	0.20	0.074	0.145	0.004	0.067	0.095	LowEA
<i>POLD1</i>	19	0.34	0.061	0.108	0.015	0.008	0.221	LowEA
<i>KIR2DL1</i>	19	0.04	0.028	0.049	0.011	0.034	0.011	LowPi
<i>KIR2DS4</i>	19	0.04	0.021	0.042	0.004	0	0	LowPi
<i>NLRP11</i>	19	0.19	0.06	0.068	0.05	0.026	0.163	MK
<i>DEFB127</i>	20	0.17	0.232	0.113	0.292	0.306	0.009	LowAA
<i>C20ORF70</i>	20	0.26	0.096	0.182	0.012	0.093	0.104	LowEA
<i>R3HDML</i>	20	0.04	0.118	0.106	0.128	0.005	0.458	MK
<i>TPTE</i>	21	0.07	0.031	0.045	0.018	0.028	0.015	LowPi
<i>SLC5A1</i>	22	0.03	0.033	0.016	0.046	0.019	0.076	LowAA
<i>APOL1</i>	22	0.32	0.155	0.141	0.143	0.152	0.162	LowPi
<i>TXN2</i>	22	0.01	0.006	0.007	0.005	0.001	0.02	LowPi
<i>APOBEC3G</i>	22	0.33	0.082	0.085	0.059	0.058	0.152	MK
<i>TRMU</i>	22	0.23	0.047	0.094	0.003	0.02	0.128	LowEA

Column headers are Gene ID, chromosome,  $F_{ST}$ , nucleotide diversity in all ESP samples, nucleotide diversity in the AA sample, nucleotide diversity in the EA sample, nucleotide diversity of nonsynonymous variants, nucleotide diversity of synonymous variants, and the characteristic that makes them an outlier (Extreme Patter = MK for McDonald-Kreitman; LowPi = low nucleotide diversity in the combined sample; LowEA or LowAA = low nucleotide diversity in EA or AA, respectively; HighFst = extreme levels of differentiation as assessed by  $F_{ST}$ ).

**Table S5.** Significantly enriched KEGG pathways among the 114 genes with signatures of positive selection.

KEGG Pathway	Number of Genes	Adjusted P-value
Olfactory transduction	10	$1.79 \times 10^{-7}$
Drug Metabolism -other enzymes	3	$5 \times 10^{-4}$
Antigen processing and presentation	3	$2.2 \times 10^{-3}$
Natural killer cell mediated cytotoxicity	3	$4.3 \times 10^{-3}$
Metabolic pathways	8	$4.3 \times 10^{-2}$

Column headers are KEGG pathway ID, Number of putatively selected genes in the pathway, and FDR adjusted p-value. Enrichment analysis was performed using the web-based tool WebGestalt(13) (<http://bioinfo.vanderbilt.edu/webgestalt>)

**Table S6.** Summary statistics of polymorphism for each of the 15,585 genes analyzed.

This is a large table that is available as a separate file. Column headers are:

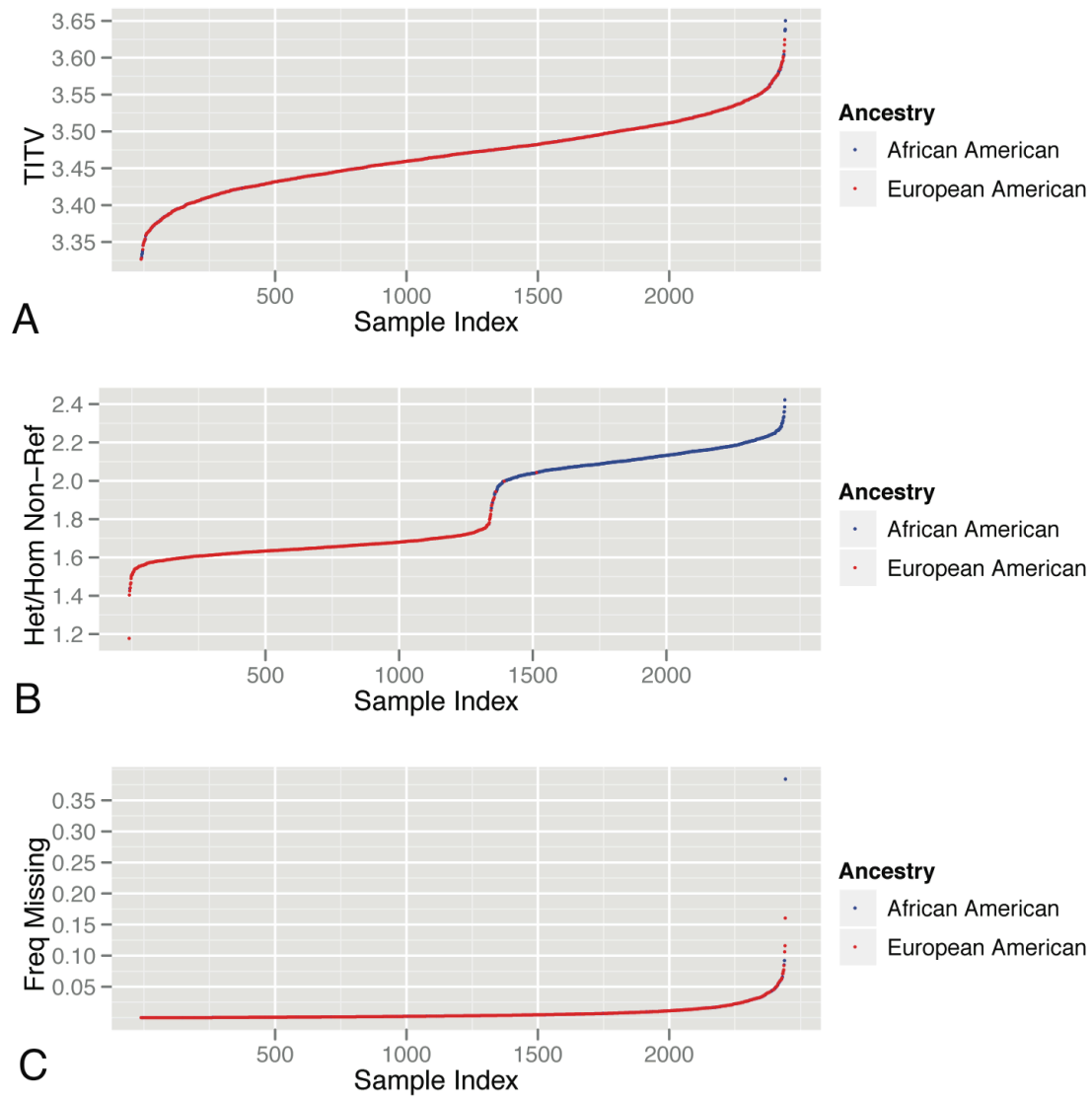
Gene (gene name)  
 Chrom (chromosome)  
 Start (start site)  
 End (end site)  
 S (number of segregating sites)  
 EA\_S (number of segregating sites in European Americans)  
 AA\_S (number of segregating sites in African Americans)  
 Length (real exonic length of the gene)  
 ObservedLength (number of sites that passed all filters, including invariant sites)  
 Pi (mean number of pairwise differences between two randomly selected chromosomes, divided by ObservedLength)  
 EA\_Pi (Pi in European Americans)  
 AA\_Pi (Pi in African Americans)  
 Rare (number of segregating sites with frequency under 0.5%)  
 Common (number of segregating sites with frequency at least 0.5%)  
 Missense (number of missense segregating sites)  
 PolyPhen2 (number of segregating sites with "probably damaging" PolyPhen2 designation)  
 HighGERP (number of segregating sites with GERP score  $\geq 5$ )  
 HighSFS (number of segregating sites with  $\geq 90\%$  probability of functional importance using SFS-based method)  
 Nons (number of nonsense segregating sites)  
 Splice (number of splice segregating sites)  
 MissenseEA (number of missense segregating sites in European Americans)  
 NonsEA (number of nonsense segregating sites in European Americans)  
 SpliceEA (number of splice segregating sites in European Americans)  
 MissenseAA (number of missense segregating sites in African Americans)  
 NonsAA (number of nonsense segregating sites in African Americans)  
 SpliceAA (number of splice segregating sites in African Americans)  
 Sprop (segregating sites divided by observed length)  
 NS\_Sprop (missense, nonsense, and splice segregating sites divided by observed length)  
 SpropEA (segregating sites in European Americans divided by observed length)  
 NS\_SpropEA (missense, nonsense, and splice segregating sites in European Americans divided by observed length)  
 SpropAA (segregating sites in African Americans divided by observed length)  
 NS\_SpropAA (missense, nonsense, and splice segregating sites in African Americans divided by observed length)

