

UC San Diego

UC San Diego Previously Published Works

Title

Analysis of protein-coding genetic variation in 60,706 humans.

Permalink

<https://escholarship.org/uc/item/04j4327s>

Journal

Nature, 536(7616)

ISSN

0028-0836

Authors

Lek, Monkol
Karczewski, Konrad J
Minikel, Eric V
[et al.](#)

Publication Date

2016-08-01

DOI

10.1038/nature19057

Peer reviewed



Published in final edited form as:

Nature. 2016 August 18; 536(7616): 285–291. doi:10.1038/nature19057.

Analysis of protein-coding genetic variation in 60,706 humans

A full list of authors and affiliations appears at the end of the article.

Summary

Large-scale reference data sets of human genetic variation are critical for the medical and functional interpretation of DNA sequence changes. We describe the aggregation and analysis of high-quality exome (protein-coding region) sequence data for 60,706 individuals of diverse ethnicities generated as part of the Exome Aggregation Consortium (ExAC). This catalogue of human genetic diversity contains an average of one variant every eight bases of the exome, and provides direct evidence for the presence of widespread mutational recurrence. We have used this catalogue to calculate objective metrics of pathogenicity for sequence variants, and to identify genes subject to strong selection against various classes of mutation; identifying 3,230 genes with near-complete depletion of truncating variants with 72% having no currently established human disease phenotype. Finally, we demonstrate that these data can be used for the efficient filtering of candidate disease-causing variants, and for the discovery of human “knockout” variants in protein-coding genes.

Background

Over the last five years, the widespread availability of high-throughput DNA sequencing technologies has permitted the sequencing of the whole genomes or exomes (the protein-coding regions of genomes) of hundreds of thousands of humans. In theory, these data represent a powerful source of information about the global patterns of human genetic variation, but in practice, are difficult to access for practical, logistical, and ethical reasons; in addition, their utility is complicated by the heterogeneity in the experimental methodologies and variant calling pipelines used to generate them. Current publicly available datasets of human DNA sequence variation contain only a small fraction of all

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:http://www.nature.com/authors/editorial_policies/license.html#terms

Correspondence to: Daniel G MacArthur.

*These authors contributed equally to this work and names appear in alphabetical order

#List of collaborators to appear at the end of manuscript

Author Contributions

M. Lek, K.J.K., E.V.M., K.E.S., E.B., T.F., A.H.O., J.S.W., A.J.H., B.B.C., T.T., D.P.B., J.A.K., L.D., K.E., F.Z., J.Z., E.P., M.J.D., D.G.M. contributed to the analysis and writing of the manuscript. M. Lek, E.B., T.F., K.J.K., E.V.M., F.Z., D.P.B., J.B., D.N.C., N.D., M.D., R.D., J.F., M.F., L.G., J.G., N.G., D.H., A.K., M.I.K., A.L.M., P.N., L.O., G.M.P., R.P., M.A.R., V.R., S.A.R., D.M.R., K.S., P.D.S., C.S., B.P.T., G.T., M.T.T., B.W., H.W., D.Y., S.B.G., M.J.D., D.G.M. contributed to the production of the ExAC data set. D.M.A., D.A., M.B., J.D., S.D., R.E., J.C.F., S.B.G., G.G., S.J.G., C.M.H., S.K., M. Laakso, S.M., M.I.M., D.M., R.M., B.M.N., A.P., S.M.P., D.S., J.S., P.S., P.F.S., J.T., M.T.T., H.C.W., J.G.W., M.J.D., D.G.M. contributed to the design and conduct of the various exome sequencing studies and critical review of manuscript.

Author Information

P.F.S is a scientific advisor to Pfizer.

ExAC data set is publicly available at <http://exac.broadinstitute.org>

sequenced samples: the Exome Variant Server, created as part of the NHLBI Exome Sequencing Project (ESP)¹, contains frequency information spanning 6,503 exomes; and the 1000 Genomes (1000G) Project, which includes individual-level genotype data from whole-genome and exome sequence data for 2,504 individuals².

Databases of genetic variation are important for our understanding of human population history and biology¹⁻⁵, but also provide critical resources for the clinical interpretation of variants observed in patients suffering from rare Mendelian diseases^{6,7}. The filtering of candidate variants by frequency in unselected individuals is a key step in any pipeline for the discovery of causal variants in Mendelian disease patients, and the efficacy of such filtering depends on both the size and the ancestral diversity of the available reference data.

Here, we describe the joint variant calling and analysis of high-quality variant calls across 60,706 human exomes, assembled by the Exome Aggregation Consortium (ExAC; exac.broadinstitute.org). This call set exceeds previously available exome-wide variant databases by nearly an order of magnitude, providing substantially increased resolution for the analysis of very low-frequency genetic variants. We demonstrate the application of this data set to the analysis of patterns of genetic variation including the discovery of widespread mutational recurrence, the inference of gene-level constraint against truncating variation, the clinical interpretation of variation in Mendelian disease genes, and the discovery of human “knockout” variants in protein-coding genes.

The ExAC Data set

Sequencing data processing, variant calling, quality control and filtering was performed on over 91,000 exomes (see Online Methods), and sample filtering was performed to produce a final data set spanning 60,706 individuals (Figure 1a). To identify the ancestry of each ExAC individual, we performed principal component analysis (PCA) to distinguish the major axes of geographic ancestry and to identify population clusters corresponding to individuals of European, African, South Asian, East Asian, and admixed American (hereafter Latino) ancestry (Figure 1b; Supplementary Information Table 3); we note that the apparent separation between East Asian and other samples reflects a deficiency of Middle Eastern and Central Asian samples in the data set. We further separated Europeans into individuals of Finnish and non-Finnish ancestry given the enrichment of this bottlenecked population; the term “European” hereafter refers to non-Finnish European individuals.

We identified 10,195,872 candidate sequence variants in ExAC. We further applied stringent depth and site/genotype quality filters to define a subset of 7,404,909 high quality (HQ) variants, including 317,381 indels (Supplementary Information Table 7), corresponding to one variant for every 8 bp within the exome intervals. The majority of these are very low-frequency variants absent from previous smaller call sets (Figure 1c): of the HQ variants, 99% have a frequency of <1%, 54% are singletons (variants seen only once in the data set), and 72% are absent from both 1000G and ESP.

The density of variation in ExAC is not uniform across the genome, and the observation of variants depends on factors such as mutational properties and selective pressures. In the

~45M well covered (80% of individuals with a minimum of 10X coverage) positions in ExAC, there are ~18M possible synonymous variants, of which we observe 1.4M (7.5%). However, we observe 63.1% of possible CpG transitions (C to T variants, where the adjacent base is G), while only observing 3% of possible transversions and 9.2% of other possible transitions (Supplementary Information Table 9). A similar pattern is observed for missense and nonsense variants, with lower proportions due to selective pressures (Figure 1D). Of 123,629 HQ insertion/deletions (indels) called in coding exons, 117,242 (95%) have length <6 bases, with shorter deletions being the most common (Figure 1E). Frameshifts are found in smaller numbers and are more likely to be singletons than in-frame indels (Figure 1F), reflecting the influence of purifying selection.

Patterns of protein-coding variation revealed by large samples

The density of protein-coding sequence variation in ExAC reveals a number of properties of human genetic variation undetectable in smaller data sets. For instance, 7.9% of HQ sites in ExAC are multiallelic (multiple different sequence variants observed at the same site), close to the Poisson expectation of 8.3% given the observed density of variation, and far higher than observed in previous data sets - 0.48% in 1000 Genomes (exome intervals) and 0.43% in ESP.

The size of ExAC also makes it possible to directly observe mutational recurrence: instances in which the same mutation has occurred multiple times independently throughout the history of the sequenced populations. For instance, among synonymous variants, a class of variation expected to have undergone minimal selection, 43% of validated *de novo* events identified in external datasets of 1,756 parent-offspring trios^{8,9} are also observed independently in our dataset (Figure 2a), indicating a separate origin for the same variant within the demographic history of the two samples. This proportion is much higher for transition variants at CpG sites, well established to be the most highly mutable sites in the human genome¹⁰: 87% of previously reported *de novo* CpG transitions at synonymous sites are observed in ExAC, indicating that our sample sizes are beginning to approach saturation of this class of variation. This saturation is detectable by a change in the discovery rate at subsets of the ExAC data set, beginning at around 20,000 individuals (Figure 2b), indicating that ExAC is the first human exome-wide dataset large enough for this effect to be directly observed.

Mutational recurrence has a marked effect on the frequency spectrum in the ExAC data, resulting in a depletion of singletons at sites with high mutation rates (Figure 2c). We observe a correlation between singleton rates (the proportion of variants seen only once in ExAC) and site mutability inferred from sequence context¹¹ ($r = -0.98$; $p < 10^{-50}$; Extended Data Figure 1d): sites with low predicted mutability have a singleton rate of 60%, compared to 20% for sites with the highest predicted rate (CpG transitions; Figure 2C). Conversely, for synonymous variants, CpG variants are approximately twice as likely to rise to intermediate frequencies: 16% of CpG variants are found in at least 20 copies in ExAC, compared to 8% of transversions and non-CpG transitions, suggesting that synonymous CpG transitions have on average two independent mutational origins in the ExAC sample. Recurrence at highly mutable sites can further be observed by examining the population sharing of doubleton

synonymous variants (variants occurring in only two individuals in ExAC). Low-mutability mutations (especially transversions), are more likely to be observed in a single population (representing a single mutational origin), while CpG transitions are more likely to be found in two separate populations (independent mutational events); as such, site mutability and probability of observation in two populations is significantly correlated ($r = 0.884$; Figure 2d).

We also explored the prevalence and functional impact of multinucleotide polymorphisms (MNPs), in cases where multiple substitutions were observed within the same codon in at least one individual. We found 5,945 MNPs (mean: 23 per sample) in ExAC (Extended Data Figure 2a) where analysis of the underlying SNPs without correct haplotype phasing would result in altered interpretation. These include 647 instances where the effect of a protein-truncating variant (PTV) variant is eliminated by an adjacent SNP (rescued PTV) and 131 instances where underlying synonymous or missense variants result in PTV MNPs (gained PTV). Additionally our analysis revealed 8 MNPs in disease-associated genes, resulting in either a rescued or gained PTV, and 10 MNPs that have previously been reported as disease causing mutations (Supplementary Information Table 10 and 11). We note that these variants would be missed by virtually all currently available variant calling and annotation pipelines.

Inferring variant deleteriousness and gene constraint

Deleterious variants are expected to have lower allele frequencies than neutral ones, due to negative selection. This theoretical property has been demonstrated previously in human population sequencing data^{12,13} and here (Figure 1d, Figure 1e). This allows inference of the degree of selection against specific functional classes of variation: however, mutational recurrence as described above indicates that allele frequencies observed in ExAC-scale samples are also skewed by mutation rate, with more mutable sites less likely to be singletons (Figure 2c and Extended Data Figure 1d). Mutation rate is in turn non-uniformly distributed across functional classes - for instance, stop lost mutations can never occur at CpG dinucleotides (Extended Data Figure 1e). We corrected for mutation rates (Supplementary Information Section 3.2) by creating a mutability-adjusted proportion singleton (MAPS) metric. This metric reflects (as expected) strong selection against predicted PTVs, as well as missense variants predicted by conservation-based methods to be deleterious (Figure 2e).

The deep ascertainment of rare variation in ExAC also allows us to infer the extent of selection against variant categories on a per-gene basis by examining the proportion of variation that is missing compared to expectations under random mutation. Conceptually similar approaches have been applied to smaller exome datasets^{11,14} but have been underpowered, particularly when analyzing the depletion of PTVs. We compared the observed number of rare (MAF <0.1%) variants per gene to an expected number derived from a selection neutral, sequence-context based mutational model¹¹. The model performs well in predicting the number of synonymous variants, which should be under minimal selection, per gene ($r = 0.98$; Extended Data Figure 3b).

