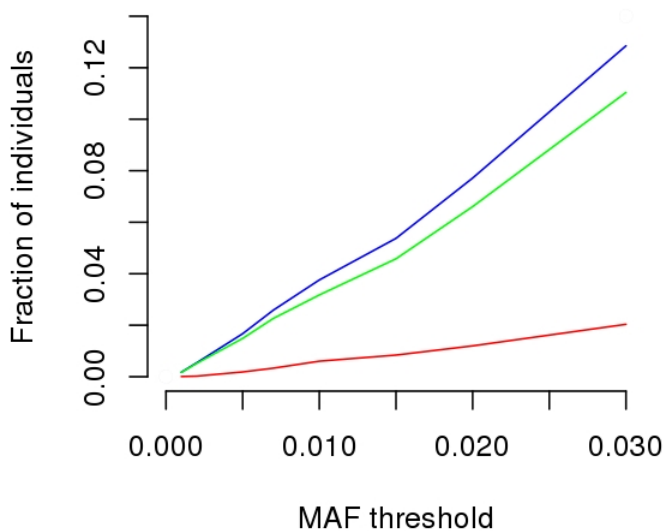
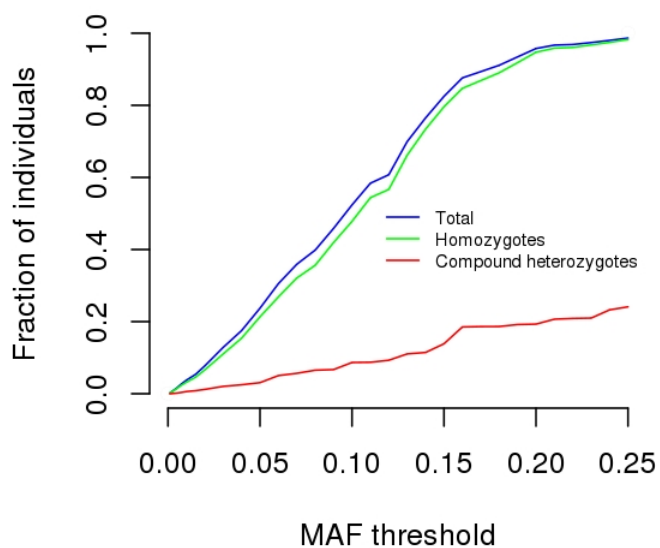


**Supplementary Figure 1**

**The probabilities of a variant being seen once and five times as a function of the number of individuals sequenced by minor allele frequency (MAF).**

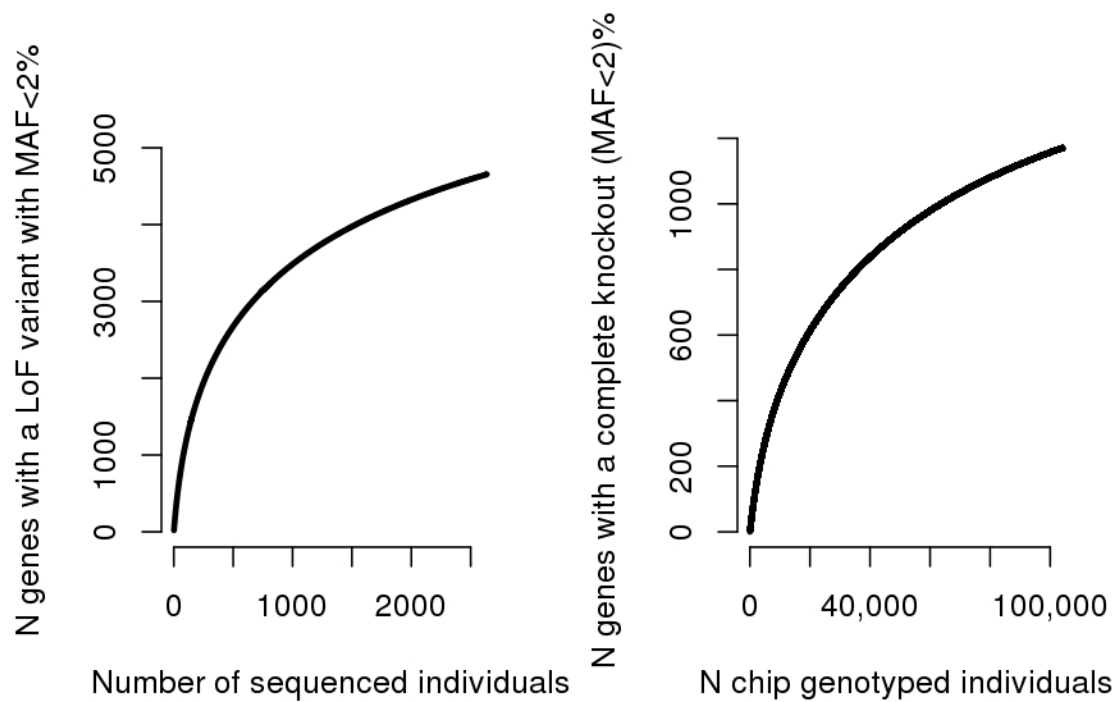
Variants that are seen at least five times, corresponding to an observed allelic frequency of 0.095%, are likely to be imputed with good quality.



**Supplementary Figure 2**

The fraction of individuals among 104,220 individuals with imputed genotypes that have genes completely knocked out by LoF variants with MAF below the given threshold.