**Table S7.** Evidence for coordinated purifying selection on synonymous and nonsynonymous variation within genes.

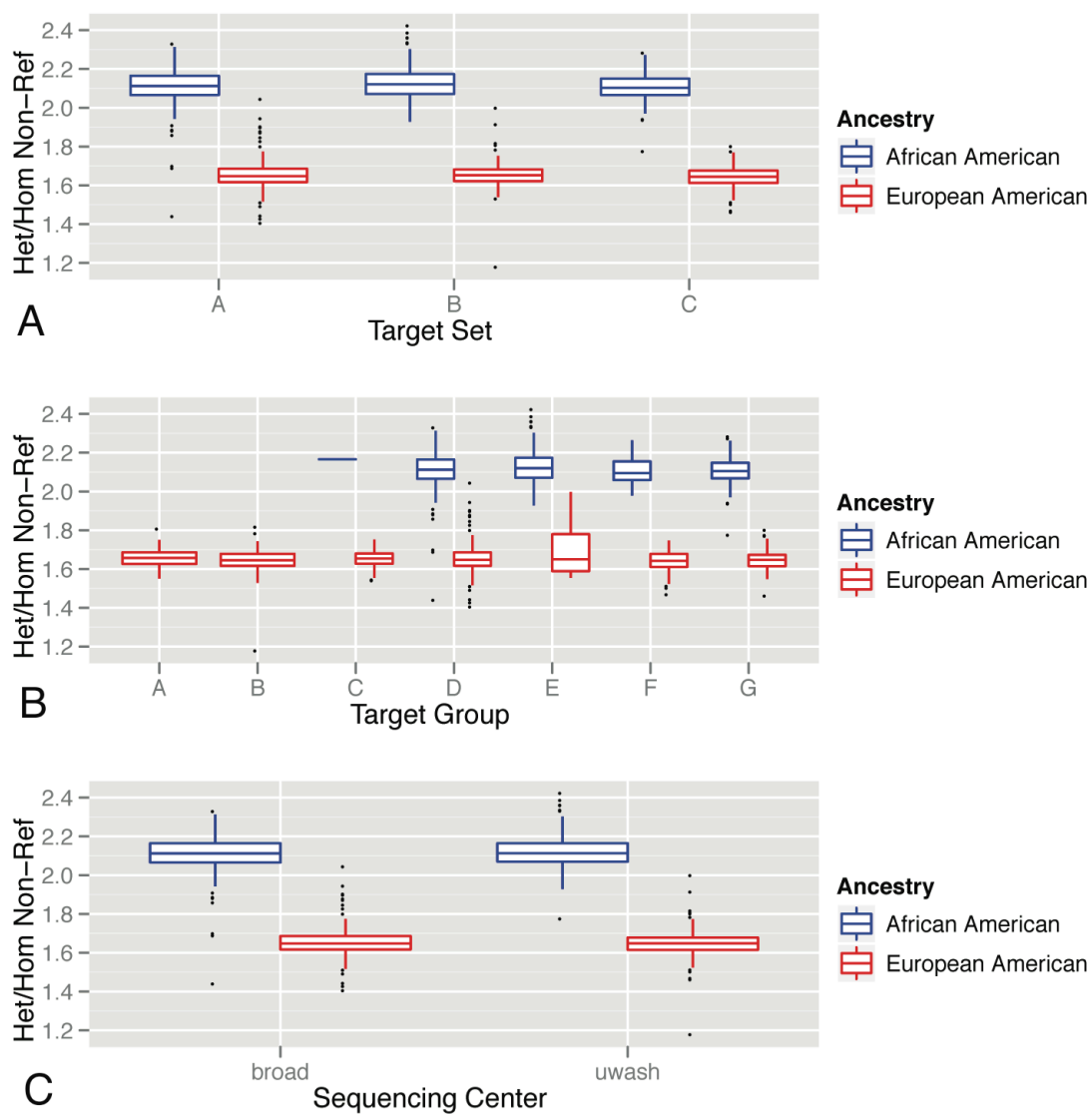
Method	AA	EA
SIFT	-0.031 (0.0002)	-0.020 (0.0194)
Polyphen2	-0.032 (0.0001)	-0.018 (0.0070)
LRT	0.010 (0.1247)	-0.001 (0.4590)
MutationTaster	-0.020 (0.0098)	-0.026 (0.0013)

Values show the rank correlation coefficients between average functional prediction scores for different methods and average w score changes of a gene. P-values from a two-sided z test are shown in parantheses.

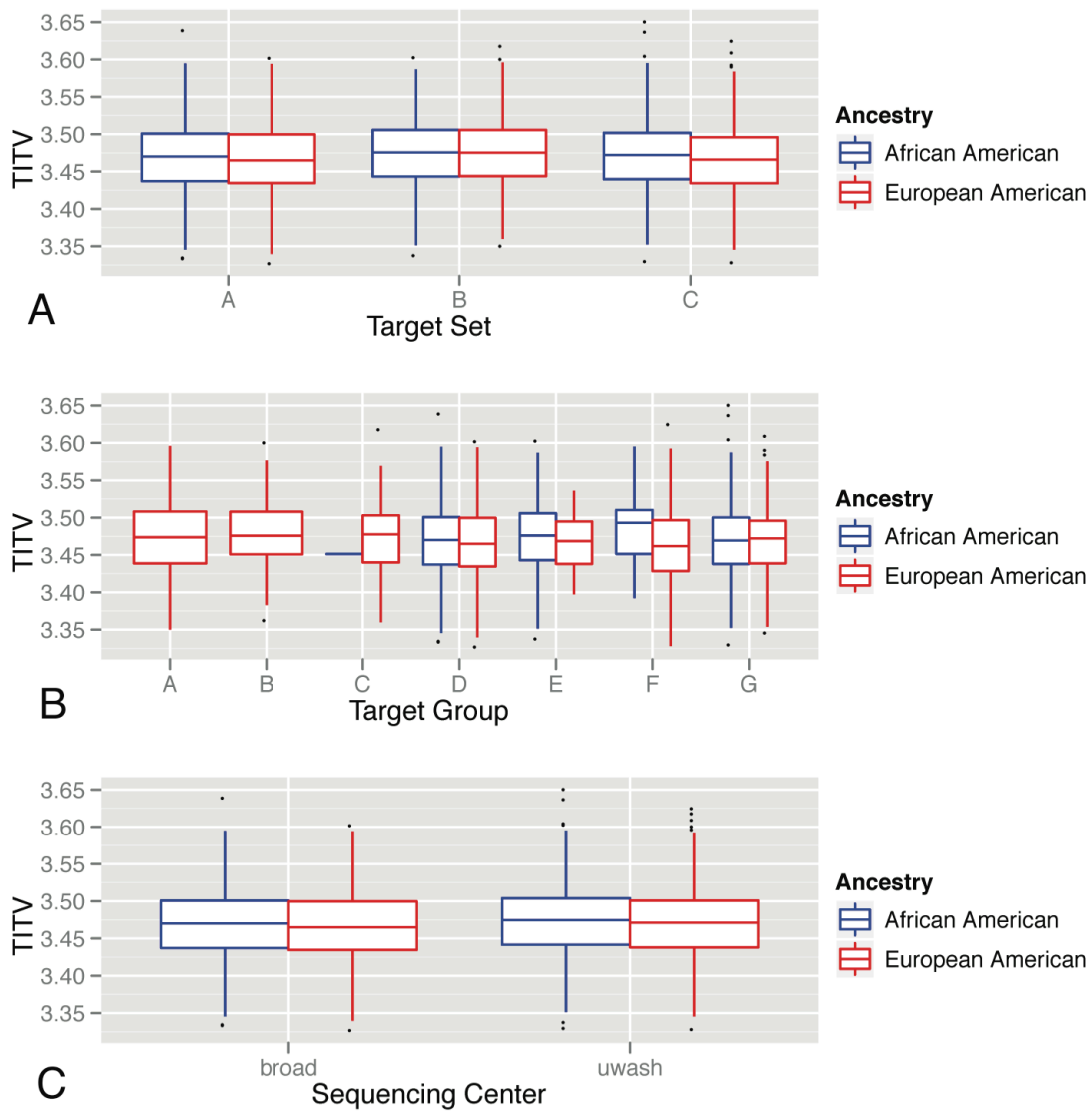
## Supplementary figures



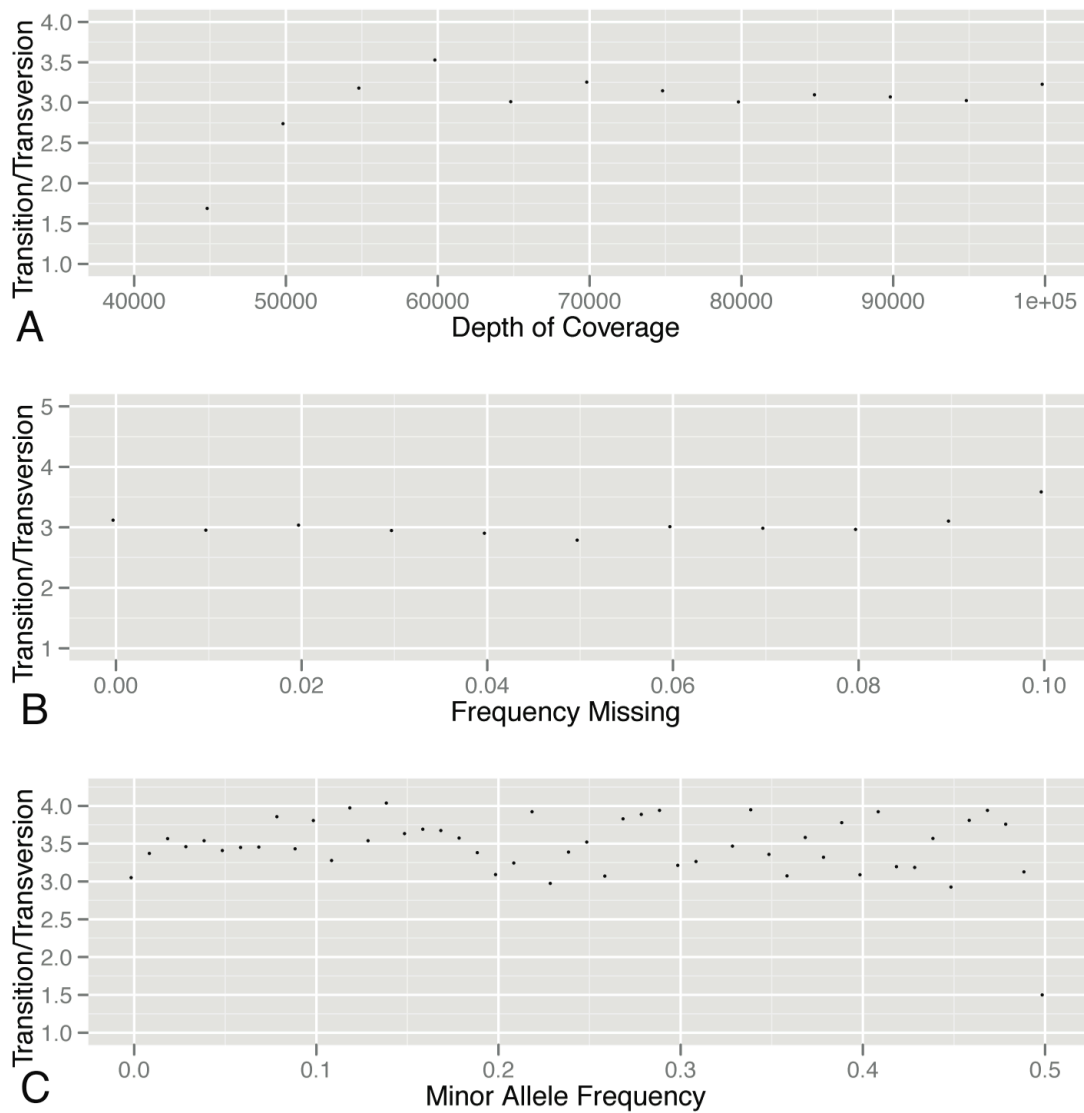
**Figure S1.** Assessment of Sample Outliers using Ti/Tv, Heterozygosity and Missingness (i.e., sites not called). A) Ti/Tv and B) Het/Hom Non-Ref and C) Frequency of Missingness vs. Sample Index. Each dot represents a single sample. Samples showed similar properties of Ti/Tv, heterozygosity and frequency of missing genotypes.



**Figure S2.** Assessment of Batch Effects between Samples Examining Heterozygosity. A) target definition and B) target Group and C) sequencing center vs. Het / Hom Non-Ref.