We quantified deviation from expectation with a Z score¹¹, which for synonymous variants is centered at zero, but is significantly shifted towards higher values (greater constraint) for both missense and PTV (Wilcoxon $p < 10^{-50}$ for both; Figure 3a). The genes on the X chromosome are significantly more constrained than those on the autosomes for missense ($p < 10^{-7}$) and loss-of-function ($p < 10^{-50}$), in line with previous work¹⁵. The high correlation between the observed and expected number of synonymous variants on the X chromosome ($r = 0.97$ vs 0.98 for autosomes) indicates that this difference in constraint is not due to a calibration issue. To reduce confounding by coding sequence length for PTVs, we developed an expectation-maximization algorithm (Supplementary Information Section 4.4) using the observed and expected PTV counts within each gene to separate genes into three categories: null (observed \approx expected), recessive (observed $\leq 50\%$ of expected), and haploinsufficient (observed $< 10\%$ of expected). This metric – the probability of being loss-of-function (LoF) intolerant (pLI) – separates genes of sufficient length into LoF intolerant (pLI ≥ 0.9 , $n=3,230$) or LoF tolerant (pLI ≤ 0.1 , $n=10,374$) categories. pLI is less correlated with coding sequence length ($r = 0.17$ as compared to 0.57 for the PTV Z score), outperforms the PTV Z score as an intolerance metric (Supplementary Information Table 15), and reveals the expected contrast between gene lists (Figure 3b). pLI is positively correlated with a gene product's number of physical interaction partners ($p < 10^{-41}$). The most constrained pathways (highest median pLI for the genes in the pathway) are core biological processes (spliceosome, ribosome, and proteasome components; KS test $p < 10^{-6}$ for all) while olfactory receptors are among the least constrained pathways (KS test $p < 10^{-16}$), demonstrated in Figure 3b and consistent with previous work^{5,16–19}.

Critically, we note that LoF-intolerant genes include virtually all known severe haploinsufficient human disease genes (Figure 3b), but that 72% of LoF-intolerant genes have not yet been assigned a human disease phenotype despite clear evidence for extreme selective constraint (Supplementary Information Table 13). We note that this extreme constraint does not necessarily reflect a lethal disease or status as a disease gene (e.g. *BRCA1* has a pLI of 0), but is likely to point to genes where heterozygous loss of function confers some non-trivial survival or reproductive disadvantage.

The most highly constrained missense (top 25% missense Z scores) and PTV (pLI ≥ 0.9) genes show higher expression levels and broader tissue expression than the least constrained genes²⁰ (Figure 3c). These most highly constrained genes are also depleted for eQTLs ($p < 10^{-9}$ for missense and PTV; Figure 3d), yet are enriched within genome-wide significant trait-associated loci ($\chi^2 p < 10^{-14}$, Figure 3e). Intuitively, genes intolerant of PTV variation are dosage sensitive: natural selection does not tolerate a 50% deficit in expression due to the loss of single allele. Unsurprisingly, these genes are also depleted of common genetic variants that have a large enough effect on expression to be detected as eQTLs with current limited sample sizes. However, smaller changes in the expression of these genes, through weaker eQTLs or functional variants, are more likely to contribute to medically relevant phenotypes.

Finally, we investigated how these constraint metrics would stratify mutational classes according to their frequency spectrum, corrected for mutability as in the previous section (Figure 3f). The effect was most dramatic when considering nonsense variants in the LoF-

intolerant set of genes. For missense variants, the missense Z score offers information additional to Polyphen2 and CADD classifications, indicating that gene-level measures of constraint offer additional information to variant-level metrics in assessing potential pathogenicity.

ExAC improves variant interpretation in Mendelian disease

We assessed the value of ExAC as a reference dataset for clinical sequencing approaches, which typically prioritize or filter potentially deleterious variants based on functional consequence and allele frequency (AF)⁶. Filtering on ExAC reduced the number of candidate protein-altering variants by 7-fold compared to ESP, and was most powerful when the highest AF in any one population (“popmax”) was used rather than average (“global”) AF (Figure 4a). ESP is not well-powered to filter at 0.1% AF without removing many genuinely rare variants, as AF estimates based on low allele counts are both upward-biased and imprecise (Figure 4b). We thus expect that ExAC will provide a very substantial boost in the power and accuracy of variant filtering in Mendelian disease projects.

Previous large-scale sequencing studies have repeatedly shown that some purported Mendelian disease-causing genetic variants are implausibly common in the population^{21–23} (Figure 4c). The average ExAC participant harbors ~54 variants reported as disease-causing in two widely-used databases of disease-causing variants (Supplementary Information Section 5.2). Most (~41) of these are high-quality genotypes but with implausibly high (>1%) popmax AF. We therefore hypothesized that most of the supposed burden of Mendelian disease alleles per person is due not to genotyping error, but rather to misclassification in the literature and/or in databases.

We manually curated the evidence of pathogenicity for 192 previously reported pathogenic variants with AF >1% either globally or in South Asian or Latino individuals, populations that are underrepresented in previous reference databases. Nine variants had sufficient data to support disease association, typically with either mild or incompletely penetrant disease effects; the remainder either had insufficient evidence for pathogenicity, no claim of pathogenicity, or were benign traits (Supplementary Information Section 5.3). It is difficult to prove the absence of any disease association, and incomplete penetrance or genetic modifiers may contribute in some cases. Nonetheless, the high cumulative AF of these variants combined with their limited original evidence for pathogenicity suggest little contribution to disease, and 163 variants met American College of Medical Genetics criteria²⁴ for reclassification as benign or likely benign (Figure 4d). 126 of these 163 have been reclassified in source databases as of December 2015 (Supplementary Information Table 20). Supporting functional data were reported for 18 of these variants, highlighting the need to review cautiously even variants with experimental support.

We also sought phenotypic data for a subset of ExAC participants homozygous for reported severe recessive disease variants, again enabling reclassification of some variants as benign. North American Indian Childhood Cirrhosis is a recessive disease of cirrhotic liver failure during childhood requiring liver transplant for survival to adulthood, previously reported to be caused by *CIRH1A* p.R565W²⁵. ExAC contains 222 heterozygous and 4 homozygous

Latino individuals, with a population AF of 1.92%. The 4 homozygotes had no history of liver disease and recontact in two individuals revealed normal liver function (Supplementary Information Table 22). Thus, despite the rigorous linkage and Sanger sequencing efforts that led to the original report of pathogenicity, the ExAC data demonstrate that this variant is either benign or insufficient to cause disease, highlighting the importance of matched reference populations.

The above curation efforts confirm the importance of AF filtering in analysis of candidate disease variants^{6,26,27}. However, literature and database errors are prevalent even at lower AFs: the average ExAC individual contains 0.89 (<1% popmax AF) reportedly Mendelian variants in well-characterized dominant disease genes²⁸ and 0.21 at <0.1% popmax AF. This inflation likely results from a combination of false reports of pathogenicity and incomplete penetrance, as we have recently shown for *PRNP*²⁹. The abundance of rare functional variation in many disease genes in ExAC is a reminder that such variants should not be assumed to be causal or highly penetrant without careful segregation or case-control analysis^{7,24}.

Impact of rare protein-truncating variants

We investigated the distribution of PTVs, variants predicted to disrupt protein-coding genes through the introduction of a stop codon or frameshift or the disruption of an essential splice site; such variants are expected to be enriched for complete loss of function of the impacted genes. Naturally-occurring PTVs in humans provide a model for the functional impact of gene inactivation, and have been used to identify many genes in which LoF causes severe disease³⁰, as well as rare cases where LoF is protective against disease³¹.

Among the 7,404,909 HQ variants in ExAC, we found 179,774 high-confidence PTVs (as defined in Supplementary Information Section 6), 121,309 of which are singletons. This corresponds to an average of 85 heterozygous and 35 homozygous PTVs per individual (Figure 5a). The diverse nature of the cohort enables the discovery of substantial numbers of novel PTVs: out of 58,435 PTVs with an allele count greater than one, 33,625 occur in only one population. However, while PTVs as a category are extremely rare, the majority of the PTVs found in any one person are common, and each individual has only ~2 singleton PTVs, of which 0.14 are found in PTV-constrained genes (pLI >0.9). ExAC recapitulates known aspects of population demographic models, including an increase in intermediate-frequency (1–5%) PTVs in Finland³² and relatively common (>1%) PTVs in Africans (Figure 5b). However, these differences are diminished when considering only LoF-constrained (pLI > 0.9) genes (Extended Data Figure 4).

Using a sub-sampling approach, we show that the discovery of both heterozygous (Figure 5c) and homozygous (Figure 5d) PTVs scales very differently across human populations, with implications for the design of large-scale sequencing studies for the ascertainment of human “knockouts” described below.

Discussion

Here we describe the generation and analysis of the most comprehensive catalogue of human protein-coding genetic variation to date, incorporating high-quality exome sequencing data from 60,706 individuals of diverse geographic ancestry. The resulting call set provides unprecedented resolution for the analysis of low-frequency protein-coding variants in human populations, as well as a public resource [exac.broadinstitute.org] for the clinical interpretation of genetic variants observed in disease patients.

The very large sample size of ExAC also provides opportunities for a high-resolution analysis of the sensitivity of human genes to functional variation. While previous sample sizes have been adequately powered for the assessment of gene-level intolerance to missense variation^{11,14}, ExAC provides for the first time sufficient power to investigate genic intolerance to PTVs, highlighting 3,230 highly LoF-intolerant genes, 72% of which have no established human disease phenotype in OMIM or ClinVar. While this extreme depletion of PTVs is likely to highlight genes where loss of a single copy has been reproductively disadvantageous over recent human history, not all high pLI genes will lead to lethal disease. Additionally, disease genes—particularly those that act after post-reproductive age—do not necessarily have high pLI values (e.g. the pLI of BRCA1 is 0). In independent work [Ruderfer et al., manuscript submitted] we show that ExAC similarly provides power to identify genes intolerant of copy number variation. Quantification of genic intolerance to both classes of variation will provide added power to disease studies.

The ExAC resource provides the largest database to date for the estimation of allele frequency for protein-coding genetic variants, providing a powerful filter for analysis of candidate pathogenic variants in severe Mendelian diseases. Frequency data from ESP¹ have been widely used for this purpose, but those data are limited by population diversity and by resolution at allele frequencies $\leq 0.1\%$. ExAC therefore provides substantially improved power for Mendelian analyses, although it is still limited in power at lower allele frequencies, emphasizing the need for more sophisticated pathogenic variant filtering strategies alongside on-going data aggregation efforts.

Finally, we show that different populations confer different advantages in the discovery of gene-disrupting PTVs, providing guidance for the identification of human “knockouts” to understand gene function. Sampling multiple populations would likely be a fruitful strategy for a researcher investigating common PTV variation. However, discovery of homozygous PTVs is markedly enhanced in the South Asian samples, which come primarily from a Pakistani cohort with 38.3% of individuals self-reporting as having closely related parents, emphasizing the extreme value of consanguineous cohorts for “human knockout” discovery^{33–35} (Figure 5d). Other approaches to enriching for homozygosity of rare PTVs, such as focusing on bottlenecked populations, have already proved fruitful^{32,33}.

Even with this large collection of jointly processed exomes, many limitations remain. Firstly, most ExAC individuals were ascertained for biomedically important disease; while we have attempted to exclude severe pediatric diseases, the inclusion of both cases and controls for several polygenic disorders means that ExAC certainly contains disease-associated

variants³⁶. Secondly, future reference databases would benefit from including a broader sampling of human diversity, especially from under-represented Middle Eastern and African populations. Thirdly, the inclusion of whole genomes will also be critical to investigate additional classes of functional variation and identify non-coding constrained regions. Finally, and most critically, detailed phenotype data are unavailable for the vast majority of ExAC samples; future initiatives that assemble sequence and clinical data from very large-scale cohorts will be required to fully translate human genetic findings into biological and clinical understanding.

While the ExAC dataset exceeds the scale of previously available frequency reference datasets, much remains to be gained by further increases in sample size. Indeed, the fact that even the rarest transversions have mutational rates¹¹ on the order of 1×10^{-9} implies that the vast majority of possible non-lethal SNVs likely exist in some living human. ExAC already includes >63% of all possible protein-coding CpG transitions at well-covered synonymous sites; orders-of-magnitude increases in sample size will eventually lead to saturation of other classes of variation.

ExAC was made possible by the willingness of multiple large disease-focused consortia to share their raw data, and by the availability of the software and computational resources required to create a harmonized variant call set on the scale of tens of thousands of samples. The creation of yet larger reference variant databases will require continued emphasis on the value of genomic data sharing.