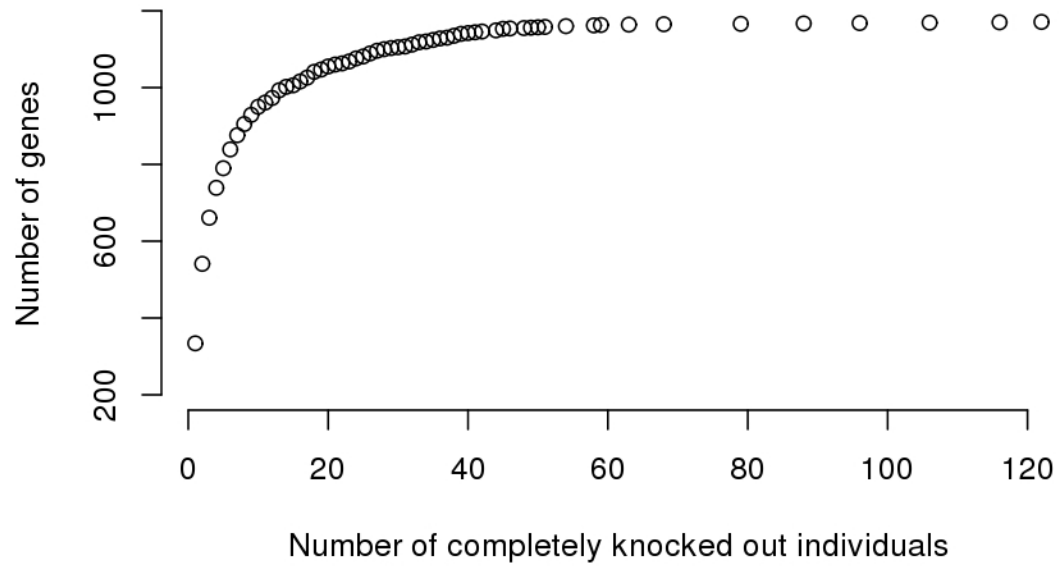
The second panel shows a magnified view of MAF below 3%.



**Supplementary Figure 3**

The number of genes that are observed to have at least one LoF variant with MAF below 2% as a function of the number of sequenced individuals and the number of genes that are completely knocked out in at least one individual by LoF variants with MAF below 2% as a function of the number of chip-genotyped individuals.

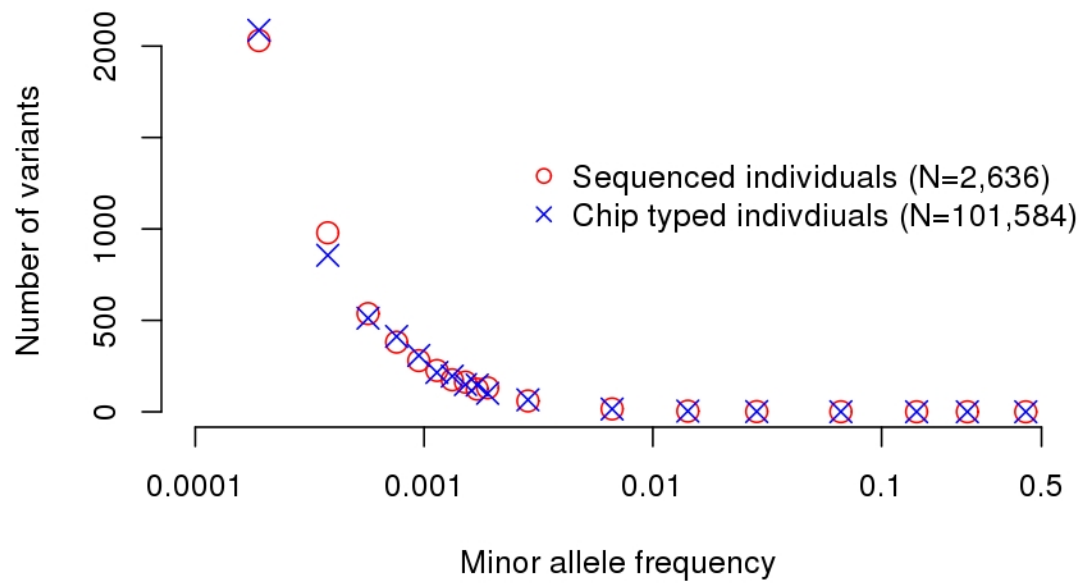
The curves are derived from the allele frequency distributions and the number of imputed complete knockouts.



**Supplementary Figure 4**

**The cumulative number of genes by the number of completely knocked out individuals.**

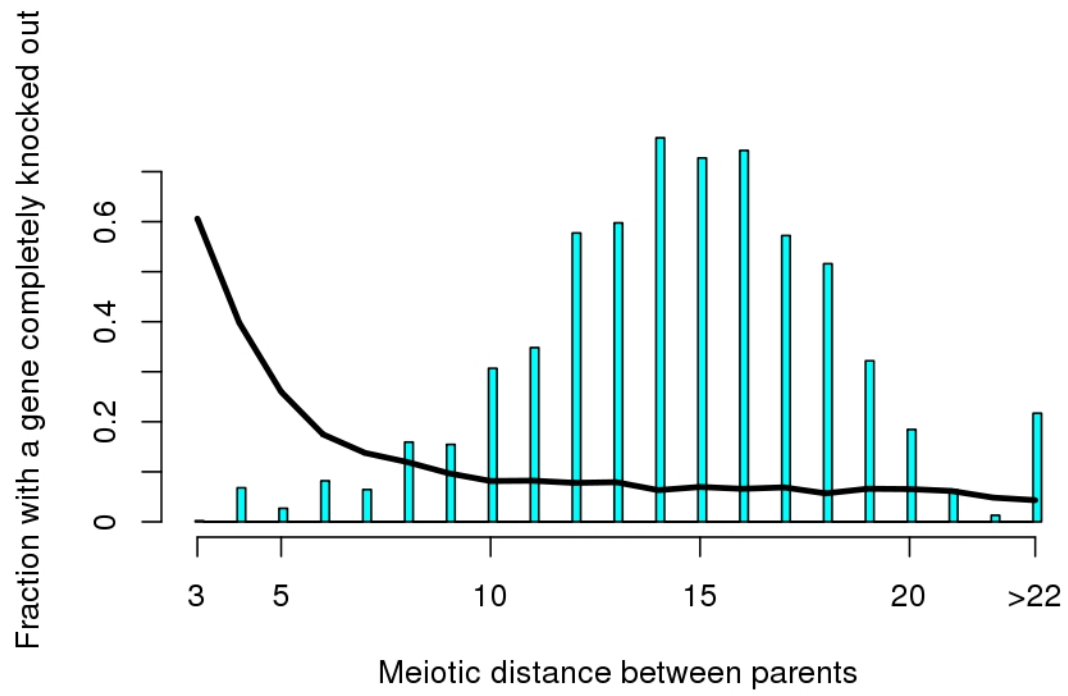
The total number of genes is 1,171.