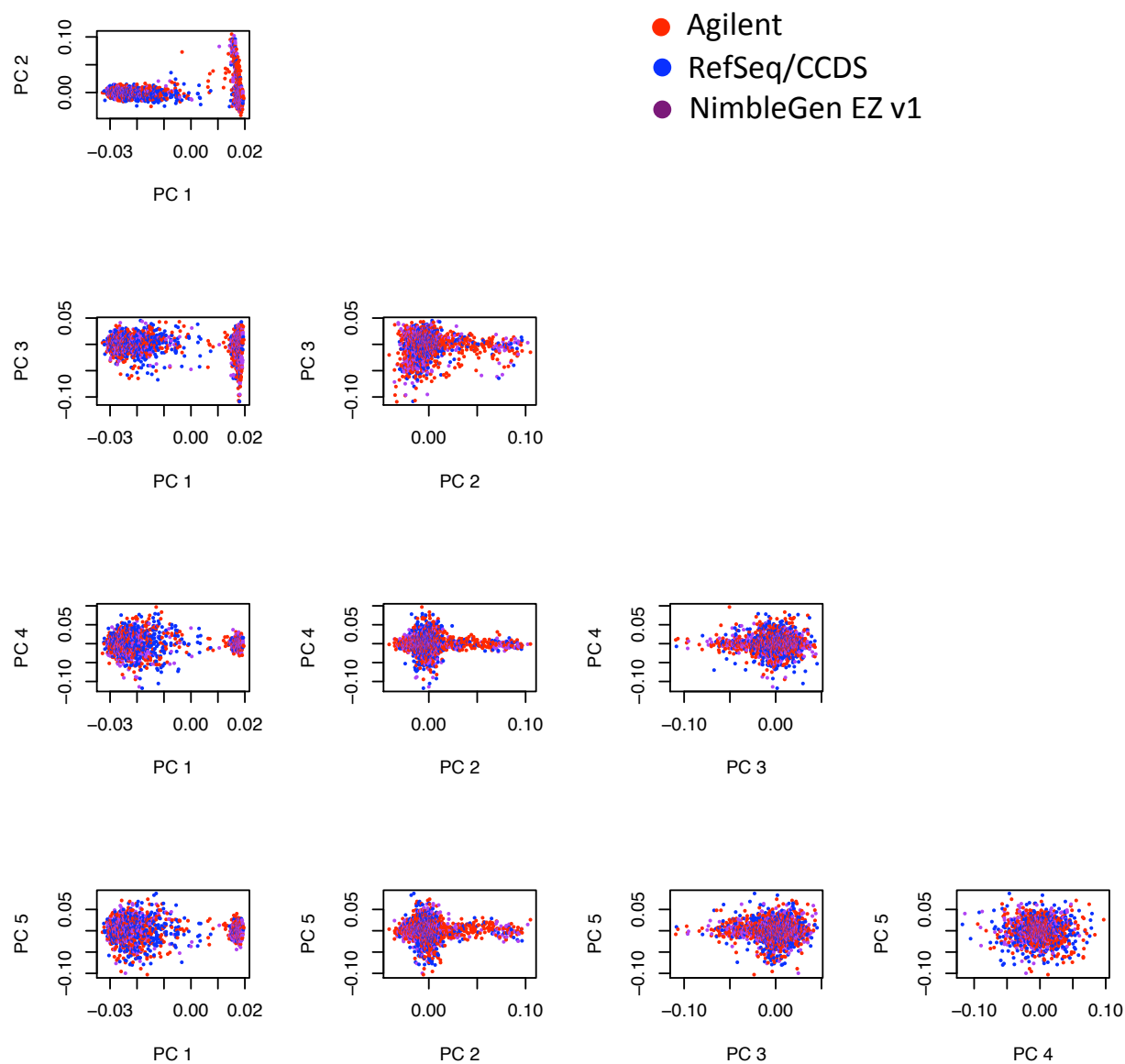


**Figure S3.** Assessment of Batch Effects between Samples using Ti/Tv. A) Target definition. B) Target group. C) Sequencing center vs. Ti/Tv.

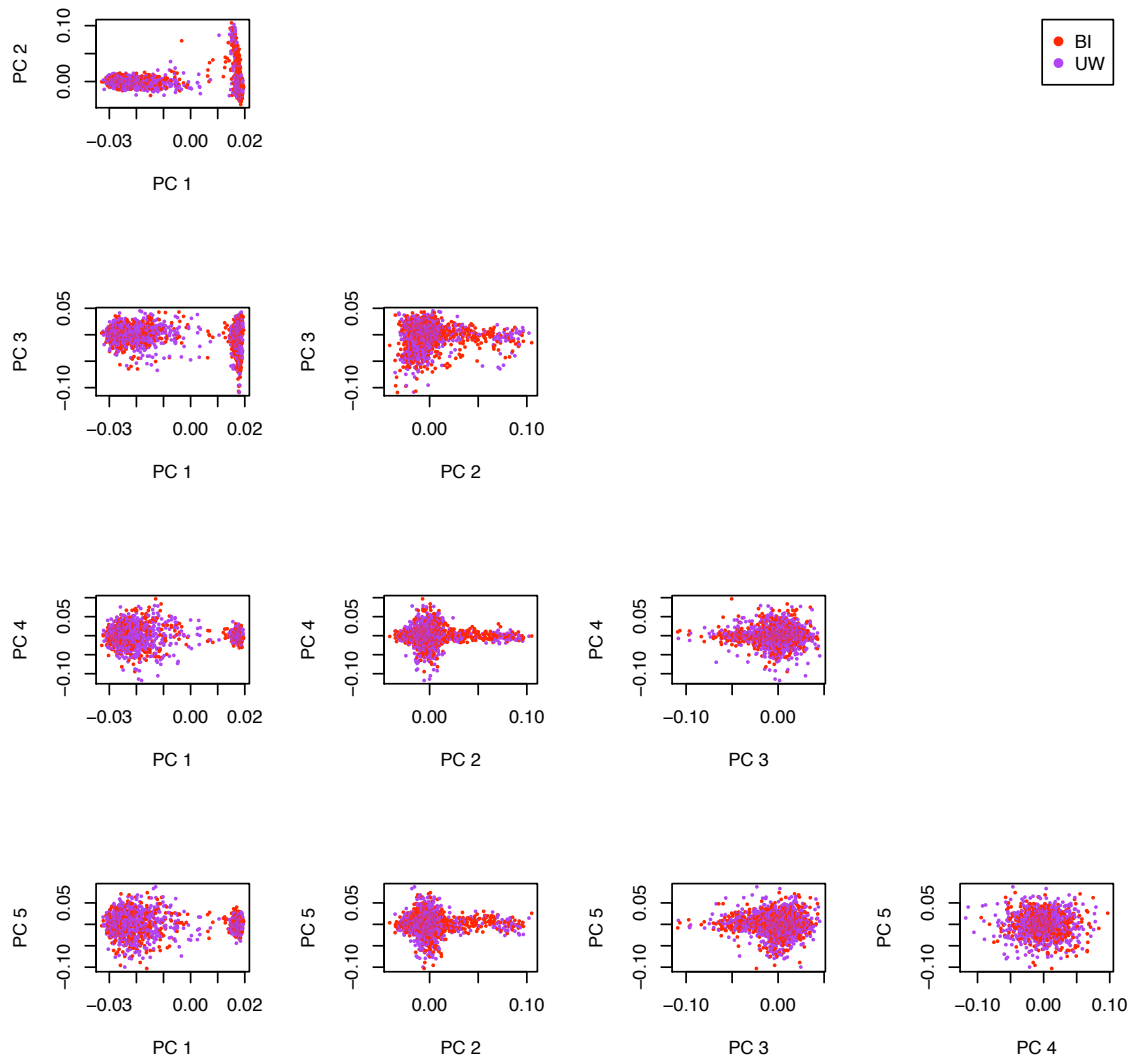


**Figure S4.** Variant QC metrics as a function of TiTv. A) depth of coverage and B) frequency missing and C) minor allele frequency. Ti/Tv is relatively constant over a broad range of coverage. Ti/Tv is similar across frequencies missing < 10% and across all minor allele frequencies.