Online Methods

Variant discovery

We assembled approximately 1 petabyte of raw sequencing data (FASTQ files) from 91,796 individual exomes drawn from a wide range of primarily disease-focused consortia (Supplementary Information Table 2). We processed these exomes through a single informatic pipeline and performed joint variant calling of single nucleotide variants (SNVs) and short insertions and deletions (indels) across all samples using a new version of the Genome Analysis Toolkit (GATK) HaplotypeCaller pipeline. Variant discovery was performed within a defined exome region that includes Gencode v19 coding regions and flanking 50 bases. At each site, sequence information from all individuals was used to assess the evidence for the presence of a variant in each individual. Full details of data processing, variant calling and resources are described in the Supplementary Information Sections 1.1–1.4.

Quality assessment

We leveraged a variety of sources of internal and external validation data to calibrate filters and evaluate the quality of filtered variants (Supplementary Information Table 7). We adjusted the standard GATK variant site filtering³⁷ to increase the number of singleton variants that pass this filter, while maintaining a singleton transmission rate of 50.1%, very near the expected 50%, within sequenced trios. We then used the remaining passing variants to assess depth and genotype quality filters compared to >10,000 samples that had been

directly genotyped using SNP arrays (Illumina HumanExome) and achieved 97–99% heterozygous concordance, consistent with known error rates for rare variants in chip-based genotyping³⁸. Relative to a “platinum standard” genome sequenced using five different technologies³⁹, we achieved sensitivity of 99.8% and false discovery rates (FDR) of 0.056% for single nucleotide variants (SNVs), and corresponding rates of 95.1% and 2.17% for insertions and deletions (indels). Lastly, we compared 13 representative Non-Finnish European exomes included in the call set with their corresponding 30x PCR-Free genome. The overall SNV and indel FDR was 0.14% and 4.71%, while for SNV singletons was 0.389%. The overall FDR by annotation classes missense, synonymous and protein truncating variants (including indels) were 0.076%, 0.055% and 0.471% respectively (Supplementary Information Table 5 and 6). Full details of quality assessments are described in the Supplementary Information Section 1.6.

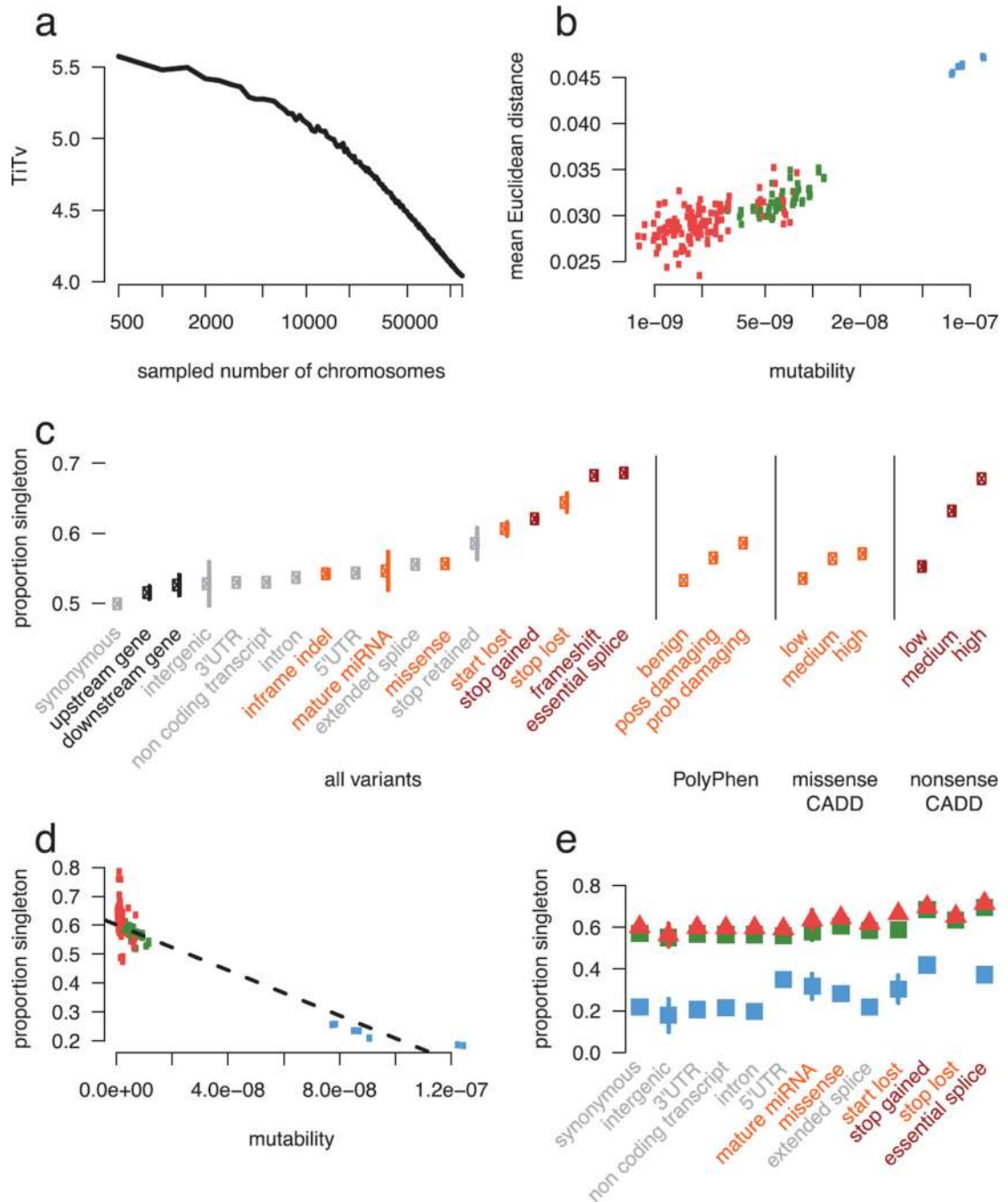
Sample filtering

The 91,796 samples were filtered based on two criteria. First, samples that were outliers for key metrics were removed (Extended Data Figure 5b). Second, in order to generate allele frequencies based on independent observations without enrichment of Mendelian disease alleles, we restricted the final release data set to unrelated adults with high-quality sequence data and without severe pediatric disease. After filtering, only 60,706 samples remained, consisting of ~77% of Agilent (33 Mb target) and ~12% of Illumina (37.7 Mb target) exome captures. Full details of the filtering process are described in the Supplementary Information Section 1.7.

ExAC data release

For each variant, summary data for genotype quality, allele depth and population specific allele counts were calculated before removing all genotype data. This variant summary file was then functionally annotated using variant effect predictor (VEP) with the LOFTEE plugin. This data set can be accessed via the ExAC Browser (<http://exac.broadinstitute.org>) or downloaded from ftp://ftp.broadinstitute.org/pub/ExAC_release/release0.3/ExAC.r0.3.sites.vep.vcf.gz. Full details regarding the annotation of the ExAC data set are described in the Supplementary Information Sections 1.9–1.10.