**Supplementary Figure 5**

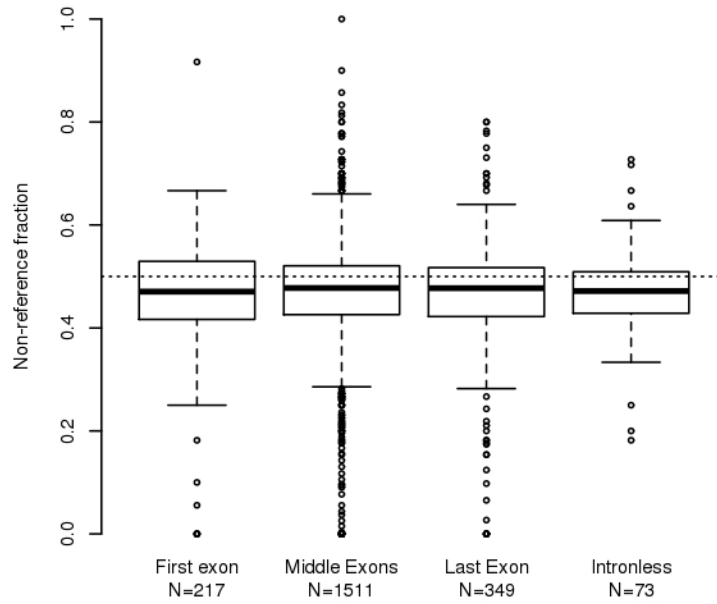
**The frequency distribution of the 6,795 LoF variants among sequenced and imputed individuals.**

The leftmost count includes all variants with a frequency less than one and half divided by twice the number of sequenced individuals, which corresponds to the number of variants seen only once in the sequenced set.



**Supplementary Figure 6**

**A histogram of the number of meiosis between the parents of the 104,220 Icelanders in our study and the fraction of individuals that have at least one gene completely knocked out by rare LoF variants.**



### Supplementary Figure 7

#### Transcriptome effect of synonymous SNPs by exon rank.

The allele-specific expression of the non-reference allele was calculated for each variant for a set of 262 individuals with blood RNA sequence data. The top, middle and bottom of the boxes are the top quartile, median and bottom quartile values calculated over the set of variants. The whiskers show the lowest and highest datum within 1.5 times the interquartile range (IQR) from the median. The dots indicate datum more than 1.5 times the IQR from the median. The  $n$  values given are the number of variants in each class.

## SUPPLEMENTARY NOTE

**The Icelandic study population.** This study is based on whole-genome sequence data from the white blood cells of 2,636 Icelanders participating in various disease projects at deCODE genetics<sup>1</sup> (Supplementary Tables 2 and 3). In addition, a total of 104,220 Icelanders have been genotyped using Illumina SNP chips<sup>1</sup>.

All participating individuals, or their guardians, gave their informed consent before blood samples were drawn. The family history of participants donating blood was incorporated into the study by including the phenotypes of first and second degree relatives and integrating over their possible genotypes.

All sample identifiers were encrypted in accordance with the regulations of the Icelandic Data Protection Authority. Approval for these studies was provided by the National Bioethics Committee and the Icelandic Data Protection Authority.

**The Icelandic Genealogy.** The Icelandic genealogical database contains 819,410 individuals back to 740 AD. Of the 471,284 Icelanders recorded to have been born in the 20<sup>th</sup> century, 91.1% had a recorded father and 93.7% had a recorded mother in the database. Similarly, of the 183,896 Icelanders recorded to have been born in the 19<sup>th</sup> century, 97.5% had a recorded father and 97.8% had a recorded mother.

**Illumina SNP Chip Genotyping.** The chip-typed samples were assayed with the Illumina HumanHap300, HumanCNV370, HumanHap610, HumanHap1M, HumanHap660, Omni-1, Omni 2.5 or Omni Express bead chips at deCODE genetics<sup>1</sup>. Chip SNPs were excluded if they had (i) yield less than 95%, (ii) minor allele frequency (MAF) less than 1% in the population or (iii) significant deviation from Hardy-Weinberg equilibrium ( $P < 0.001$ ), (iv) if they produced an excessive inheritance error rate (over 0.001), (v) if there was substantial difference in allele frequency between chip types (from just a single chip if that resolved all differences, but from all chips otherwise). All samples with a call rate below 97% were excluded from the analysis. The final set used for long-range phasing comprised 676,913 autosomal SNPs.

**Long range phasing.** Long range phasing of all chip-genotyped individuals was performed with methods described previously<sup>2,3</sup>. In brief, phasing is achieved using an iterative algorithm which phases a single proband at a time given the available phasing information about everyone else that shares a long haplotype identically by state with the proband. Given the large fraction of the Icelandic population that has been chip-typed, accurate long range



phasing is available genome-wide for all chip-typed Icelanders. For long range phased haplotype association analysis, we then partitioned the genome into non-overlapping fixed 0.3 cM bins. Within each bin, we observed the haplotype diversity described by the combination of all chip-typed markers in the bin.