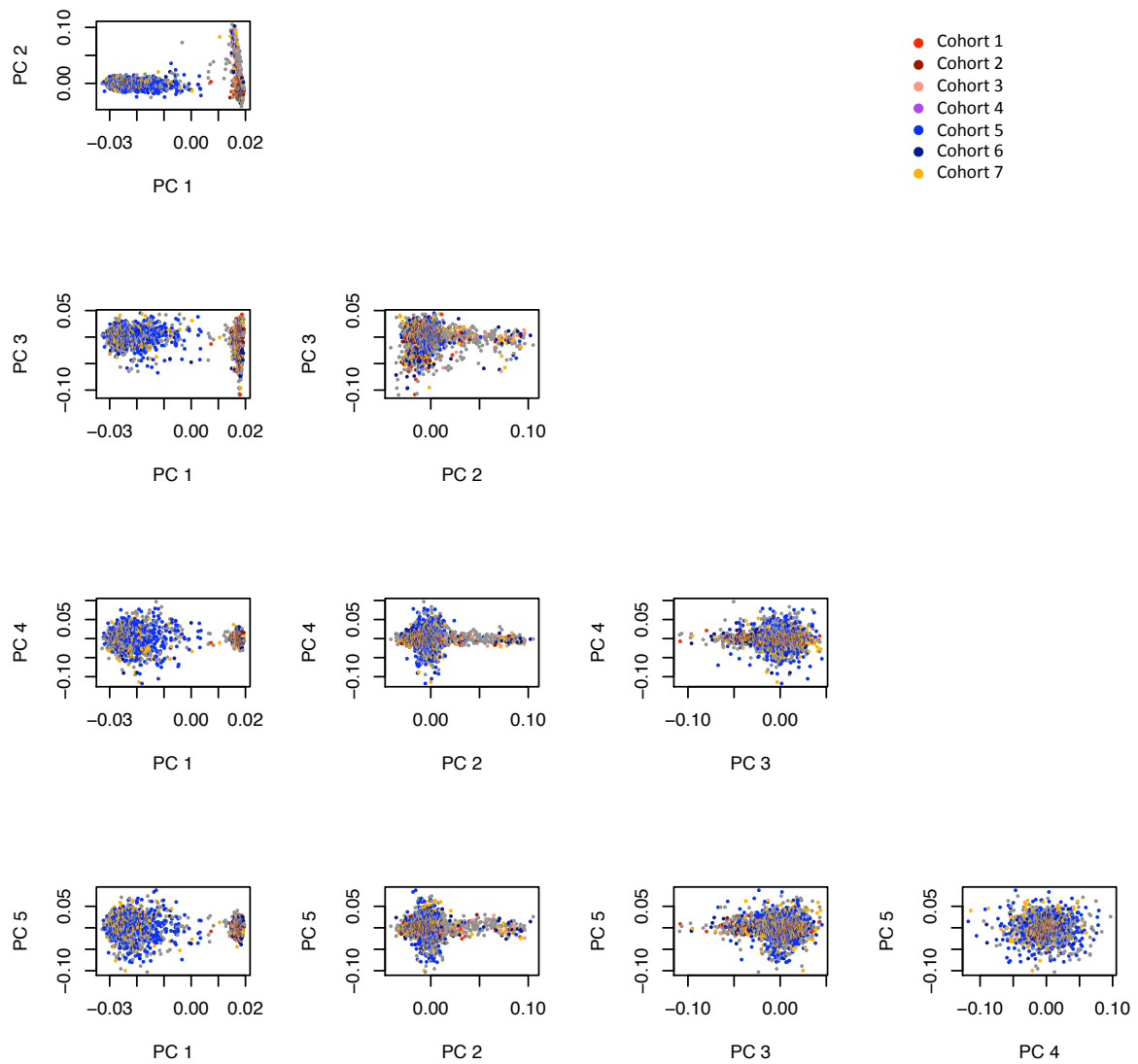




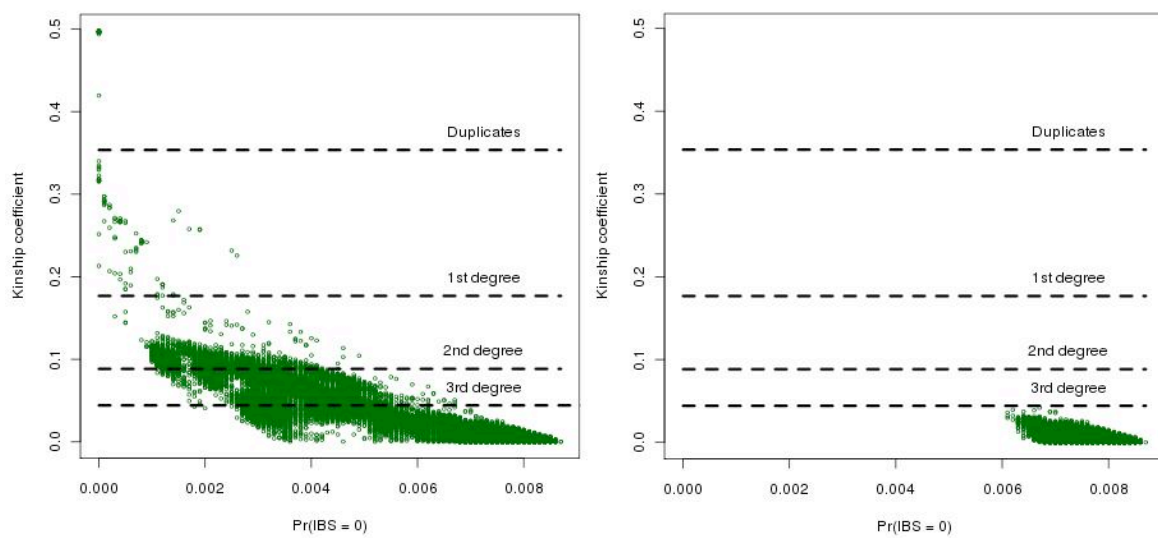
**Figure S5.** PCA plots with samples labeled by exome capture target.



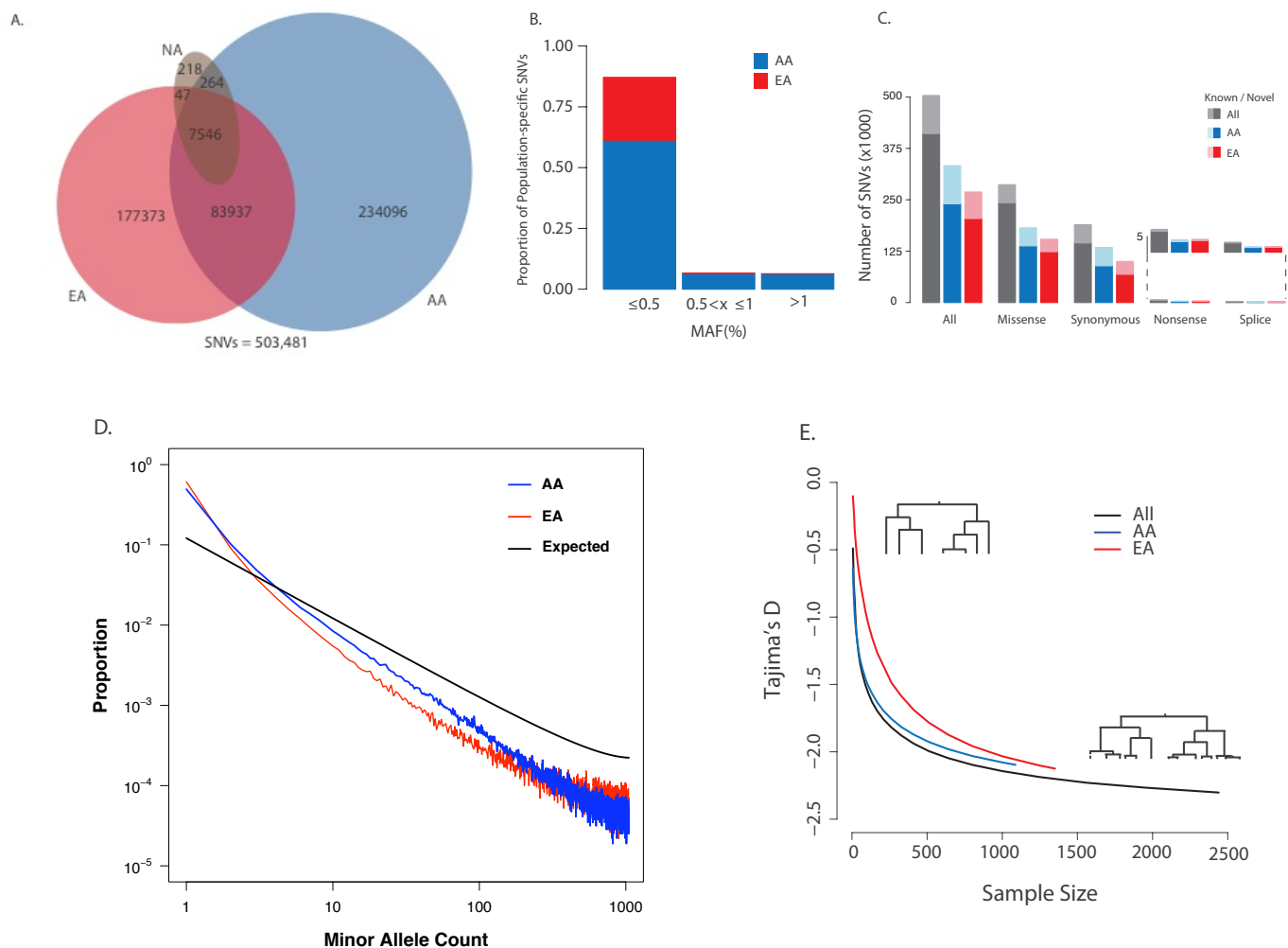
**Figure S6.** PCA plots with samples labeled by sequencing center. BI and WI denote the Broad Institute and University of Washington, respectively.



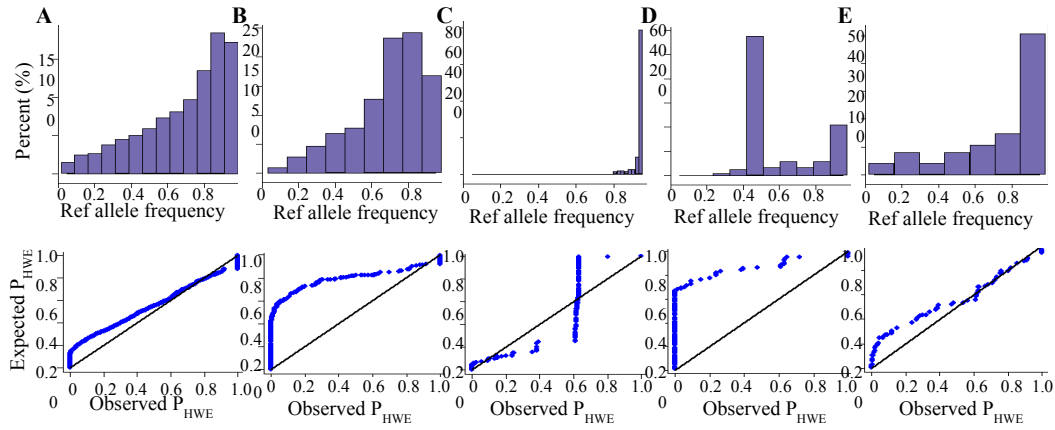
**Figure S7.** PCA plots with samples labeled by phenotype cohort.