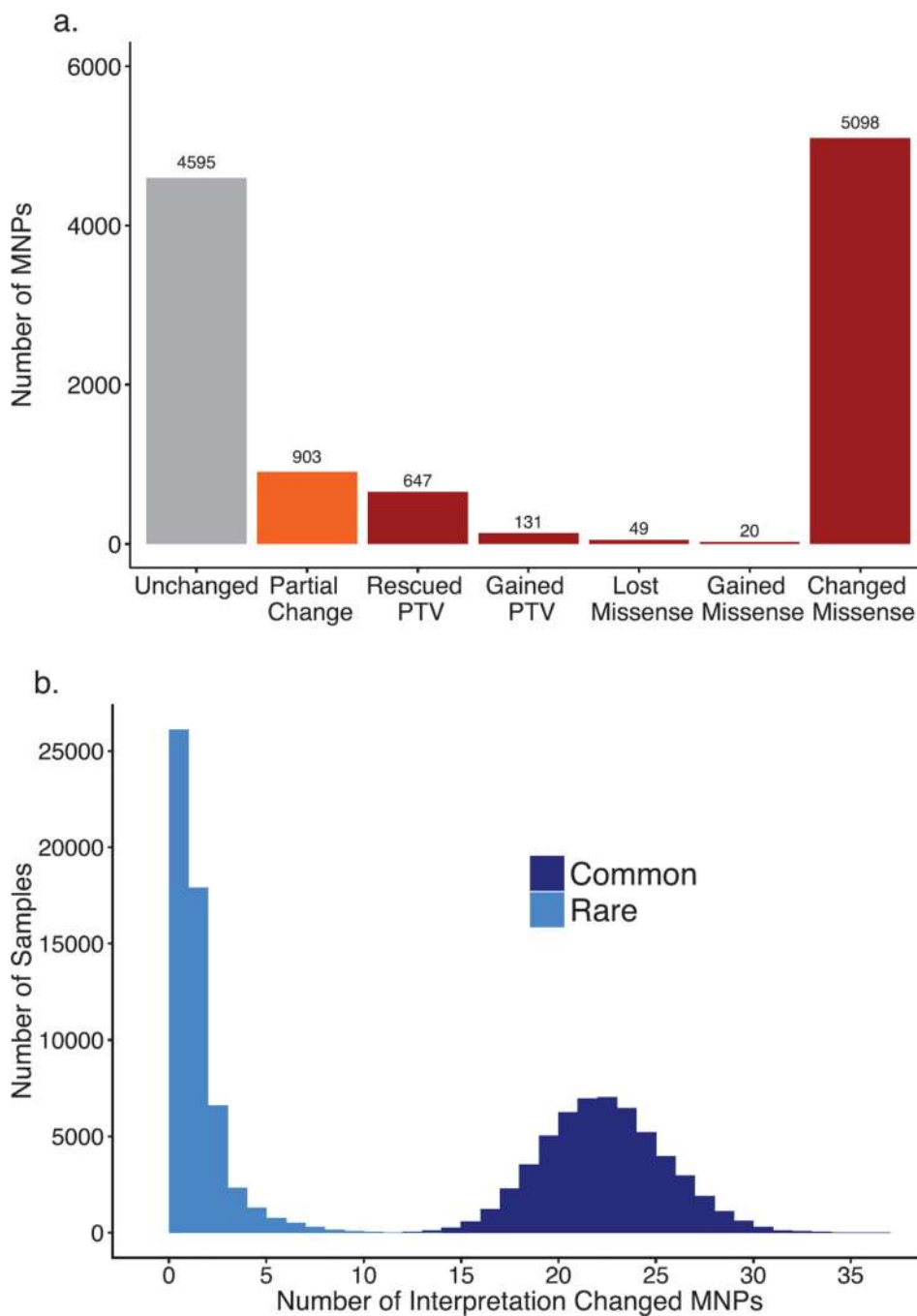
Extended Data



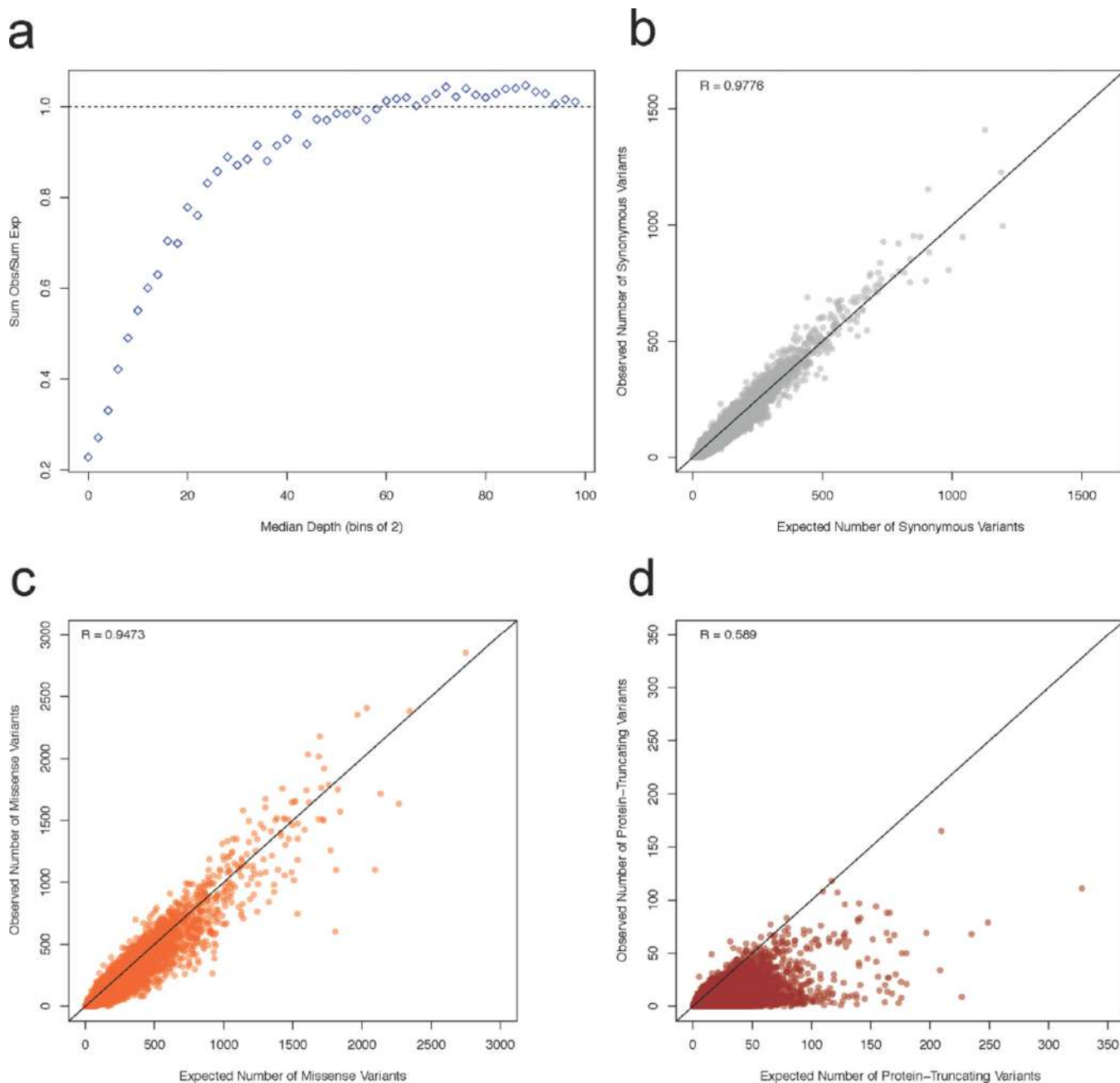
Extended Data Figure 1. The impact of recurrence across different mutation and functional classes

a) TiTv (Transition to transversion) ratio of synonymous variants at downsampled intervals of ExAC. The TiTv is relatively stable at previous sample sizes (<5000) but changes drastically at larger sample sizes. b) For synonymous doubleton variants, mutability of each trinucleotide context is correlated with mean Euclidean distance of individuals that share the doubleton. Transversion (red) and non-CpG transition (green) doubletons are more likely to

be found in closer PCA space (i.e. more similar ethnicities) than CpG transitions (blue) c) The proportion singleton among various functional categories. The functional category stop lost has a higher singleton rate than nonsense. Error bars represent standard error of the mean. d) Among synonymous variants, mutability of each trinucleotide context is correlated with proportion singleton, suggesting CpG transitions (blue) are more likely to have multiple independent origins driving their allele frequency up. e) The proportion singleton metric from c) broken down by transversions, non-CpG transitions, and CpG variants. Notably, there is a wide variation in singleton rates among mutational contexts in functional classes, and there are no stop-lost CpG transitions. Error bars represent standard error of the mean.

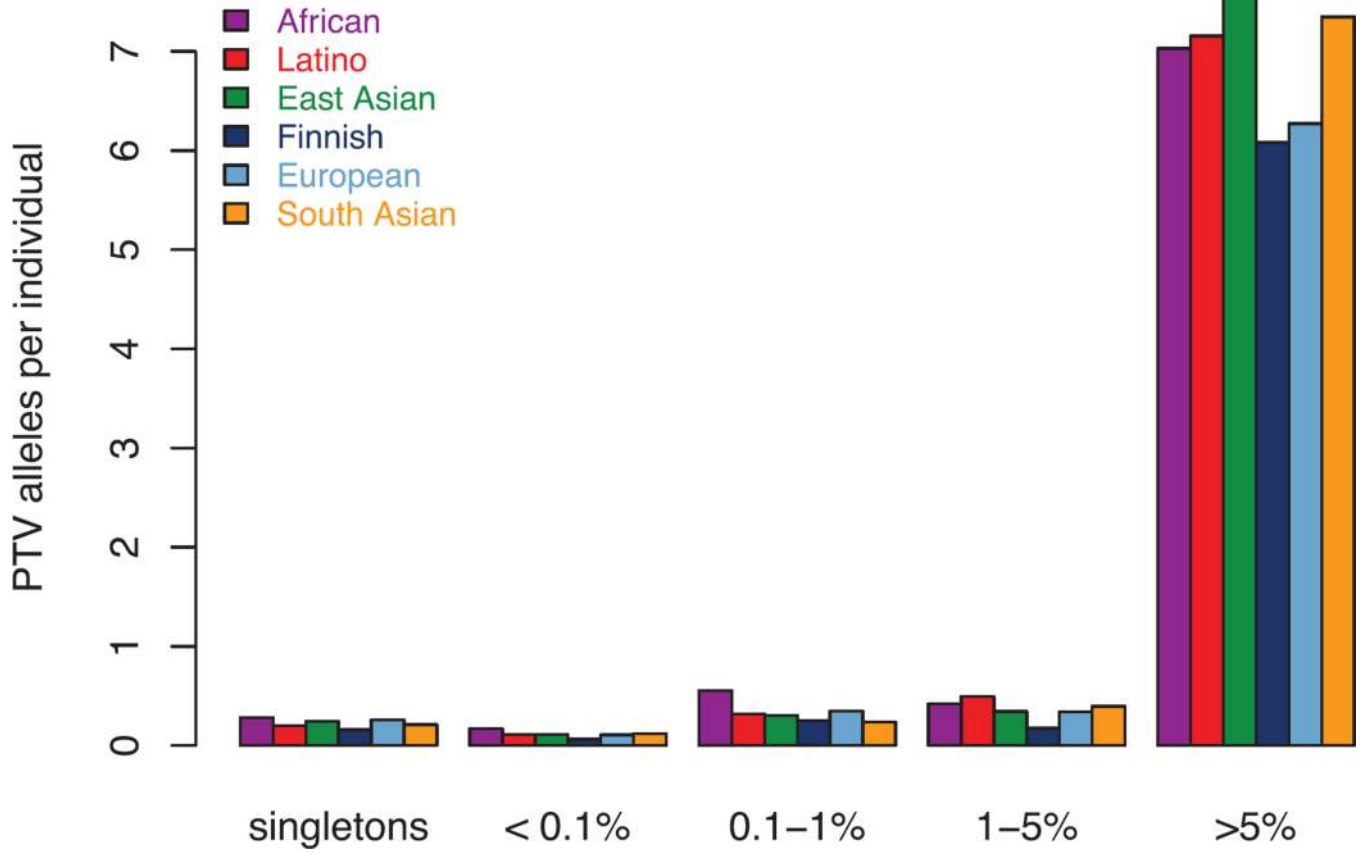


Extended Data Figure 2. Multi-nucleotide variants discovered in the ExAC data set
 a) Number of MNPs per impact on the variant interpretation. b) Distribution of the number of MNPs per sample where phasing changes interpretation, separated by allele frequency. Common > 1%, Rare < 1%. MNPs comprised of a rare and common allele are considered rare as this defines the frequency of the MNP.



Extended Data Figure 3. Relationships between depth and observed vs expected variants as well as correlations between observed and expected variant counts for synonymous, missense, and protein-truncating

a) The relationship between the median depth of exons (bins of 2) and the sum of all observed synonymous variants in those exons divided by the sum of all expected synonymous variants. The curve was used to determine the appropriate depth adjustment for expected variant counts. For the rest of the panels, the correlation between the depth-adjusted expected variants counts and observed are depicted for synonymous (b), missense (c), and protein-truncating (d). The black line indicates a perfect correlation (slope = 1). Axes have been trimmed to remove *TTN*.



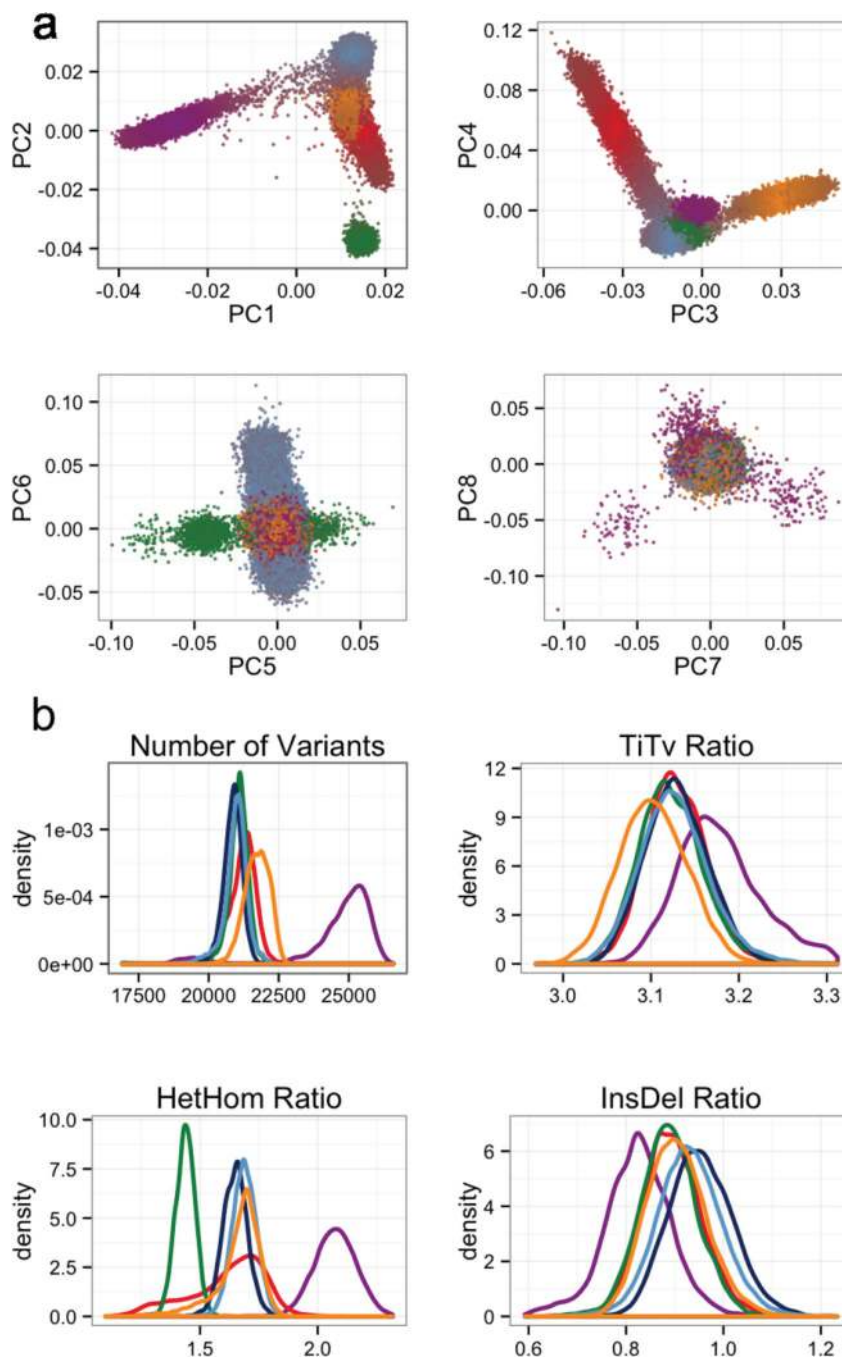
Extended Data Figure 4. Number of protein-truncating variants in constrained genes per individual by allele frequency bin
 Equivalent to Figure 5b limited to constrained (pLI ≥ 0.9) genes.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



Extended Data Figure 5. Principal component analysis (PCA) and key metrics used to filter samples

a) Principal component analysis using a set of 5,400 common exome SNPs. Individuals are colored by their distance from each of the population cluster centers using the first 4 principal components. b) The metrics number of variants, TiTv, alternate heterozygous/homozygous (HetHom) ratio and Insertion/Deletion (InsDel) ratio. Populations are their respective colors: Latino (red), African (purple), European (blue), South Asian (yellow) and East Asian (green).