**Whole-genome sequencing sample preparation.** Paired-end libraries for sequencing were prepared according to the manufacturer's instructions (Illumina, TruSeq™). In short, approximately 1 µg of genomic DNA, isolated from frozen blood samples, was fragmented to a mean target size of 300 bp using a Covaris E210 instrument. The resulting fragmented DNA was end repaired using T4 and Klenow polymerases and T4 polynucleotide kinase with 10 mM dNTP followed by addition of an 'A' base at the ends using Klenow exo fragment (3' to 5'-exo minus) and dATP (1 mM). Sequencing adaptors containing 'T' overhangs were ligated to the DNA products followed by agarose (2%) gel electrophoresis. Fragments of about 400-500 bp were isolated from the gels (QIAGEN Gel Extraction Kit), and the adaptor-modified DNA fragments were PCR enriched for ten cycles using Phusion DNA polymerase (Finnzymes Oy) and a PCR primer cocktail (Illumina). Enriched libraries were further purified using AMPure XP beads (Beckman-Coulter). The quality and concentration of the libraries were assessed with the Agilent 2100 Bioanalyzer using the DNA 1000 LabChip (Agilent). Barcoded libraries were stored at -20 °C. All steps in the workflow were monitored using an in-house laboratory information management system with barcode tracking of all samples and reagents.

**Whole-genome sequencing.** Template DNA fragments were hybridized to the surface of flow cells (GA PE cluster kit (v2) or HiSeq PE cluster kits (v2.5 or v3)) and amplified to form clusters using the Illumina cBot. In brief, DNA (2.5–12 pM) was denatured, followed by hybridization to grafted adaptors on the flow cell. Isothermal bridge amplification using Phusion polymerase was then followed by linearization of the bridged DNA, denaturation, blocking of 3' ends and hybridization of the sequencing primer. Sequencing-by-synthesis (SBS) was performed on Illumina GAIIx and/or HiSeq 2000 instruments. Paired-end libraries were sequenced at 2 × 101 (HiSeq) or 2 × 120 (GAIIx) cycles of incorporation and imaging using the appropriate TruSeq™ SBS kits. Each library or sample was initially run on a single GAIIx lane for QC validation followed by further sequencing on either GAIIx (≥4 lanes) or HiSeq (≥1 lane) with targeted raw cluster densities of 500–800 k/mm<sup>2</sup>, depending on the version of the data imaging and analysis packages (SCS2.6-2-9/RTA1.6-1.9, HCS1.3.8-

1.4.8/RTA1.10.36-1.12.4.2). Real-time analysis involved conversion of image data to base-calling in real-time.

**Whole-genome alignment.** Reads were aligned to NCBI Build 36 (hg18) of the human reference sequence using Burrows-Wheeler Aligner (BWA) 0.5.7-0.5.9<sup>4</sup>. Alignments were merged into a single BAM file and marked for duplicates using Picard 1.55 (<http://picard.sourceforge.net/>). Only non-duplicate reads were used for the downstream analyses. Resulting BAM files were realigned and recalibrated using GATK version 1.2-29-g0acaf2d<sup>5,6</sup>.

**Whole-genome SNP and INDEL calling.** Multi-sample variant calling was performed with GATK version 2.3.9 using all the 2,636 BAM files together.

Genotype calls made solely on the basis of next generation sequence data yield errors at a rate that decreases as a function of sequencing depth. Thus, for example, if sequence reads at a heterozygous SNP position carry one copy of the alternative allele and seven copies of the reference allele, then without further information the genotype would be called homozygous for the reference allele. To minimize the number of such errors, we used information about haplotype sharing, taking advantage of the fact that all the sequenced individuals had also been chip-typed and long range phased<sup>1,7</sup>. Extending the previous example, if the individual shares a haplotype with another who is heterozygous given his sequence reads, then the ambiguous individual would be called as heterozygous. Conversely, if the individual shares both his haplotypes with others who are homozygous for the major allele his genotype would be called homozygous. In order to improve genotype quality and to phase the sequencing genotypes, an iterative algorithm based on the IMPUTE HMM model<sup>8</sup> which uses the LRP haplotypes was employed. Assume a SNP with alleles 0 and 1 is being phased. We let  $H$  be the long range phased haplotypes of the sequenced individuals and applied the following hidden Markov model (HMM) based algorithm.

Assuming that at each marker  $i$  the haplotype  $h$  has a common ancestor with a haplotype in  $H \setminus \{h\}$  and denote the variable indicating this with the latent variable  $z_i \in H \setminus \{h\}$ , the hidden variable in the HMM. Then

$$\gamma_{h,k,i} = P(z_i = k | \text{all LRP markers}),$$

for all  $k \in H \setminus \{h\}$ . Given a haplotype  $h$  in  $H$ ,  $\gamma_{h,k}$  are calculated simultaneously for all  $k \in H \setminus \{h\}$  using the same HMM model as IMPUTE<sup>8</sup>. Given the Markov assumptions of the HMM, the model is fully specified by emission and transition probabilities.