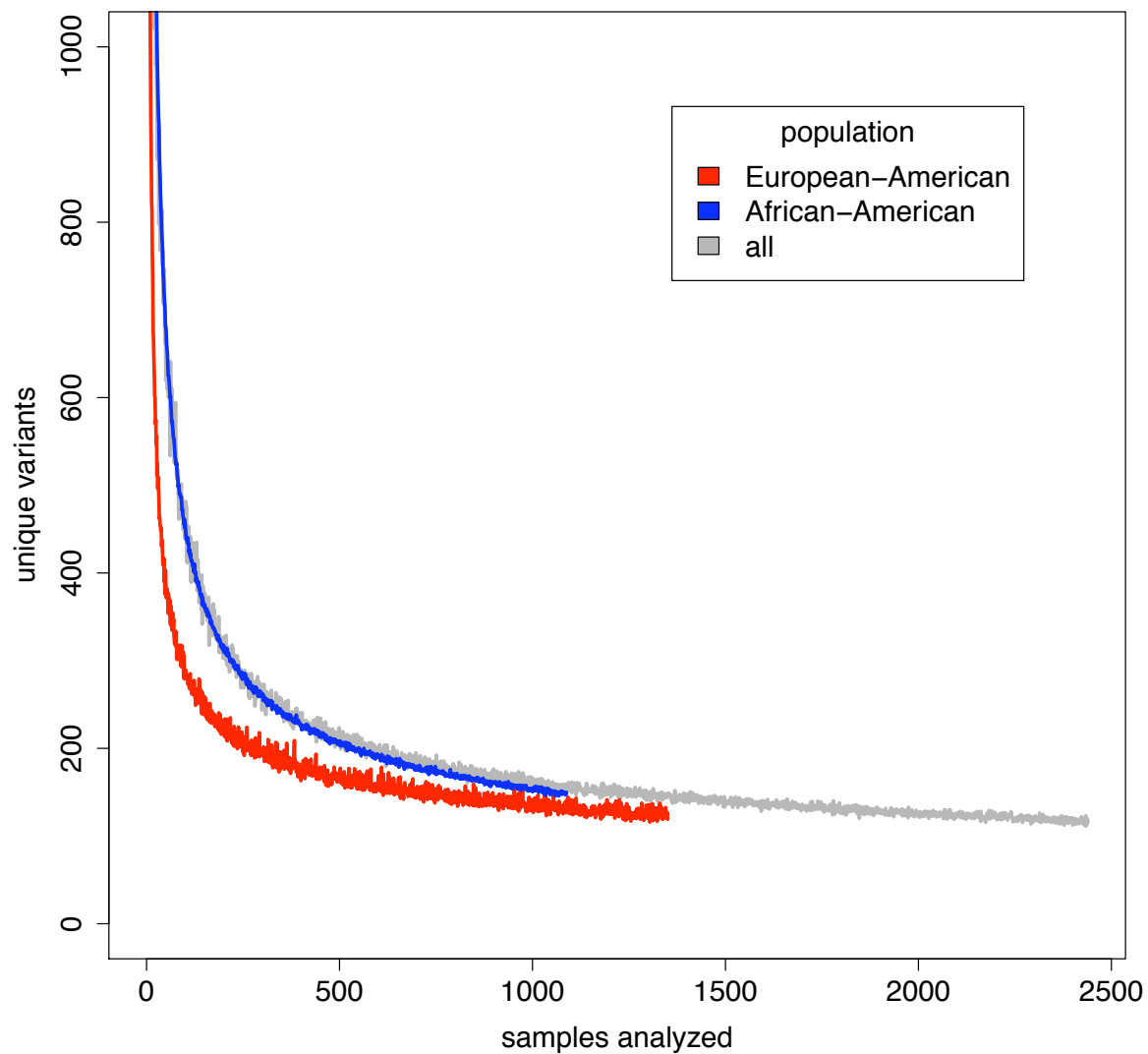
**Figure S8.** Kinship vs. IBS before (left) and after (right) filtering of related individuals and low quality samples.



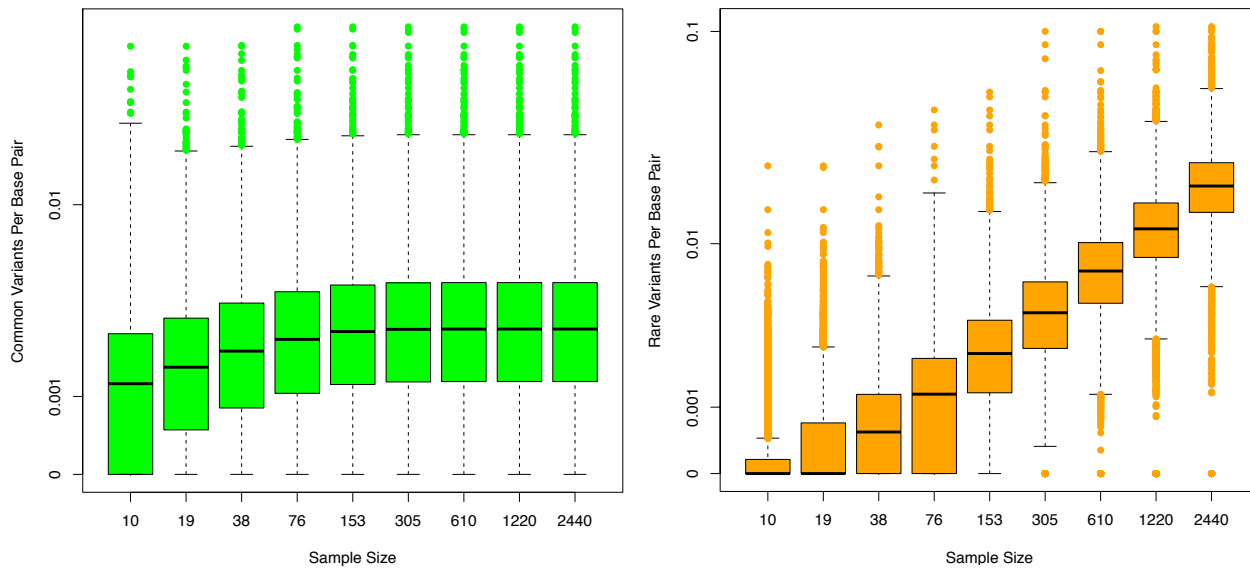
**Figure S9.** Characteristics of human protein-coding variation. A) Venn diagram showing the number of SNVs discovered in each sample, including the single Native American (NA) individual. Note, circles are not drawn to scale. B) Proportion of all non-singleton population-specific SNVs for different bins of MAF. Each bar also shows the relative proportion of population-specific SNVs observed in the AA and EA samples. C) Number of known and novel SNVs discovered for different site types. D) SFS of 1000 randomly selected individuals per sample. Note the large excess of rare SNVs relative to what is expected under a standard neutral model. E) The value of Tajima's  $D$  as a function of sample size. As sample size increases, the value of Tajima's  $D$  becomes sharply negative in each sample due to the increased resolution of observing recent population expansions (notice the primary consequence of increased sample size on genealogies is the addition of recent coalescent events).



**Figure S10.** The performance of 5,561 SNPs that were only observed in European ancestral populations from 1KGP but not observed in EAs from ESP. The frequency distribution for reference allele were illustrated in the upper row; and the Q-Q plot for  $P$  values based upon HWE test were illustrated in the bottom row. **A** for C90, **B** for HWD, **C** for INV, **D** for NoC, **E** for OTH.

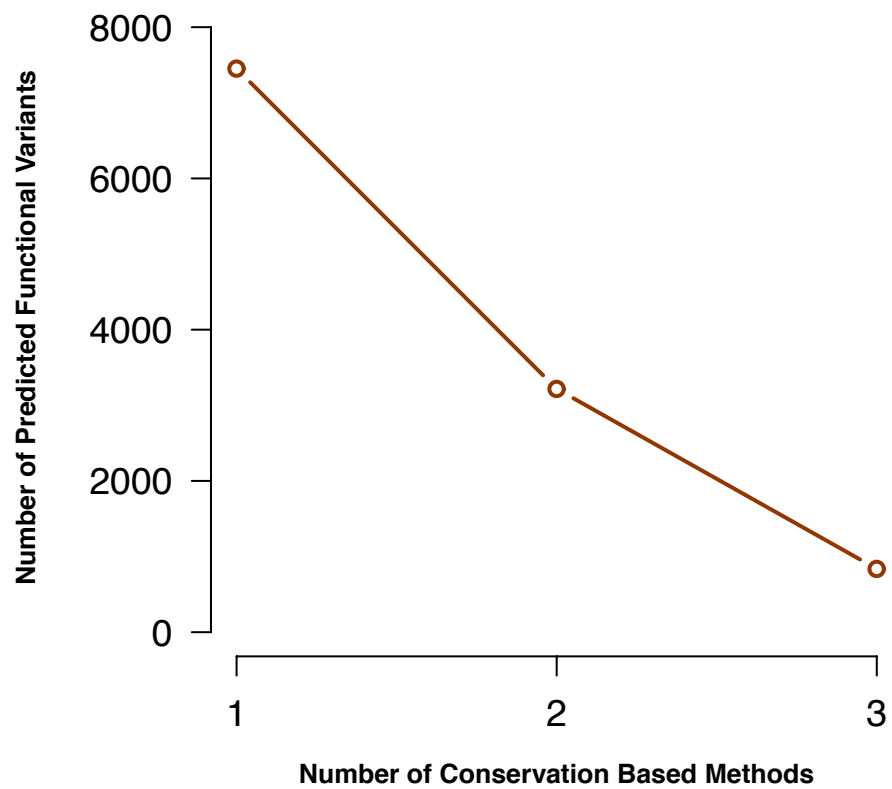


**Figure S11.** Number of unique variants identified as a function of sample size. Each point is the mean across 100 random samples of the desired number of individuals. Calculation of the number of unique variants was done for EA and AA populations individually and with both populations combined.