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Authors

Exome Aggregation Consortium[#], Monkol Lek^{1,2,3,4}, Konrad J Karczewski^{1,2,*}, Eric V Minikel^{1,2,5,*}, Kaitlin E Samocha^{1,2,6,5,*}, Eric Banks², Timothy Fennell², Anne H O'Donnell-Luria^{1,2,7}, James S Ware^{2,8,9,10,11}, Andrew J Hill^{1,2,12}, Beryl B Cummings^{1,2,5}, Taru Tukiainen^{1,2}, Daniel P Birnbaum², Jack A Kosmicki^{1,2,6,13}, Laramie E Duncan^{1,2,6}, Karol Estrada^{1,2}, Fengmei Zhao^{1,2}, James Zou², Emma Pierce-Hoffman^{1,2}, Joanne Berghout^{14,15}, David N Cooper¹⁶, Nicole Deflaux¹⁷, Mark DePristo¹⁸, Ron Do^{19,20,21,22}, Jason Flannick^{2,23}, Menachem Fromer^{1,6,24,19,20}, Laura Gauthier¹⁸, Jackie Goldstein^{1,2,6}, Namrata Gupta², Daniel Howrigan^{1,2,6}, Adam Kiezun¹⁸, Mitja I Kurki^{2,25}, Ami Levy Moonshine¹⁸, Pradeep Natarajan^{2,26,27,28}, Lorena Orozco²⁹, Gina M Peloso^{2,27,28}, Ryan Poplin¹⁸, Manuel A Rivas², Valentin Ruano-Rubio¹⁸, Samuel A Rose⁶, Douglas M Ruderfer^{24,19,20}, Khalid Shakir¹⁸, Peter D Stenson¹⁶, Christine Stevens², Brett P Thomas^{1,2}, Grace Tiao¹⁸, Maria T Tusie-Luna³⁰, Ben Weisburd², Hong-Hee Won³¹, Dongmei Yu^{6,27,25,32}, David M Altshuler^{2,33}, Diego Ardisino³⁴, Michael Boehnke³⁵, John Danesh³⁶, Stacey Donnelly², Roberto Elosua³⁷, Jose C Florez^{2,26,27}, Stacey B Gabriel², Gad Getz^{18,26,38}, Stephen J Glatt^{39,40,41}, Christina M Hultman⁴², Sekar Kathiresan^{2,26,27,28}, Markku Laakso⁴³, Steven McCarroll^{6,8}, Mark I McCarthy^{44,45,46}, Dermot McGovern⁴⁷, Ruth McPherson⁴⁸, Benjamin M Neale^{1,2,6}, Aarno Palotie^{1,2,5,49}, Shaun M Purcell^{24,19,20}, Danish Saleheen^{50,51,52}, Jeremiah M Scharf^{2,6,27,25,32}, Pamela Sklar^{24,19,20,53,54}, Patrick F Sullivan^{55,56}, Jaakko Tuomilehto⁵⁷, Ming T Tsuang⁵⁸, Hugh C Watkins^{59,44}, James G Wilson⁶⁰, Mark J Daly^{1,2,6}, and Daniel G MacArthur^{1,2}

Affiliations

¹Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA ²Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA ³School of Paediatrics and Child Health, University of Sydney, Sydney, NSW, Australia ⁴Institute for Neuroscience and Muscle Research, Childrens Hospital at Westmead, Sydney, NSW, Australia ⁵Program in Biological and Biomedical Sciences, Harvard Medical School, Boston, MA, USA ⁶Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA, USA ⁷Division of Genetics and Genomics, Boston Children's Hospital, Boston, MA, USA ⁸Department of Genetics, Harvard Medical School, Boston, MA, USA ⁹National Heart and Lung Institute, Imperial College London, London, UK ¹⁰NIHR Royal Brompton Cardiovascular Biomedical Research Unit, Royal Brompton Hospital, London, UK ¹¹MRC Clinical Sciences Centre, Imperial College London, London, UK ¹²Genome Sciences, University of Washington, Seattle, WA, USA ¹³Program in Bioinformatics and Integrative Genomics, Harvard Medical School, Boston, MA, USA ¹⁴Mouse Genome Informatics, Jackson Laboratory, Bar Harbor, ME, USA ¹⁵Center for Biomedical Informatics and

Biostatistics, University of Arizona, Tucson, AZ, USA ¹⁶Institute of Medical Genetics, Cardiff University, Cardiff, UK ¹⁷Google Inc, Mountain View, CA, USA ¹⁸Broad Institute of MIT and Harvard, Cambridge, MA, USA ¹⁹Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA ²⁰Institute for Genomics and Multiscale Biology, Icahn School of Medicine at Mount Sinai, New York, NY, USA ²¹The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA ²²The Center for Statistical Genetics, Icahn School of Medicine at Mount Sinai, New York, NY, USA ²³Department of Molecular Biology, Massachusetts General Hospital, Boston, MA, USA ²⁴Department of Psychiatry, Icahn School of Medicine at Mount Sinai, New York, NY, USA ²⁵Psychiatric and Neurodevelopmental Genetics Unit, Massachusetts General Hospital, Boston, MA, USA ²⁶Harvard Medical School, Boston, MA, USA ²⁷Center for Human Genetic Research, Massachusetts General Hospital, Boston, MA, USA ²⁸Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA ²⁹Immunogenomics and Metabolic Disease Laboratory, Instituto Nacional de Medicina Genómica, Mexico City, Mexico ³⁰Molecular Biology and Genomic Medicine Unit, Instituto Nacional de Ciencias Médicas y Nutrición, Mexico City, Mexico ³¹Samsung Advanced Institute for Health Sciences and Technology (SAIHST), Sungkyunkwan University, Samsung Medical Center, Seoul, Republic of Korea ³²Department of Neurology, Massachusetts General Hospital, Boston, MA, USA ³³Vertex Pharmaceuticals, Boston, MA, USA ³⁴Department of Cardiology, University Hospital, Parma, Italy ³⁵Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA ³⁶Department of Public Health and Primary Care, Strangeways Research Laboratory, Cambridge, UK ³⁷Cardiovascular Epidemiology and Genetics, Hospital del Mar Medical Research Institute, Barcelona, Spain ³⁸Department of Pathology and Cancer Center, Massachusetts General Hospital, Boston, MA, USA ³⁹Psychiatric Genetic Epidemiology & Neurobiology Laboratory, State University of New York, Upstate Medical University, Syracuse, NY, USA ⁴⁰Department of Psychiatry and Behavioral Sciences, State University of New York, Upstate Medical University, Syracuse, NY, USA ⁴¹Department of Neuroscience and Physiology, State University of New York, Upstate Medical University, Syracuse, NY, USA ⁴²Department of Medical Epidemiology and Biostatistics, Karolinska Institute, Stockholm, Sweden ⁴³Department of Medicine, University of Eastern Finland and Kuopio University Hospital, Kuopio, Finland ⁴⁴Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford, UK ⁴⁵Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Oxford, UK ⁴⁶Oxford NIHR Biomedical Research Centre, Oxford University Hospitals Foundation Trust, Oxford, UK ⁴⁷Inflammatory Bowel Disease and Immunobiology Research Institute, Cedars-Sinai Medical Center, Los Angeles, CA, USA ⁴⁸Atherogenomics Laboratory, University of Ottawa Heart Institute, Ottawa, ON, Canada ⁴⁹Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland ⁵⁰Department of Biostatistics and Epidemiology, Perelman School of Medicine at the University of Pennsylvania, Philadelphia, PA, USA ⁵¹Department of Medicine, Perelman School of

Medicine at the University of Pennsylvania, Philadelphia, PA, USA ⁵²Center for Non-Communicable Diseases, Karachi, , Pakistan ⁵³Friedman Brain Institute, Icahn School of Medicine at Mount Sinai, New York, NY, USA ⁵⁴Department of Neuroscience, Icahn School of Medicine at Mount Sinai, New York, NY, USA ⁵⁵Department of Genetics, University of North Carolina, Chapel Hill, NC, USA ⁵⁶Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden ⁵⁷Department of Public Health, University of Helsinki, Helsinki, Finland ⁵⁸Department of Psychiatry, University of California, San Diego, CA, USA ⁵⁹Radcliffe Department of Medicine, University of Oxford, Oxford, UK ⁶⁰Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, MS, USA

Acknowledgments

We would like to thank the reviewers and editor for their time, valuable comments and suggestions. The scientific community for their support and comments on biorxiv, twitter and other public forums. Brendan Bulik-Sullivan, Jon Bloom and Raymond Walters for their help with mathematical notation. The full acknowledgements are detailed in Supplementary Information Section 8.

Collaborators (alphabetical order)

Hanna E Abboud⁶¹, Goncalo Abecasis³⁵, Carlos A Aguilar-Salinas⁶², Olimpia Arellano-Campos⁶², Gil Atzmon^{63,64}, Ingvild Aukrust^{65,66,67}, Cathy L Barr^{68,69}, Graeme I Bell⁷⁰, Graeme I Bell^{70,71}, Sarah Bergen⁴², Lise Bjørkhaug^{66,67}, John Blangero^{72,73}, Donald W Bowden^{74,75,76}, Cathy L Budman⁷⁷, Noël P Burt², Federico Centeno-Cruz⁷⁸, John C Chambers^{79,80,81}, Kimberly Chambert⁶, Robert Clarke⁸², Rory Collins⁸², Giovanni Coppola⁸³, Emilio J Córdova⁷⁸, Maria L Cortes¹⁸, Nancy J Cox⁸⁴, Ravindranath Duggirala⁸⁵, Martin Farrall^{59,44}, Juan C Fernandez-Lopez⁷⁸, Pierre Fontanillas², Timothy M Frayling⁸⁶, Nelson B Freimer⁸³, Christian Fuchsberger³⁵, Humberto García-Ortiz⁷⁸, Anuj Goel^{59,44}, María J Gómez-Vázquez⁶², María E González-Villalpando⁸⁷, Clicerio González-Villalpando⁸⁷, Marco A Grados⁸⁸, Leif Groop⁸⁹, Christopher A Haiman⁹⁰, Craig L Hanis⁹¹, Craig L Hanis⁹¹, Andrew T Hattersley⁸⁶, Brian E Henderson⁹², Jemma C Hopewell⁸², Alicia Huerta-Chagoya⁹³, Sergio Islas-Andrade⁹⁴, Suzanne BR Jacobs², Shapour Jalilzadeh^{59,44}, Christopher P Jenkinson⁶¹, Jennifer Moran², Silvia Jiménez-Morale⁷⁸, Anna Kähler⁴², Robert A King⁹⁵, George Kirov⁹⁶, Jaspal S Kooner^{80,9,81}, Theodosios Kyriakou^{59,44}, Jong-Young Lee⁹⁷, Donna M Lehman⁶¹, Gholson Lyon⁹⁸, William MacMahon⁹⁹, Patrik KE Magnusson⁴², Anubha Mahajan¹⁰⁰, Jaume Marrugat³⁷, Angélica Martínez-Hernández⁷⁸, Carol A Mathews¹⁰¹, Gilean McVean¹⁰⁰, James B Meigs^{102,26}, Thomas Meitinger^{103,104}, Elvia Mendoza-Caamal⁷⁸, Josep M Mercader^{2,105,106}, Karen L Mohlke⁵⁵, Hortensia Moreno-Macías¹⁰⁷, Andrew P Morris^{108,100,109}, Laeya A Najmi^{65,110}, Pål R Njølstad^{65,66}, Michael C O'Donovan⁹⁶, Maria L Ordóñez-Sánchez⁶², Michael J Owen⁹⁶, Taesung Park^{111,112}, David L Pauls²⁵, Danielle Posthuma^{113,114,115}, Cristina Revilla-Monsalve⁹⁴, Laura Riba⁹³, Stephan Ripke⁶, Rosario Rodríguez-Guillén⁶², Maribel Rodríguez-Torres⁶², Paul Sandor^{116,68}, Mark Seielstad^{117,118}, Rob Sladek^{119,120,121}, Xavier Soberón⁷⁸, Timothy D Spector¹²², Shyong E Tai^{123,124,125}, Tanya M Teslovich³⁵, Geoffrey Walford^{105,26}, Lynne R Wilkens⁹², Amy L Williams^{2,126}