We define the emission probabilities of the HMM at each marker  $i$  as:

$$P(z_i = k | \text{marker } i) = \begin{cases} 1 - \lambda, & \text{if } h \text{ and } k \text{ match at } i \\ \lambda, & \text{if } h \text{ and } k \text{ mismatch at } i \end{cases}$$

where  $\lambda$  can be thought of as a penalty for a mismatch. We used  $\lambda = 10^{-7}$  in our implementation. We define the transmission probabilities of the HMM model as:

$$P(z_i | z_{i-1}, \text{markers } 1, \dots, i-1) = \begin{cases} e^{-\frac{\rho_i}{N}} + \frac{1-e^{-\frac{\rho_i}{N}}}{N}, & \text{if } z_i = z_{i-1} \\ \frac{1-e^{-\frac{\rho_i}{N}}}{N}, & \text{if } z_i \neq z_{i-1} \end{cases}$$

Where  $N$  is the number of haplotypes in  $k \in H \setminus \{h\}$ , which for autosomal chromosomes is  $2(2,636 - 1)$  here and  $\rho_i = 4N_e r_i$ , where  $r_i$  is the genetic distance between markers  $i - 1$  and  $i$  according to the most recent version of the deCODE genetic map<sup>9</sup> and  $N_e$  was originally meant to be an estimate of the effective number of haplotypes in the population that our sample comes from, we used  $N_e = 7,000$ . These definitions fully specify the probability distribution  $P(z_i | \text{all markers})$ . Calculating  $\gamma_{h,k}$  for a single haplotype requires  $O(MN)$  operations, where  $N$  is the number of haplotypes and  $M$  is the number of markers. Since these calculations can be performed for one haplotype at a time, the calculations can be parallelized across a computer cluster for efficiency. In practice most of the  $\gamma_{h,k}$  will be close to zero and can be safely ignored (we used a threshold of  $10^{-6}$  of the largest value at each marker for each  $h$ ) greatly reducing storage requirements.

Now we are set to describe an iterative algorithm for the actual phasing. For every  $h$  in  $H$ , initialize the parameter  $\theta_h$ , which specifies how likely the one allele of the SNP is to occur on the background of  $h$  from the genotype likelihoods obtained from sequencing. The genotype likelihood  $L_g$  is the probability of the observed sequencing data at the SNP for a given individual assuming  $g$  is the true genotype at the SNP. If  $L_0$ ,  $L_1$  and  $L_2$  are the likelihoods of the genotypes 0, 1 and 2 in the individual that carries  $h$ , then set  $\theta_h$ :

$$\theta_h = \frac{L_2 + \frac{1}{2}L_1}{L_2 + L_1 + L_0}.$$

For every pair of haplotypes  $h$  and  $k$  in  $H$  that are carried by the same individual, use the other haplotypes in  $H$  to predict the genotype of the SNP on the backgrounds of  $h$  and  $k$ :

$$\tau_h = \sum_{l \in H \setminus \{h\}} \gamma_{h,l} \theta_l \text{ and}$$

$$\tau_k = \sum_{l \in H \setminus \{k\}} \gamma_{k,l} \theta_l.$$

Combining these predictions with the genotype likelihoods from sequencing gives un-normalized updated phased genotype probabilities:

$$P_{00} = (1 - \tau_h)(1 - \tau_k)L_0,$$

$$P_{10} = \tau_h(1 - \tau_k)\frac{1}{2}L_1,$$

$$P_{01} = (1 - \tau_h)\tau_k\frac{1}{2}L_1,$$

$$\text{and } P_{11} = \tau_h\tau_kL_2.$$

Now use these values to update  $\theta_h$  and  $\theta_k$  to:

$$\theta_h = \frac{P_{10} + P_{11}}{P_{00} + P_{01} + P_{10} + P_{11}} \text{ and}$$

$$\theta_k = \frac{P_{01} + P_{11}}{P_{00} + P_{01} + P_{10} + P_{11}}.$$

Iterate until the maximum difference between iterations is less than a convergence threshold  $\varepsilon$ . We used  $\varepsilon=10^{-7}$ .

**Genotype imputation.** Given the long range phased haplotypes of the sequenced individuals,  $H$ , and  $\theta$ , the carrier probabilities of the haplotypes in the sequenced set, the probability that a new haplotype  $n$ , not in the set of sequenced haplotypes  $H$ , is imputed as  $\sum_{l \in H} \gamma_{n,l} \theta_l$ , where  $\gamma_{n,l}$  is calculated as above for every  $l \in H$ .

In order to test the quality of our imputations, we compared chip genotypes to those imputed based on the sequence reads. Chip-based genotypes were available from 28,204 of the SNPs detected by sequencing that affect exons and splice regions, which proved to be fully concordant with genotypes from whole-genome sequencing. To assess imputation quality, we compared the imputed genotypes with those from the same set of 28,204 chip SNPs<sup>1</sup>. The concordance between imputed and chip-typed genotypes was high. For example, 76.4% of SNPs with a DAF of 1% were imputed almost perfectly ( $r^2 > 0.99$ ) and 98.4% of SNPs were imputed accurately ( $r^2 > 0.9$ ). Even for the SNPs that have only two copies of the minor allele observed in the set of sequenced individuals (DAF=0.04%) 47% were imputed almost perfectly ( $r^2 > 0.99$ ) and 64% of SNPs were imputed accurately ( $r^2 > 0.9$ ). We estimate the probability of a variant with MAF  $> 0.1\%$  being observed at least twice among the 2,636 sequenced individuals to be 96.7%.

**Genotype imputation information.** The informativeness of genotype imputation is estimated by the ratio of the variance of imputed expected allele counts and the variance of the actual allele counts:

$$\frac{\text{Var}(E(\theta|chip\ data))}{\text{Var}(\theta)},$$

where  $\theta$  is the allele count.  $\text{Var}(E(\theta|chip\ data))$  is estimated by the observed variance of the imputed expected counts.

In the sequenced set the information was calculated on a genotype level so that  $\theta \in \{0, 1, 2\}$  and  $\text{Var}(\theta) = 2p(1 - p)$ , assuming Hardy-Weinberg equilibrium, where  $p$  is the allele frequency. In the case of correct and accurate genotype calls the estimated information will be close to 1. Ambiguous genotype calls will usually lead to information estimates below 1. However, in some cases, due to artifacts, all observed allele calls will be 0 or 2 and the estimated information can go up to 2.