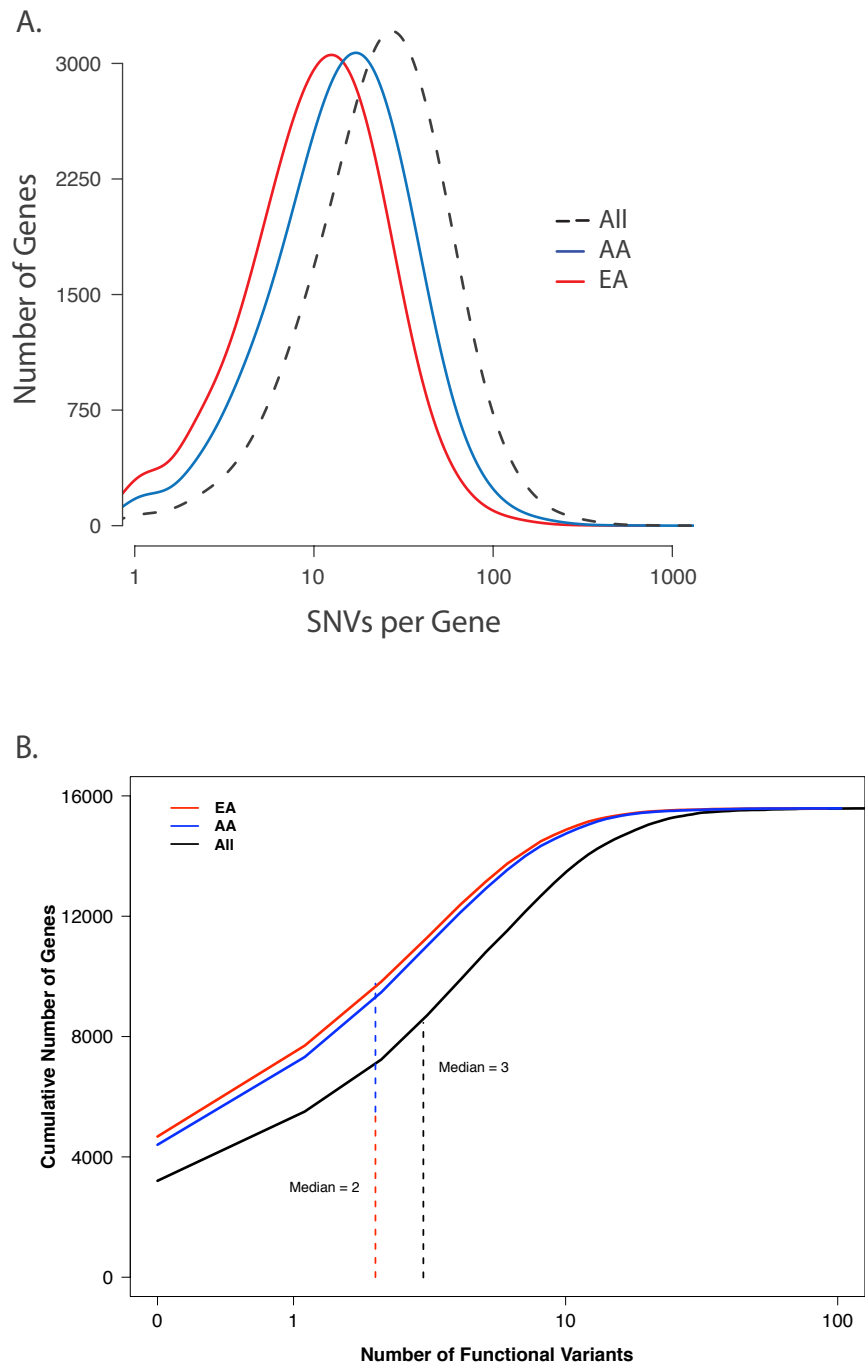


**Figure S12.** Number of common (left) and rare (right) variants per base pair as a function of sample size. A) Box plots showing the number of variants per bp with a MAF  $>0.5\%$  as a function of sample size, estimated for each gene. Box edges represent the first and third quartile, while dark horizontal lines represent the median. Whiskers represent 1.5 times the interquartile range. Points represent outlier genes beyond the range of the whiskers. The number of common variants per base pair increases from a median of 0.0012 for a sample size of ten to a median of 0.0026 in the full dataset of 2440 individuals. Samples of approximately 100 individuals or more capture essentially all common variants. B) Box plots showing the number of rare ( $<0.5\%$ ) variants per bp as a function of sample size. Layout is the same as in (A). The number of rare variants per base pair increases from a median of zero for sample sizes under twenty to a median of 0.019 in the full dataset of 2440 individuals. Each increase in sample size reveals a linear increase of rare variants, with no sign of leveling off.

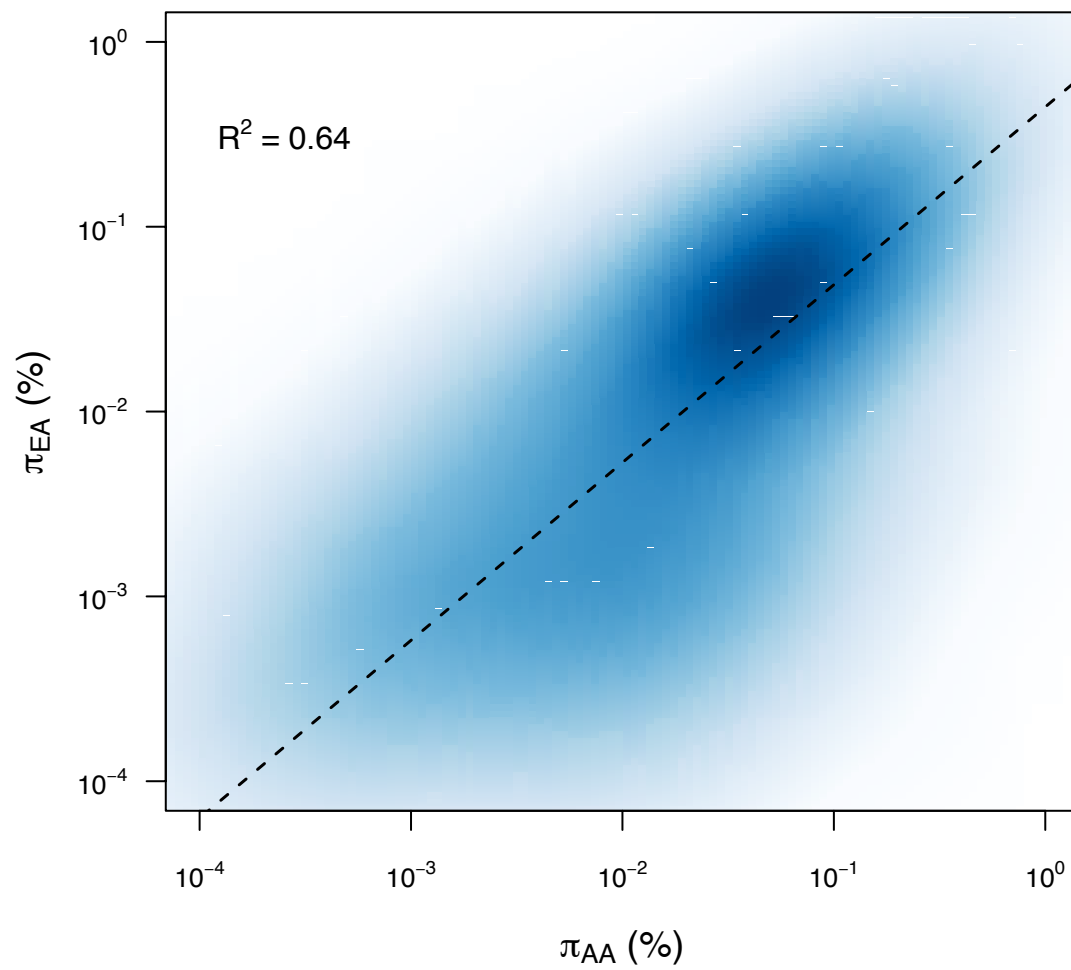




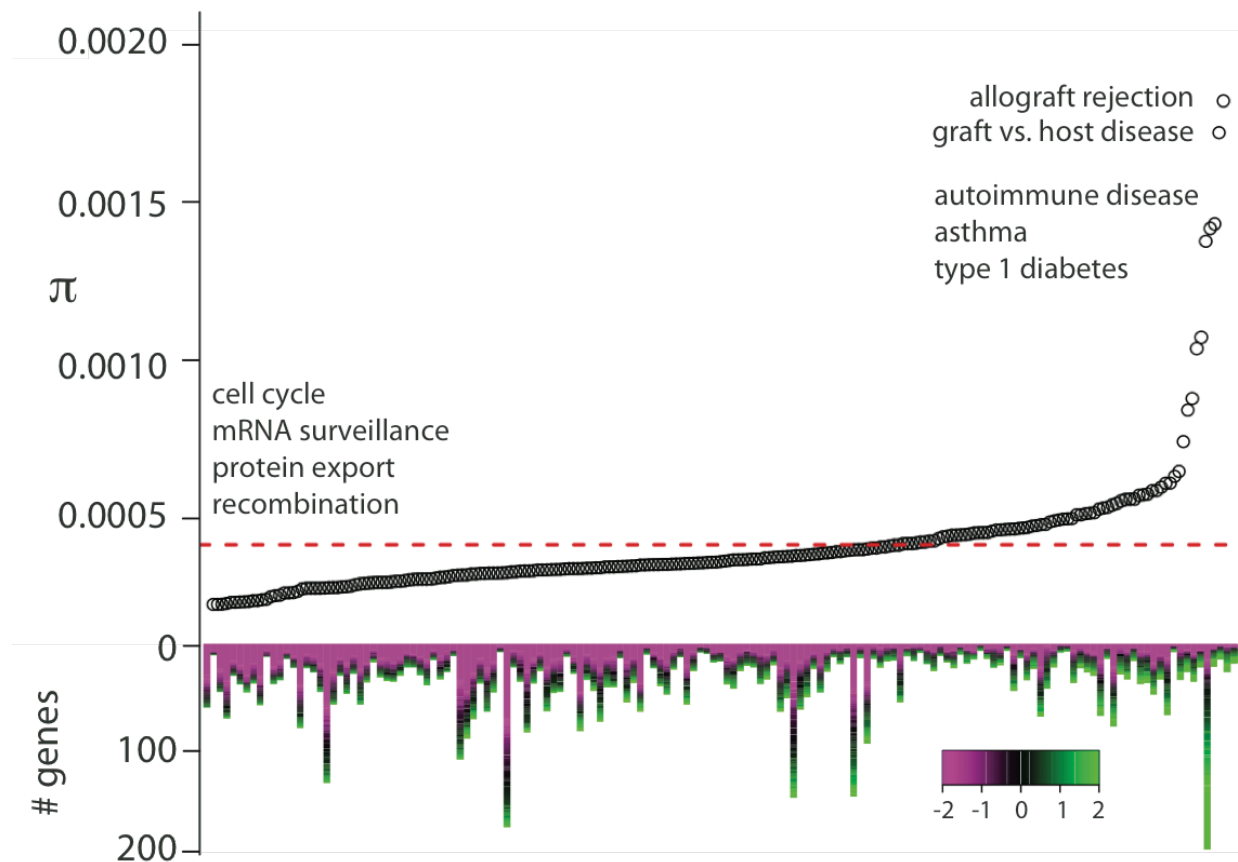
**Figure S13.** Number of synonymous SNVs predicted to be functionally important as a function of three different methods.