- ⁶¹Department of Medicine, University of Texas Health Science Center, San Antonio, TX, USA
- ⁶²Instituto Nacional de Ciencias Médicas y Nutrición Salvador Zubirán, Mexico City, Mexico
- ⁶³Departments of Medicine and Genetics, Albert Einstein College of Medicine, New York City, NY, USA
- ⁶⁴Department of Natural Science, University of Haifa, Haifa, Israel
- ⁶⁵Department of Clinical Science, University of Bergen, Bergen, Norway
- ⁶⁶Department of Pediatrics, Haukeland University Hospital, Bergen, Norway
- ⁶⁷Department of Biomedicine, University of Bergen, Bergen, Norway
- ⁶⁸The Toronto Western Research Institute, University Health Network, Toronto, Canada
- ⁶⁹The Hospital for Sick Children, Toronto, Canada
- ⁷⁰Departments of Medicine and Human Genetics, University of Chicago, Chicago, IL, USA
- ⁷¹Department of Medicine, University of Chicago, Chicago, IL, USA
- ⁷²South Texas Diabetes and Obesity Institute, University of Texas Health Science Center, San Antonio, TX, USA
- ⁷³University of Texas Rio Grande Valley, Brownsville, TX, USA
- ⁷⁴Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ⁷⁵Center for Genomics and Personalized Medicine Research, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ⁷⁶Center for Diabetes Research, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ⁷⁷North Shore-Long Island Jewish Health System, Manhasset, NY, USA
- ⁷⁸Instituto Nacional de Medicina Genómica, Mexico City, Mexico
- ⁷⁹Department of Epidemiology and Biostatistics, Imperial College London, London, UK
- ⁸⁰Department of Cardiology, Ealing Hospital NHS Trust, Southall, UK
- ⁸¹Imperial College Healthcare NHS Trust, Imperial College London, London, UK
- ⁸²Nuffield Department of Population Health, University of Oxford, Oxford, UK
- ⁸³Center for Neurobehavioral Genetics, University of California, Los Angeles, CA, USA

- ⁸⁴Vanderbilt Genetics Institute, Vanderbilt University School of Medicine, Nashville, TN, USA
- ⁸⁵Department of Genetics, Texas Biomedical Research Institute, San Antonio, TX, USA
- ⁸⁶University of Exeter Medical School, University of Exeter, Exeter, UK
- ⁸⁷Instituto Nacional de Salud Publica, Mexico City, Mexico
- ⁸⁸Department of Psychiatry and Behavioral Sciences, Johns Hopkins University School of Medicine, Baltimore, MD, USA
- ⁸⁹Department of Clinical Sciences, Lund University Diabetes Centre, Malm_, Sweden
- ⁹⁰Department of Preventive Medicine, University of Southern California, Los Angeles, CA, USA
- ⁹¹Human Genetics Center, The University of Texas Health Science Center, Houston, TX, USA
- ⁹²Epidemiology Program, University of Hawaii Cancer Center, Honolulu, HI, USA
- ⁹³Instituto de Investigaciones Biomédicas, Mexico City, Mexico
- ⁹⁴Instituto Mexicano del Seguro Social, Mexico City, Mexico
- ⁹⁵Department of Genetics, Yale University School of Medicine, New Haven, CT, USA
- ⁹⁶MRC Centre for Neuropsychiatric Genetics and Genomics, Cardiff University, Cardiff, UK
- ⁹⁷Center for Genome Science, Korea National Institute of Health, Chungcheongbuk-do, Republic of Korea
- ⁹⁸Stanley Institute for Cognitive Genomics, Cold Spring Harbor Laboratory, Woodbury, NY, USA
- ⁹⁹Department of Psychiatry, University of Utah, Salt Lake City, UT, USA
- ¹⁰⁰Nuffield Department of Medicine, University of Oxford, Oxford, UK
- ¹⁰¹Department of Psychiatry, University of Florida, Gainesville, FL, USA
- ¹⁰²General Medicine Division, Massachusetts General Hospital, Boston, MA, USA
- ¹⁰³Institute of Human Genetics, Technische Universität München, Munich, Germany
- ¹⁰⁴Institute of Human Genetics, German Research Center for Environmental Health, Neuherberg, Germany

- ¹⁰⁵Diabetes Research Center (Diabetes Unit), Massachusetts General Hospital, Boston, MA, USA
- ¹⁰⁶Research Program in Computational Biology, Barcelona Supercomputing Center, Barcelona, Spain
- ¹⁰⁷Universidad Autónoma Metropolitana, Mexico City, Mexico
- ¹⁰⁸Estonian Genome Centre, University of Tartu, Tartu, Estonia, University of Tartu, Tartu, Estonia
- ¹⁰⁹Department of Biostatistics, University of Liverpool, Liverpool, UK
- ¹¹⁰Center for Medical Genetics and Molecular Medicine, Haukeland University Hospital, Bergen, Norway
- ¹¹¹Interdisciplinary Program in Bioinformatics, Seoul National University, Seoul, Republic of Korea
- ¹¹²Department of Statistics, Seoul National University, Seoul, Republic of Korea
- ¹¹³Department of Functional Genomics, University of Amsterdam, Amsterdam, The Netherlands
- ¹¹⁴Department of Clinical Genetics, VU Medical Centre, Amsterdam, The Netherlands
- ¹¹⁵Department of Child and Adolescent Psychiatry, Erasmus University Medical Centre, Rotterdam, The Netherlands
- ¹¹⁶Department of Psychiatry, University of Toronto, Toronto, Canada
- ¹¹⁷Department of Laboratory Medicine, University of California, San Francisco, CA, USA
- ¹¹⁸Blood Systems Research Institute, San Francisco, CA, USA
- ¹¹⁹Department of Human Genetics, McGill University, Montreal, Canada
- ¹²⁰Department of Medicine, McGill University, Montreal, Canada
- ¹²¹McGill University and Genome Quebec Innovation Centre, Montreal, Canada
- ¹²²Department of Twin Research and Genetic Epidemiology, King's College London, London, UK
- ¹²³Saw Swee Hock School of Public Health, National University of Singapore, Singapore, Singapore
- ¹²⁴Department of Medicine, National University of Singapore, Singapore, Singapore
- ¹²⁵Cardiovascular & Metabolic Disorders Program, Duke-NUS Graduate Medical School Singapore, Singapore, Singapore

¹²⁶Department of Biological Sciences, Columbia University, New York, NY, USA

References

1. Fu W, et al. Analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants. *Nature*. 2013; 493:216–220. [PubMed: 23201682]
2. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature*. 2015; 526:68–74. [PubMed: 26432245]
3. Li H, Durbin R. Inference of human population history from individual whole-genome sequences. *Nature*. 2011; 475:493–496. [PubMed: 21753753]
4. Stoneking M, Krause J. Learning about human population history from ancient and modern genomes. *Nat. Rev. Genet.* 2011; 12:603–614. [PubMed: 21850041]
5. MacArthur DG, et al. A systematic survey of loss-of-function variants in human protein-coding genes. *Science*. 2012; 335:823–828. [PubMed: 22344438]
6. Bamshad MJ, et al. Exome sequencing as a tool for Mendelian disease gene discovery. *Nat. Rev. Genet.* 2011; 12:745–755. [PubMed: 21946919]
7. MacArthur DG, et al. Guidelines for investigating causality of sequence variants in human disease. *Nature*. 2014; 508:469–476. [PubMed: 24759409]
8. Deciphering Developmental Disorders Study. Large-scale discovery of novel genetic causes of developmental disorders. *Nature*. 2015; 519:223–228. [PubMed: 25533962]
9. Fromer M, et al. De novo mutations in schizophrenia implicate synaptic networks. *Nature*. 2014; 506:179–184. [PubMed: 24463507]
10. Cooper DN, Youssoufian H. The CpG dinucleotide and human genetic disease. *Hum. Genet.* 1988; 78:151–155. [PubMed: 3338800]
11. Samocha KE, et al. A framework for the interpretation of de novo mutation in human disease. *Nat. Genet.* 2014
12. Tennessen, Ja, et al. Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science*. 2012; 337:64–69. [PubMed: 22604720]
13. Gudbjartsson DF, et al. Large-scale whole-genome sequencing of the Icelandic population. *Nat. Genet.* 2015; 47:435–444. [PubMed: 25807286]
14. Petrovski S, Wang Q, Heinzen EL, Allen AS, Goldstein DB. Genic intolerance to functional variation and the interpretation of personal genomes. *PLoS Genet.* 2013; 9:e1003709. [PubMed: 23990802]
15. Vicoso B, Charlesworth B. Evolution on the X chromosome: unusual patterns and processes. *Nat. Rev. Genet.* 2006; 7:645–653. [PubMed: 16847464]
16. Jeong H, Mason SP, Barabási, a L, Oltvai ZN. Lethality and centrality in protein networks. *Nature*. 2001; 411:41–42. [PubMed: 11333967]
17. Goh K-I, et al. The human disease network. *Proc. Natl. Acad. Sci. U. S. A.* 2007; 104:8685–8690. [PubMed: 17502601]
18. Rolland T, et al. Resource A Proteome-Scale Map of the Human Interactome Network. *Cell*. 2014; 159:1212–1226. [PubMed: 25416956]
19. Itan Y, et al. The human gene damage index as a gene-level approach to prioritizing exome variants. *Proc. Natl. Acad. Sci. U. S. A.* 2015; 112:13615–13620. [PubMed: 26483451]
20. GTEx Consortium. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science*. 2015; 348:648–660. [PubMed: 25954001]
21. Bell CJ, et al. Carrier testing for severe childhood recessive diseases by next-generation sequencing. *Sci. Transl. Med.* 2011; 3:65ra4.
22. Xue Y, et al. Deleterious- and disease-allele prevalence in healthy individuals: Insights from current predictions, mutation databases, and population-scale resequencing. *Am. J. Hum. Genet.* 2012; 91:1022–1032. [PubMed: 23217326]
23. Piton A, Redin C, Mandel J-L. XLID-Causing Mutations and Associated Genes Challenged in Light of Data From Large-Scale Human Exome Sequencing. *Am. J. Hum. Genet.* 2013; 93:368–383. [PubMed: 23871722]

24. Richards S, et al. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* 2015; 17:405–423. [PubMed: 25741868]
25. Chagnon P, et al. A missense mutation (R565W) in cirhin (FLJ14728) in North American Indian childhood cirrhosis. *Am. J. Hum. Genet.* 2002; 71:1443–1449. [PubMed: 12417987]
26. Stenson PD, et al. The Human Gene Mutation Database: Building a comprehensive mutation repository for clinical and molecular genetics, diagnostic testing and personalized genomic medicine. *Hum. Genet.* 2014; 133:1–9. [PubMed: 24077912]
27. Dewey FE, et al. Sequence to Medical Phenotypes: A Framework for Interpretation of Human Whole Genome DNA Sequence Data. *PLOS Genet.* 2015; 11:e1005496. [PubMed: 26448358]
28. Blekhman R, et al. Natural Selection on Genes that Underlie Human Disease Susceptibility. *Curr. Biol.* 2008; 18:883–889. [PubMed: 18571414]
29. Minikel EV, et al. Quantifying prion disease penetrance using large population control cohorts. *Sci. Transl. Med.* 2016; 8:322ra9–322ra9.
30. Chong JX, et al. The Genetic Basis of Mendelian Phenotypes: Discoveries, Challenges, and Opportunities. *Am. J. Hum. Genet.* 2015:1–17.
31. Kathiresan S. Developing Medicines That Mimic the Natural Successes of the Human Genome. *J. Am. Coll. Cardiol.* 2015; 65:1562–1566. [PubMed: 25881938]
32. Lim ET, et al. Distribution and Medical Impact of Loss-of-Function Variants in the Finnish Founder Population. *PLoS Genet.* 2014; 10:e1004494. [PubMed: 25078778]
33. Sulem P, et al. Identification of a large set of rare complete human knockouts. *Nat. Genet.* 2015; 47:448–452. [PubMed: 25807282]
34. Narasimhan VM, et al. Health and population effects of rare gene knockouts in adult humans with related parents. *Science (80-.)*. 2016; 8624:1–8.
35. Saleheen D, et al. Human knockouts in a cohort with a high rate of consanguinity. *bioRxiv*. 2015
36. Freischmidt A, et al. Haploinsufficiency of TBK1 causes familial ALS and fronto-temporal dementia. *Nat. Neurosci.* 2015; 18
37. DePristo, Ma, et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* 2011; 43:491–498. [PubMed: 21478889]
38. Voight BF, et al. The MetaboChip, a Custom Genotyping Array for Genetic Studies of Metabolic, Cardiovascular, and Anthropometric Traits. *PLoS Genet.* 2012; 8:e1002793. [PubMed: 22876189]
39. Zook JM, et al. Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. *Nat. Biotechnol.* 2014; 32:246–251. [PubMed: 24531798]