The imputations into chip typed individuals is done on a haplotype basis so that  $\theta \in \{0, 1\}$  and the estimated information is bounded between 0 and 1, where 1 represents exact imputations and ambiguous imputations yield values below 1.

**Whole-genome variant quality filtering.** The variants identified by GATK were filtered using thresholds on GATK variant call annotations<sup>10</sup>. SNPs were discarded if at least one of the following thresholds for their GATK call annotation parameter was violated: QD (Variant confidence/quality by depth) < 2.0, MQ (RMS mapping quality) < 40.0, FS (Fisher strand) > 60.0, HaploTypeScore > 13.0, MQRankSum < -12.5, and ReadPosRankSum < -8.0. Indels were discarded if at least one of the inequalities QD < 2.0, FS > 200.0, or ReadPosRankSum < -20.0 were satisfied. The thresholds for these parameters were adopted from GATK Best Practices (see URLs). In addition, SNPs and indels were discarded if one of the following thresholds were violated: DP (sequencing coverage/depth) > 110,000, AN (total number of alleles in called genotypes) < 4,200, and HW (Hardy-Weinberg P among sequenced samples) <  $10^{-7}$ , SI (genotype information among sequenced individuals) > 1.4, or if SI < 0.6 for SNPs and SI < 0.9 for indels. The GATK call annotation AN corresponds to the total number of chromosomes in called genotypes which equals  $2 \times 2,636 = 5,272$  when all chromosomes can be called. The additional filtering removes 2% of the remaining SNPs but 50% of the remaining indels. It is primarily the SI condition that removes indels, which usually indicates

that the failing indel is not called with a high degree of certainty and that its calling cannot be resolved in a coherent manner based on haplotype sharing.

Simple repeat regions were defined by combining the entire Simple Tandem Repeats by TRF track in UCSC hg18 with all homopolymer regions in hg18 of length 6 bp or more<sup>11</sup>. Variants called in these regions were ignored in the analysis.

**Liftover between hg18 and hg19.** Coordinates of variants and regions were converted between hg18 and hg19 using the liftOver tool from UCSC<sup>12</sup>.

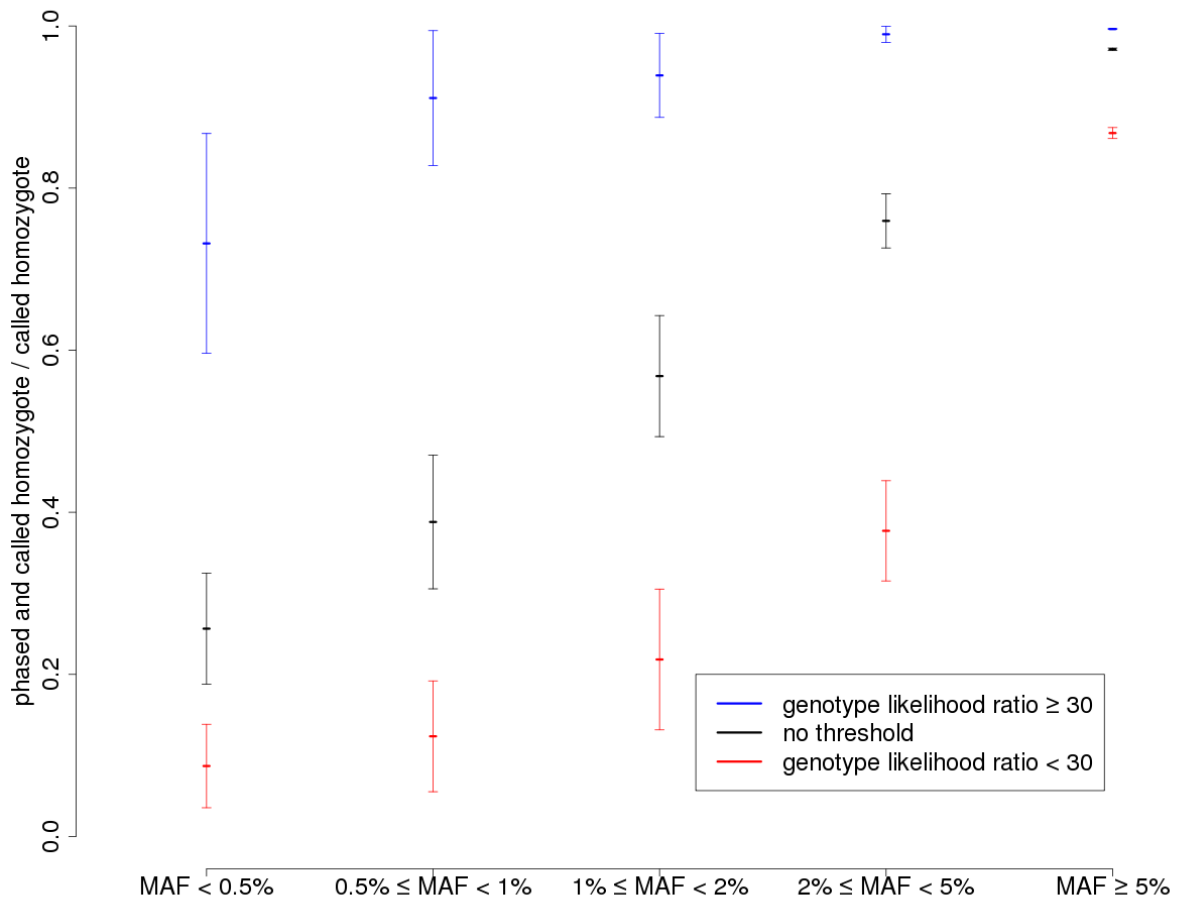
**Gene and variant annotation.** Variants were annotated with information from Ensembl release 70 using Variant Effect Predictor (VEP) version 2.8<sup>13,14</sup>. Only protein coding transcripts from RefSeq Release 56<sup>15</sup> were considered. Variants were annotated with the classification categories of impact LoF, MODERATE, LOW, and OTHER based on the Sequence Ontology (SO)<sup>16</sup> annotation from VEP (see Supplementary Table 1 for the definition of the classification of impact categories). Sequence variants that could be assigned to more than one category (primarily due to their impact on more than one gene transcript) were assigned to the most severe of the applicable categories.