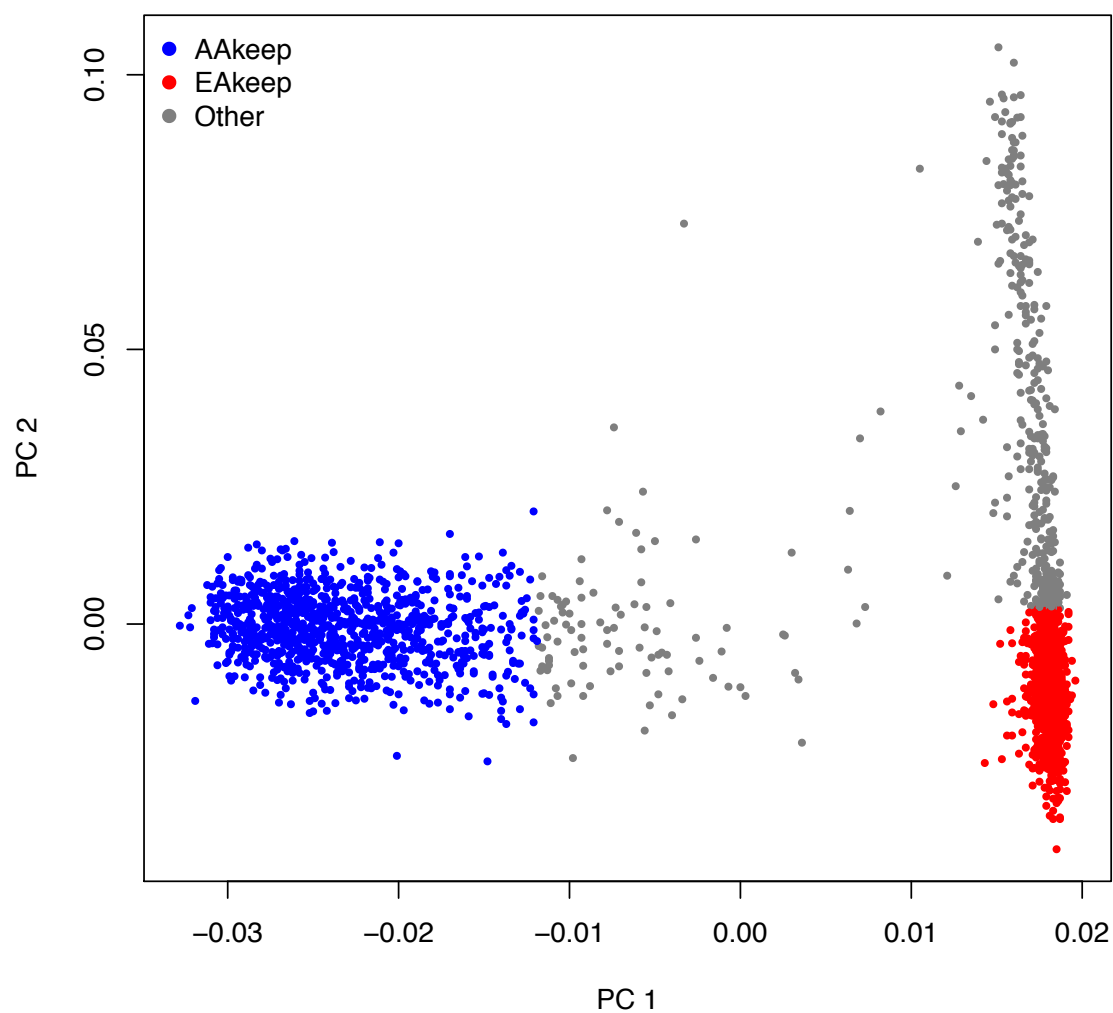
**Figure S14. Distribution of protein-coding variation in genes, pathways, and across the genome.** A) Number of SNVs per gene (note, to facilitate comparisons 1000 individuals were selected at random for each sample). B) Number of predicted functionally important SNVs per gene.



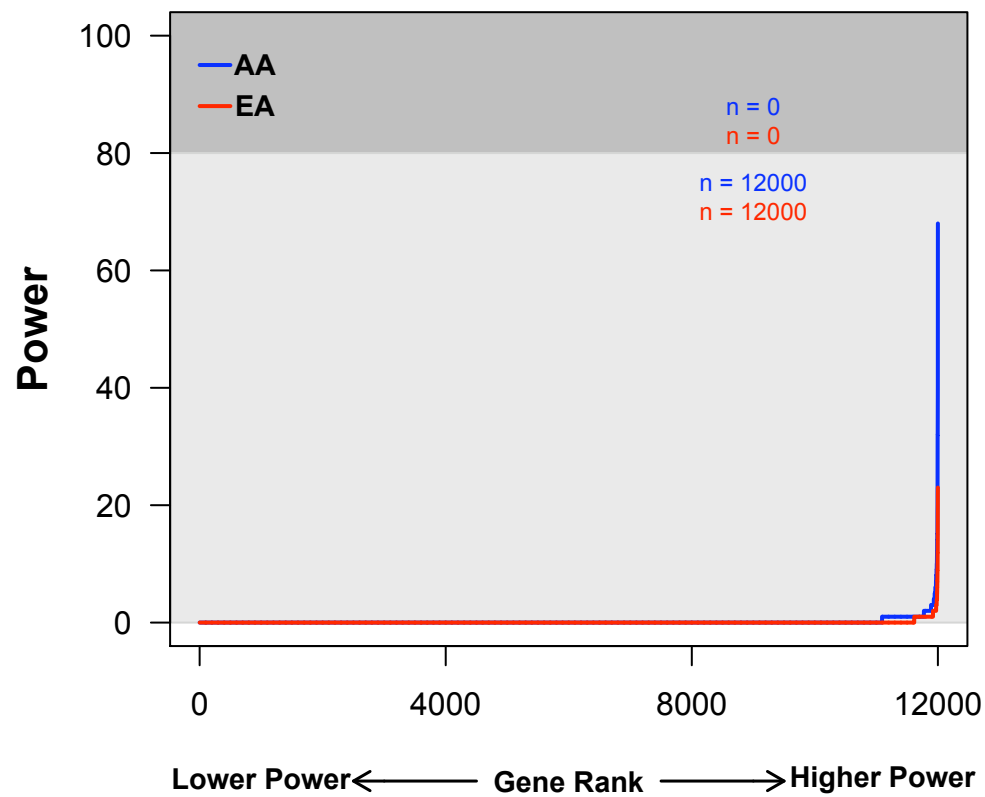
**Figure S15.** Smoothed scatter plot of nucleotide diversity in the EA versus AA samples. Darker shading increased a higher density of points.



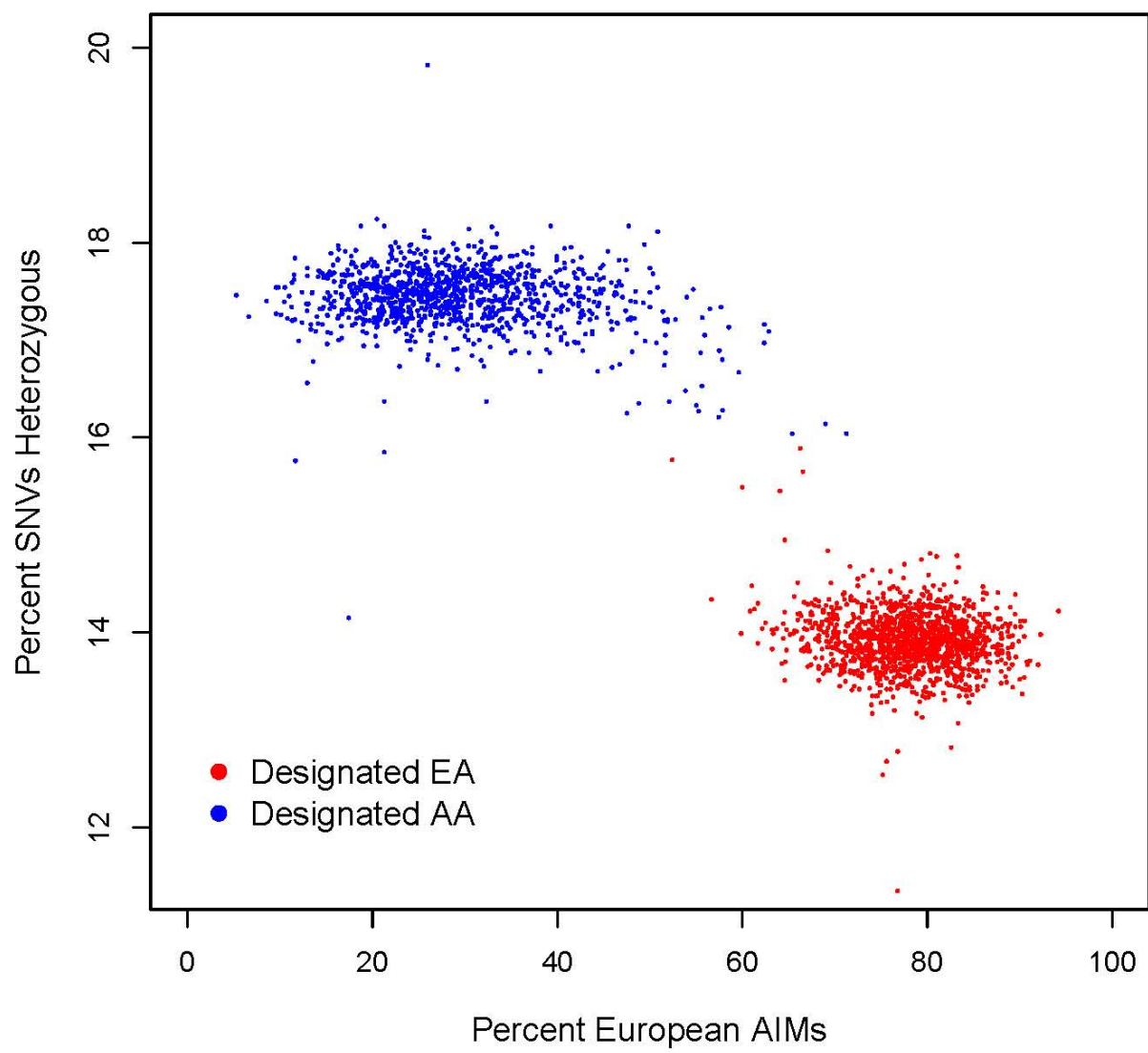
**Figure S16.** Average  $\pi$  for KEGG Pathways. The dashed red line denotes the average  $\pi$  across all pathways. The rows of the heatmap correspond to individual genes in each pathway, which are colored by their value of  $\pi$  (expressed as the log<sub>2</sub> fold change relative to the average  $\pi$  across pathways).



**Figure S17.** Summary of the 1,000 individuals selected in each sample (denoted EAkeep and AAkeep) for the simulations to evaluate the power of rare variant association studies. The goal of selecting these individuals was to mitigate spurious associations due to population stratification.



**Figure S18.** Distribution of gene-specific estimates of power to map causal rare variants across 12,000 protein-coding genes with at least three SNVs in the EA or AA samples assuming an OR for causal variants of 1.5.