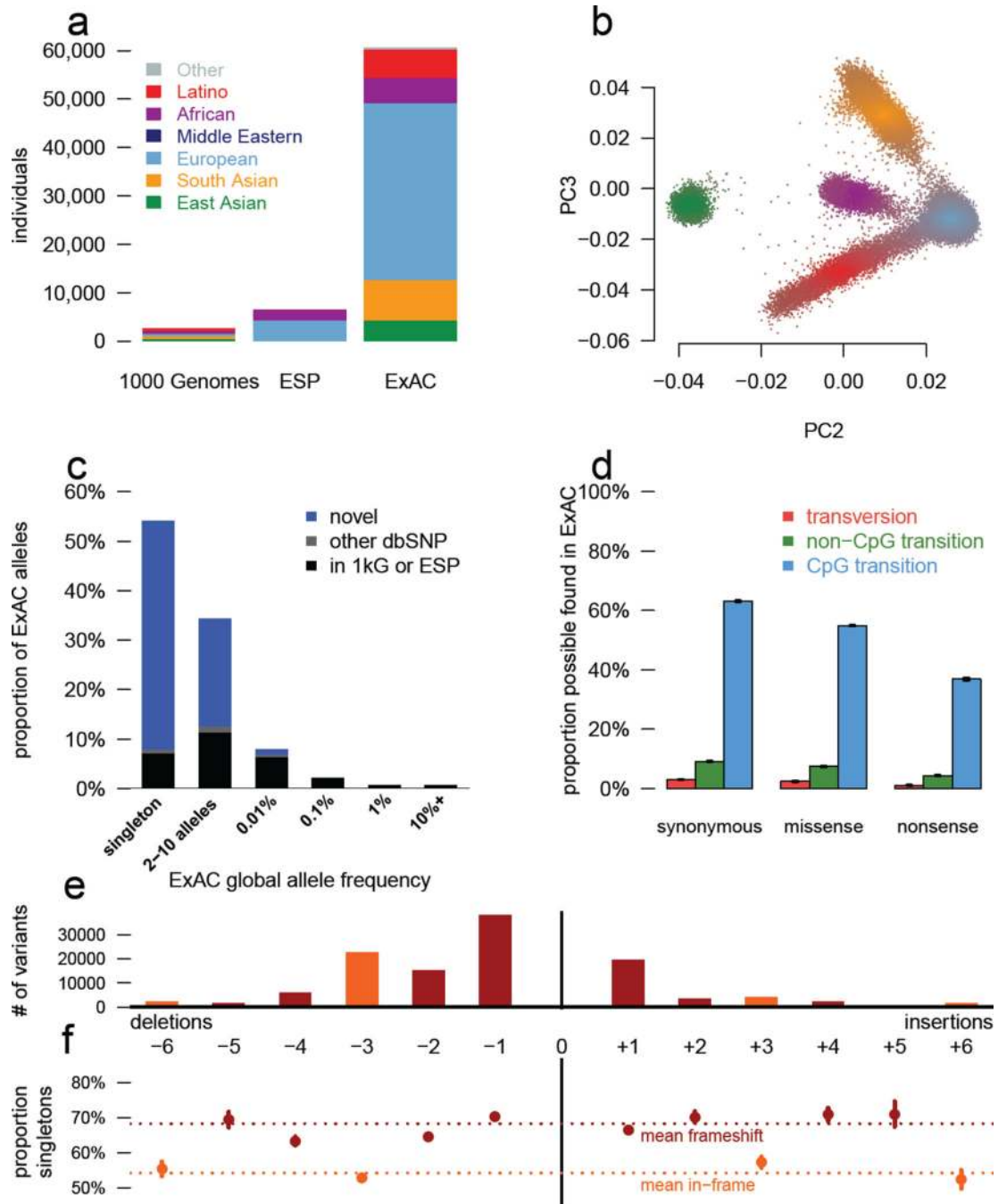


Figure 1. Patterns of genetic variation in 60,706 humans

a) The size and diversity of public reference exome datasets. ExAC exceeds previous datasets in size for all studied populations. b) Principal component analysis (PCA) dividing ExAC individuals into five continental populations. PC2 and PC3 are shown; additional PCs are in Extended Data Figure 5a. c) The allele frequency spectrum of ExAC highlights that the majority of genetic variants are rare and novel. d) The proportion of possible variation observed by mutational context and functional class. Over half of all possible CpG transitions are observed. Error bars represent standard error of the mean. e-f) The number (e)

and frequency distribution (proportion singleton; f) of indels, by size. Compared to in-frame indels, frameshift variants are less common (have a higher proportion of singletons, a proxy for predicted deleteriousness on gene product). Error bars indicate 95% confidence intervals.

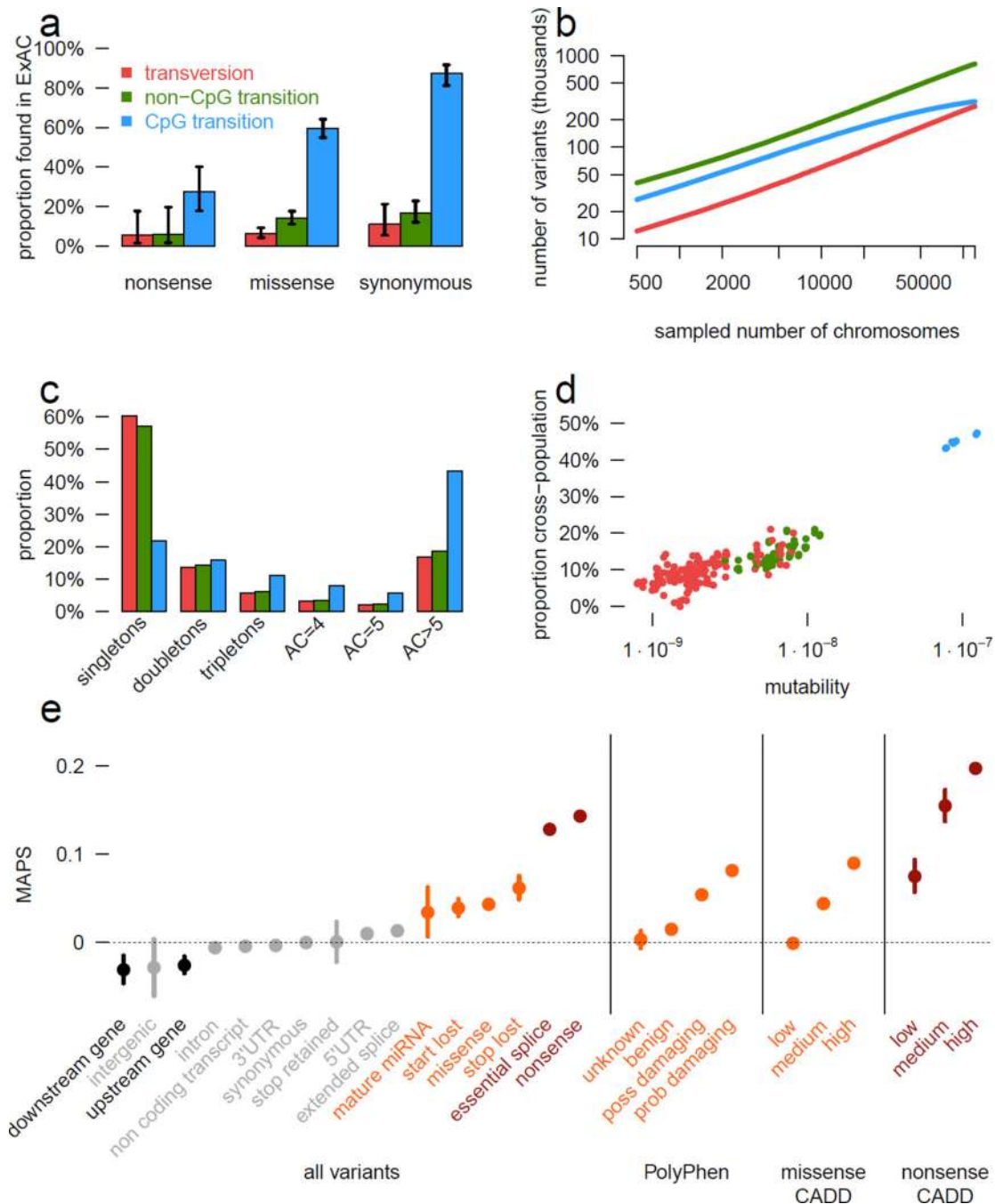


Figure 2. Mutational recurrence at large sample sizes

a) Proportion of validated *de novo* variants from two external datasets that are independently found in ExAC, separated by functional class and mutational context. Error bars represent standard error of the mean. Colors are consistent in a-d. b) Number of unique variants observed, by mutational context, as a function of number of individuals (down-sampled from ExAC). CpG transitions, the most likely mutational event, begin reaching saturation at ~20,000 individuals. c) The site frequency spectrum is shown for each mutational context. d) For doubletons (variants with an allele count of 2), mutation rate is positively correlated with

the likelihood of being found in two individuals of different continental populations. e) The mutability-adjusted proportion of singletons (MAPS) is shown across functional classes. Error bars represent standard error of the mean of the proportion of singletons.

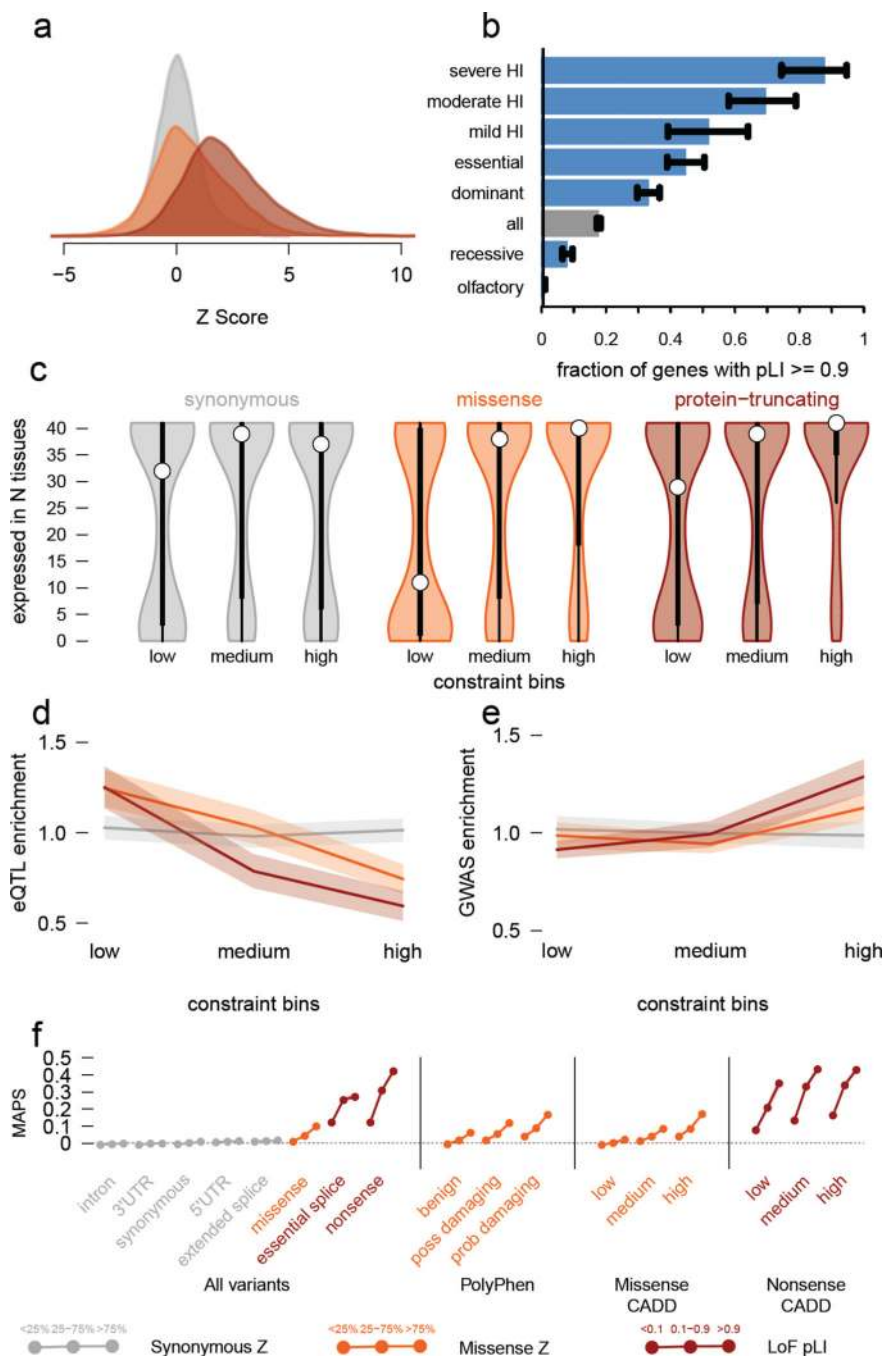


Figure 3. Quantifying intolerance to functional variation in genes and gene sets

a) Histograms of constraint Z scores for 18,225 genes. This measure of departure of number of variants from expectation is normally distributed for synonymous variants, but right-shifted (higher constraint) for missense and protein-truncating variants (PTVs), indicating that more genes are intolerant to these classes of variation. b) The proportion of genes that are very likely intolerant of loss-of-function variation ($pLI \geq 0.9$) is highest for ClinGen haploinsufficient genes, and stratifies by the severity and age of onset of the haploinsufficient phenotype. Genes essential in cell culture and dominant disease genes are