**Estimation of carrier status based on imputation data.** An imputed individual was counted as a carrier of the minor allele of a variant if the imputation probability was greater than 0.9 for at least one of his two haplotypes and as homozygous carrier if the imputation probability was greater than 0.9 for both of his haplotypes.

**Exon position.** We divided genes into two categories: multi-exon genes, those with two or more exons, and intronless genes, those with only one exon based on the RefSeq set. The exons of multi-exon genes were further divided into three groups: (i) first exons (those that are the first exon in at least one multi-exon transcript), (ii) last exons (those that are the last exon in at least one multi-exon transcript), and (iii) middle exons (those that are never the first, the last or the only exon in a transcript).

**Overlap with ESP and dbSNP.** We assessed the overlap between the variants we discovered in Iceland with those in dbSNP<sup>17</sup> and with those reported by the ESP<sup>14,15</sup>. The Icelandic variants were counted as existing in the ESP or dbSNP if a variant at the same position and with the same allele was present in the database.

**Calling rare homozygotes.** With the depth of sequencing that we have employed, there is a chance that only one allele of a sequence variant will be represented in the sequencing reads from a truly heterozygous individual<sup>18</sup>. The probability that the homozygous minor allele genotype is the most likely genotype given only the sequencing reads of the individual when he is truly heterozygous goes up with lower sequencing depth at the variant and lower MAF. This has the practical consequence that for low MAFs, a large proportion of individuals who are most likely minor allele homozygotes given the sequencing data are truly heterozygotes. We do however call genotypes based on combining genotype likelihoods and haplotype sharing with other sequenced individuals which addresses this issue.



**Relatedness of parents in Iceland.** The Icelandic genealogical database was used to estimate the shortest distance between the parents of the 104,220 genotyped Icelanders. The distance between the parents was measured in the number of meiosis between them.

**Transmission.** Sharing between parent offspring pairs was tested and pairs with less than 99% haplotype sharing were excluded from the analysis. Haplotype sharing with other close family members was also tested and inconsistencies removed. We restricted our analysis to chip typed and long range phased individuals and excluded sequenced individuals to avoid observation bias. For example, if only the offspring of a triad has been sequenced, a variant that is transmitted is more likely to be imputed into the parent than a variant that was not transmitted. After exclusion, 35,024 father offspring, 47,769 mother offspring pairs and 26,188 triads remained. All the individuals are long-range phased which allows identification of recombinations. We excluded regions within 1Mb of recombinations in the offspring since these regions are most likely to suffer from phasing uncertainty and errors. For markers within these regions where both parents were imputed to be heterozygous (each allele being imputed to be the major or allele with 99% probability) and counted how many times the minor allele was transmitted to the offspring. This method does not rely on imputation into the offspring and should be robust to errors in phasing and genotyping, as such errors will shuffle the transmitted allele and are not likely to cause a bias.

Confidence intervals and significance were estimated using studentized bootstrap sampling<sup>19</sup> of parent offspring pairs and trios for the heterozygous and double heterozygous transmissions, respectively. The standard error of the estimated transmission probability was estimated by creating 10,000 bootstrap samples and the distribution of the transmission probability estimate was assumed to be approximately normally distributed around the true transmission probability. The 95% confidence interval around the observed transmission probability  $t$  with estimated standard error  $\sigma_B$  is then  $(t - 1.96\sigma_B, t + 1.96\sigma_B)$ .

**Validation of rare genotypes.** In order to validate some of the imputed homozygous rare LOF genotypes, we attempted to Sanger sequence over 211 SNP or indel variants in up to 302 individuals. Dye-terminator Sanger sequencing was performed with the Applied Biosystems BigDye Terminator v3.1 Cycle Sequencing Kit, with Agencourt Ampure XP and Agencourt CleanSeq for cleanup of the PCR and cycle sequencing product, respectively. AmpureXP and CleanSeq bead cleaning was performed on a Zymark Sciclone ALH-500 liquid-handling robot system. Tray dilutions for PCR setup and cycle-sequencing setup were prepared on a Packard MultiprobeII HTEX liquid-handling robot system, and genomic DNA and PCR product were transferred into their respective trays using the Zymark Sciclone ALH-500. PCR and cycle sequencing reactions were performed on MJ Research PTC-225 thermal cyclers. For signal detection, Applied Biosystems 3730xl DNA analyzers were used.



**Quantifying allele specific expression in white blood cells for LoF variants.** We estimated allele specific expression in blood for 262 individuals that been RNA sequenced. All of these individuals had imputed and phased genotypes.

**Preparation of Poly-A cDNA sequencing libraries.** The quality and quantity of isolated total RNA samples was assessed using the Total RNA 6000 Nano chip for the Agilent 2100 Bioanalyzer. cDNA libraries derived from Poly-A mRNA were generated using Illumina's TruSeq™ RNA Sample Prep Kit. Briefly, Poly-A mRNA was isolated from total RNA samples (1-4 µg input) using hybridization to Poly-T beads. The Poly-A mRNA was fragmented at 94°C, and first-strand cDNA was prepared using random hexamers and the SuperScript II reverse transcriptase (Invitrogen). Following second-strand cDNA synthesis, end repair, addition of a single A base, adaptor ligation, AMPure® bead purification, and PCR amplification, the resulting cDNA was measured on a Bioanalyzer using the DNA 1000 Lab Chip.