**Figure S19.** Assignment of AA and EA ancestry based on heterozygosity and percent European AIMS.

## References and Notes

1. M. J. Bamshad *et al.*, Exome sequencing as a tool for Mendelian disease gene discovery. *Nat. Rev. Genet.* **12**, 745 (2011). [doi:10.1038/nrg3031](https://doi.org/10.1038/nrg3031) [Medline](#)
2. S. S. Ajay, S. C. Parker, H. O. Abaan, K. V. Fajardo, E. H. Margulies, Accurate and comprehensive sequencing of personal genomes. *Genome Res.* **21**, 1498 (2011). [doi:10.1101/gr.123638.111](https://doi.org/10.1101/gr.123638.111) [Medline](#)
3. N. L. Sobreira *et al.*, Whole-genome sequencing of a single proband together with linkage analysis identifies a Mendelian disease gene. *PLoS Genet.* **6**, e1000991 (2010). [doi:10.1371/journal.pgen.1000991](https://doi.org/10.1371/journal.pgen.1000991) [Medline](#)
4. International HapMap Consortium, A haplotype map of the human genome. *Nature* **437**, 1299 (2005). [doi:10.1038/nature04226](https://doi.org/10.1038/nature04226) [Medline](#)
5. K. A. Frazer *et al.*; International HapMap Consortium, A second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**, 851 (2007). [doi:10.1038/nature06258](https://doi.org/10.1038/nature06258) [Medline](#)
6. J. Z. Li *et al.*, Worldwide human relationships inferred from genome-wide patterns of variation. *Science* **319**, 1100 (2008). [doi:10.1126/science.1153717](https://doi.org/10.1126/science.1153717) [Medline](#)
7. Y. X. Fu, Statistical properties of segregating sites. *Theor. Popul. Biol.* **48**, 172 (1995). [doi:10.1006/tpbi.1995.1025](https://doi.org/10.1006/tpbi.1995.1025) [Medline](#)
8. G. T. Marth *et al.*; the 1000 Genomes Project, The functional spectrum of low-frequency coding variation. *Genome Biol.* **12**, R84 (2011). [doi:10.1186/gb-2011-12-9-r84](https://doi.org/10.1186/gb-2011-12-9-r84) [Medline](#)
9. S. B. Ng *et al.*, Targeted capture and massively parallel sequencing of 12 human exomes. *Nature* **461**, 272 (2009). [doi:10.1038/nature08250](https://doi.org/10.1038/nature08250) [Medline](#)
10. B. J. O’Roak *et al.*, Exome sequencing in sporadic autism spectrum disorders identifies severe de novo mutations. *Nat. Genet.* **43**, 585 (2011). [doi:10.1038/ng.835](https://doi.org/10.1038/ng.835) [Medline](#)
11. S. B. Ng *et al.*, Exome sequencing identifies MLL2 mutations as a cause of Kabuki syndrome. *Nat. Genet.* **42**, 790 (2010). [doi:10.1038/ng.646](https://doi.org/10.1038/ng.646) [Medline](#)
12. S. B. Ng *et al.*, Exome sequencing identifies the cause of a mendelian disorder. *Nat. Genet.* **42**, 30 (2010). [doi:10.1038/ng.499](https://doi.org/10.1038/ng.499) [Medline](#)
13. J. A. Tennessen, J. Madeoy, J. M. Akey, Signatures of positive selection apparent in a small sample of human exomes. *Genome Res.* **20**, 1327 (2010). [doi:10.1101/gr.106161.110](https://doi.org/10.1101/gr.106161.110) [Medline](#)
14. X. Yi *et al.*, Sequencing of 50 human exomes reveals adaptation to high altitude. *Science* **329**, 75 (2010). [doi:10.1126/science.1190371](https://doi.org/10.1126/science.1190371) [Medline](#)
15. J. McClellan, M. C. King, Genetic heterogeneity in human disease. *Cell* **141**, 210 (2010). [doi:10.1016/j.cell.2010.03.032](https://doi.org/10.1016/j.cell.2010.03.032) [Medline](#)
16. T. A. Manolio *et al.*, Finding the missing heritability of complex diseases. *Nature* **461**, 747 (2009). [doi:10.1038/nature08494](https://doi.org/10.1038/nature08494) [Medline](#)



17. G. Gibson, Rare and common variants: twenty arguments. *Nat. Rev. Genet.* **13**, 135 (2011). [doi:10.1038/nrg3118](https://doi.org/10.1038/nrg3118) [Medline](#)
18. Supporting material is available on *Science Online*.
19. M. Kimura, Evolutionary rate at the molecular level. *Nature* **217**, 624 (1968). [doi:10.1038/217624a0](https://doi.org/10.1038/217624a0) [Medline](#)
20. A. Ramírez-Soriano, R. Nielsen, Correcting estimators of theta and Tajima's D for ascertainment biases caused by the single-nucleotide polymorphism discovery process. *Genetics* **181**, 701 (2009). [doi:10.1534/genetics.108.094060](https://doi.org/10.1534/genetics.108.094060) [Medline](#)
21. J. M. Akey *et al.*, Population history and natural selection shape patterns of genetic variation in 132 genes. *PLoS Biol.* **2**, e286 (2004). [doi:10.1371/journal.pbio.0020286](https://doi.org/10.1371/journal.pbio.0020286) [Medline](#)
22. A. Coventry *et al.*, Deep resequencing reveals excess rare recent variants consistent with explosive population growth. *Nat Commun* **1**, 131 (2010). [doi:10.1038/ncomms1130](https://doi.org/10.1038/ncomms1130) [Medline](#)
23. S. Gravel *et al.*; 1000 Genomes Project, Demographic history and rare allele sharing among human populations. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 11983 (2011). [doi:10.1073/pnas.1019276108](https://doi.org/10.1073/pnas.1019276108) [Medline](#)
24. M. dos Reis, R. Savva, L. Wernisch, Solving the riddle of codon usage preferences: a test for translational selection. *Nucleic Acids Res.* **32**, 5036 (2004). [doi:10.1093/nar/gkh834](https://doi.org/10.1093/nar/gkh834) [Medline](#)
25. J. H. McDonald, M. Kreitman, Adaptive protein evolution at the Adh locus in *Drosophila*. *Nature* **351**, 652 (1991). [doi:10.1038/351652a0](https://doi.org/10.1038/351652a0) [Medline](#)
26. J. M. Akey, Constructing genomic maps of positive selection in humans: where do we go from here? *Genome Res.* **19**, 711 (2009). [doi:10.1101/gr.086652.108](https://doi.org/10.1101/gr.086652.108) [Medline](#)
27. J. Asimit, E. Zeggini, Rare variant association analysis methods for complex traits. *Annu. Rev. Genet.* **44**, 293 (2010). [doi:10.1146/annurev-genet-102209-163421](https://doi.org/10.1146/annurev-genet-102209-163421) [Medline](#)
28. G. V. Kryukov, A. Shpunt, J. A. Stamatoyannopoulos, S. R. Sunyaev, Power of deep, all-exon resequencing for discovery of human trait genes. *Proc. Natl. Acad. Sci. U.S.A.* **106**, 3871 (2009). [doi:10.1073/pnas.0812824106](https://doi.org/10.1073/pnas.0812824106) [Medline](#)
29. K. E. Lohmueller *et al.*, Proportionally more deleterious genetic variation in European than in African populations. *Nature* **451**, 994 (2008). [doi:10.1038/nature06611](https://doi.org/10.1038/nature06611) [Medline](#)
30. L. Yu, F. D. Martinez, W. T. Klimecki, Automated high-throughput sex-typing assay. *Biotechniques* **37**, 662 (2004). [Medline](#)
31. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754 (2009). [doi:10.1093/bioinformatics/btp324](https://doi.org/10.1093/bioinformatics/btp324) [Medline](#)

32. M. A. DePristo *et al.*, A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491 (2011).  
[doi:10.1038/ng.806](https://doi.org/10.1038/ng.806) [Medline](#)
33. S. Fisher *et al.*, A scalable, fully automated process for construction of sequence-ready human exome targeted capture libraries. *Genome Biol.* **12**, R1 (2011).  
[doi:10.1186/gb-2011-12-1-r1](https://doi.org/10.1186/gb-2011-12-1-r1) [Medline](#)
34. Y. Li, C. Sidore, H. M. Kang, M. Boehnke, G. R. Abecasis, Low-coverage sequencing: implications for design of complex trait association studies. *Genome Res.* **21**, 940 (2011). [doi:10.1101/gr.117259.110](https://doi.org/10.1101/gr.117259.110) [Medline](#)
35. H. Li, J. Ruan, R. Durbin, Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome Res.* **18**, 1851 (2008).  
[doi:10.1101/gr.078212.108](https://doi.org/10.1101/gr.078212.108) [Medline](#)
36. R. M. Durbin *et al.*; 1000 Genomes Project Consortium, A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061 (2010).  
[doi:10.1038/nature09534](https://doi.org/10.1038/nature09534) [Medline](#)
37. A. Manichaikul *et al.*, Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867 (2010). [doi:10.1093/bioinformatics/btq559](https://doi.org/10.1093/bioinformatics/btq559) [Medline](#)
38. I. A. Adzhubei *et al.*, A method and server for predicting damaging missense mutations. *Nat. Methods* **7**, 248 (2010). [doi:10.1038/nmeth0410-248](https://doi.org/10.1038/nmeth0410-248) [Medline](#)
39. P. Kumar, S. Henikoff, P. C. Ng, Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm. *Nat. Protoc.* **4**, 1073 (2009). [doi:10.1038/nprot.2009.86](https://doi.org/10.1038/nprot.2009.86) [Medline](#)
40. S. Chun, J. C. Fay, Identification of deleterious mutations within three human genomes. *Genome Res.* **19**, 1553 (2009). [doi:10.1101/gr.092619.109](https://doi.org/10.1101/gr.092619.109) [Medline](#)
41. J. M. Schwarz, C. Rödelberger, M. Schuelke, D. Seelow, MutationTaster evaluates disease-causing potential of sequence alterations. *Nat. Methods* **7**, 575 (2010).  
[doi:10.1038/nmeth0810-575](https://doi.org/10.1038/nmeth0810-575) [Medline](#)
42. G. M. Cooper *et al.*, Single-nucleotide evolutionary constraint scores highlight disease-causing mutations. *Nat. Methods* **7**, 250 (2010). [doi:10.1038/nmeth0410-250](https://doi.org/10.1038/nmeth0410-250) [Medline](#)
43. G. M. Cooper *et al.*, NISC Comparative Sequencing Program, Distribution and intensity of constraint in mammalian genomic sequence. *Genome Res.* **15**, 901 (2005). [doi:10.1101/gr.3577405](https://doi.org/10.1101/gr.3577405) [Medline](#)
44. M. E. Zwick, D. J. Cutler, A. Chakravarti, Patterns of genetic variation in Mendelian and complex traits. *Annu. Rev. Genomics Hum. Genet.* **1**, 387 (2000).  
[doi:10.1146/annurev.genom.1.1.387](https://doi.org/10.1146/annurev.genom.1.1.387) [Medline](#)
45. R. N. Gutenkunst, R. D. Hernandez, S. H. Williamson, C. D. Bustamante, Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genet.* **5**, e1000695 (2009).  
[doi:10.1371/journal.pgen.1000695](https://doi.org/10.1371/journal.pgen.1000695) [Medline](#)

46. S. Gravel *et al.*; 1000 Genomes Project, Demographic history and rare allele sharing among human populations. *Proc. Natl. Acad. Sci. U.S.A.* **108**, 11983 (2011).  
[doi:10.1073/pnas.1019276108](https://doi.org/10.1073/pnas.1019276108) [Medline](#)
47. Y. Y. Waldman, T. Tuller, T. Shlomi, R. Sharan, E. Ruppin, Translation efficiency in humans: tissue specificity, global optimization and differences between developmental stages. *Nucleic Acids Res.* **38**, 2964 (2010).  
[doi:10.1093/nar/gkq009](https://doi.org/10.1093/nar/gkq009) [Medline](#)