likewise enriched for intolerant genes, while recessive disease genes and olfactory receptors have fewer intolerant genes. Black error bars indicate 95% confidence intervals (CI). c) Synonymous Z scores show no correlation with the number of tissues in which a gene is expressed, but the most missense- and PTV-constrained genes tend to be expressed in more tissues. Thick black bars indicate the first to third quartiles, with the white circle marking the median. d) Highly missense- and PTV-constrained genes are less likely to have eQTLs discovered in GTEx as the average gene. Shaded regions around the lines indicate 95% CI. e) Highly missense- and PTV-constrained genes are more likely to be adjacent to GWAS signals than the average gene. Shaded regions around the lines indicate 95% CI. f) MAPS (Figure 2d) is shown for each functional category, broken down by constraint score bins as shown. Missense and PTV constraint score bins provide information about natural selection at least partially orthogonal to MAPS, PolyPhen, and CADD scores, indicating that this metric should be useful in identifying variants associated with deleterious phenotypes. Shaded regions around the lines indicate 95% CI. For panels a,c-f: synonymous shown in gray, missense in orange, and protein-truncating in maroon.

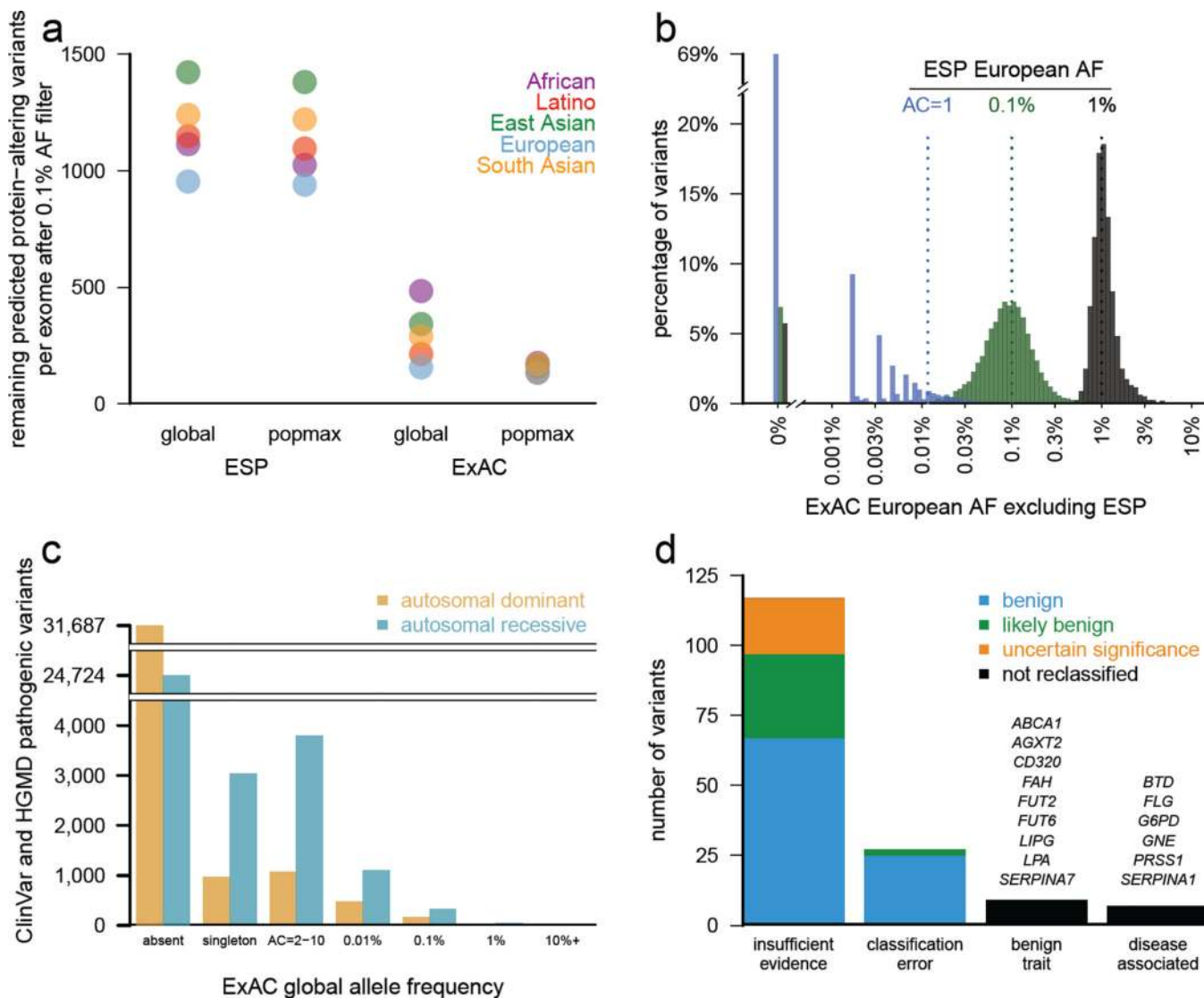


Figure 4. Filtering for Mendelian variant discovery

a) Predicted missense and protein-truncating variants in 500 randomly chosen ExAC individuals were filtered based on allele frequency information from ESP, or from the remaining ExAC individuals. At a 0.1% allele frequency (AF) filter, ExAC provides greater power to remove candidate variants, leaving an average of 154 variants for analysis, compared to 1090 after filtering against ESP. Popmax AF also provides greater power than global AF, particularly when populations are unequally sampled. b) Estimates of allele frequency in Europeans based on ESP are more precise at higher allele frequencies. Sampling variance and ascertainment bias make AF estimates unreliable, posing problems for Mendelian variant filtration. 69% of ESP European singletons are not seen a second time in ExAC (tall bar at left), illustrating the dangers of filtering on very low allele counts. c) Allele frequency spectrum of disease-causing variants in the Human Gene Mutation Database (HGMD) and/or pathogenic or likely pathogenic variants in ClinVar for well characterized autosomal dominant and autosomal recessive disease genes²⁸. Most are not found in ExAC; however, many of the reportedly pathogenic variants found in ExAC are at

too high a frequency to be consistent with disease prevalence and penetrance. d) Literature review of variants with >1% global allele frequency or >1% Latin American or South Asian population allele frequency confirmed there is insufficient evidence for pathogenicity for the majority of these variants. Variants were reclassified by ACMG guidelines²⁴.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

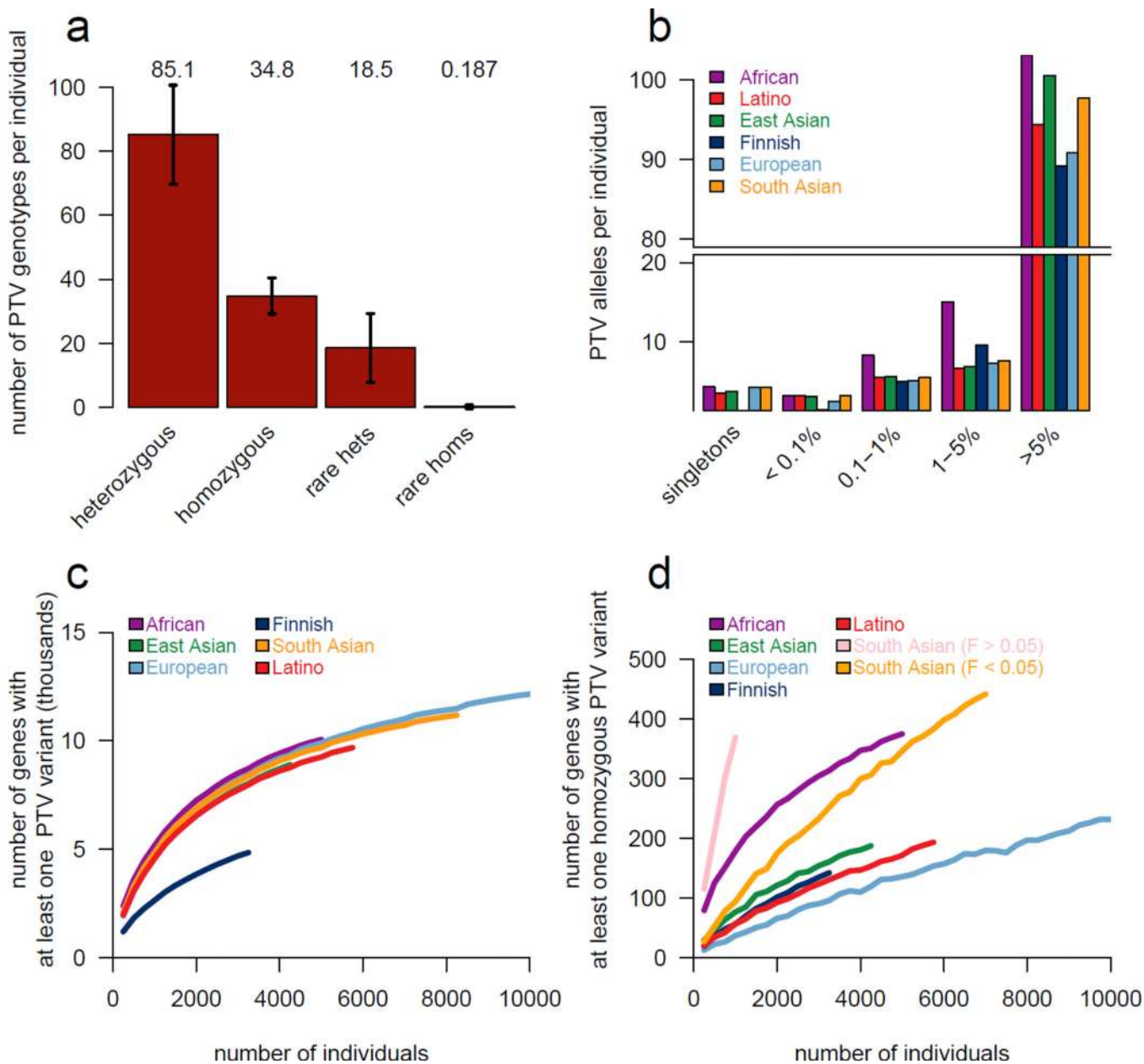


Figure 5. Protein-truncating variation in ExAC

a) The average ExAC individual has 85 heterozygous and 35 homozygous protein-truncating variants (PTVs), of which 18 and 0.19 are rare ($<0.1\%$ popmax AF), respectively. Error bars represent standard deviation. b) Breakdown of PTVs per individual (a) by popmax AF bin. Across all populations, most PTVs found in a given individual are common ($>5\%$ popmax AF). c-d) Number of genes with at least one PTV (c) or homozygous PTV (d) as a function of number of individuals, downsampled from ExAC. South Asian population is broken down by consanguinity (Inbreeding coefficient, F). At 60,000 individuals for ExAC, the plots in c) and d) extends to 15,750 with at least one PTV and 1,550 genes with at least one homozygous PTV.