**RNA Sequencing.** Samples were clustered on to flowcells using Illumina's cBot and the TruSeq PE cluster kits v2, respectively. Paired-end sequencing (2x76 cycles) was performed on GAIIX instruments, equipped with paired-end modules using the TruSeq SBS kits v5 from Illumina. Approximately 125-175 million forward reads (250-350 M total reads) were sequenced per sample.

**Protocol and software used for RNA read alignment.** Read alignment and gene expression estimation were based on the recently published Tuxedo protocol described in<sup>20</sup>. The Tuxedo protocol describes how to use the open-source software tools TopHat<sup>21</sup> and Cufflinks<sup>22</sup> which have both been developed specifically for analysis of RNA-seq data. Sequencing reads were aligned to Homo sapiens Build 36 with TopHat version 1.4.1 with a supplied set of known transcripts in GTF format (RefSeq hg18; Homo sapiens, NCBI, build36.3). With this option, TopHat will first try to align reads to the provided transcriptome and only the reads that do not fully map to the transcriptome will then be mapped on the genome.

**Quantifying allele specific expression based on phased heterozygous loci.** To estimate the allele specific expression of a gene in which an individual was heterozygous for a LoF variant, reads for each allele were counted (using samtools mpileup<sup>23</sup>). Loci covered by skipped reads were ignored.

**Complete knockouts by human expression pattern.** Gene counts for tissue-specific expression in Table 3 and Supplementary Table 5 were based on publicly available data that

classifies tissue-specific expression of genes across twenty-seven tissues<sup>24</sup> (Supplementary Dataset 1). Counts in Table 3 are based on genes and tissues that have FPKM greater than or equal to 20 and are in the ‘Expressed in all high’ category in Fagerberg et al.<sup>24</sup>; the counts in Supplementary Table 6 are based the list of specific and enriched genes are autosomal genes in the categories ‘Highly tissue enriched’ and ‘Moderately tissue enriched’ from Fagerberg et al.<sup>24</sup>. The ‘Highly tissue enriched’ (‘Moderately tissue enriched’) category consists of genes-tissue pairs where the FPKM level for the gene expression in one of the tissues is 50-fold (5-fold) higher than for any of the other 26 tissues.

**Complete knockouts by phenotypes of mouse orthologs.** Genes were annotated with phenotypes from the Mouse Genome Database (MGD)<sup>25</sup> at the Mouse Genome Informatics (MGI) website. The MGD gene annotation was based on the tables in the files HMD\_HumanPhenotype.rpt and VOC\_MammalianPhenotype.rpt, by joining the tables on the Mammalian Phenotype (MP) accession number. Genes were annotated as recessive based on the Human Phenotype Ontology (HPO)<sup>26</sup> recessive gene list (based on Online Mendelian Inheritance in Man (OMIM)) and a recently published list for severe recessive diseases with pediatric presentations<sup>27</sup>; the genes on the latter list were annotated as ‘recessive-severe-genes’. Twenty eight mouse phenotype classes were considered. The ‘normal’ phenotype was not included in our analysis.

**Confidence intervals for expression patterns and mouse phenotype orthologs.** Confidence intervals for the number of completely knocked out genes by class were estimated using studentized bootstrap sampling (B=10,000)<sup>19</sup> of the genes in the class and estimated as described above.

<b>SO term</b>	<b>SO description</b>	<b>SO accession</b>	<b>Impact</b>
transcript ablation	A feature ablation whereby the deleted region includes a transcript feature	<a href="#">SO:0001893</a>	LoF
splice donor variant	A splice variant that changes the 2 base region at the 5' end of an intron	<a href="#">SO:0001575</a>	LoF
splice acceptor variant	A splice variant that changes the 2 base region at the 3' end of an intron	<a href="#">SO:0001574</a>	LoF
stop gained	A sequence variant whereby at least one base of a codon is changed, resulting in a premature stop codon, leading to a shortened transcript	<a href="#">SO:0001587</a>	LoF
frameshift variant	A sequence variant which causes a disruption of the translational reading frame, because the number of nucleotides inserted or deleted is not a multiple of three	<a href="#">SO:0001589</a>	LoF
stop lost	A sequence variant where at least one base of the terminator codon (stop) is changed, resulting in an elongated transcript	<a href="#">SO:0001578</a>	LoF
initiator codon variant	A codon variant that changes at least one base of the first codon of a transcript	<a href="#">SO:0001582</a>	LoF
inframe insertion	An inframe non synonymous variant that inserts bases into in the coding sequence	<a href="#">SO:0001821</a>	MODERATE
inframe deletion	An inframe non synonymous variant that deletes bases from the coding sequence	<a href="#">SO:0001822</a>	MODERATE
missense variant	A sequence variant, that changes one or more bases, resulting in a different amino acid sequence but where the length is preserved	<a href="#">SO:0001583</a>	MODERATE
transcript amplification	A feature amplification of a region containing a transcript	<a href="#">SO:0001889</a>	MODERATE
splice region variant	A sequence variant in which a change has occurred within the region of the splice site, either within 1-3 bases of the exon or 3-8 bases of the intron	<a href="#">SO:0001630</a>	MODERATE
incomplete terminal codon variant	A sequence variant where at least one base of the final codon of an incompletely annotated transcript is changed	<a href="#">SO:0001626</a>	MODERATE

**Supplementary Table 1.** Classification of the impact of sequence variant annotations from VEP<sup>13</sup>. The sequence variant annotation is based on the Sequence Ontology (SO)<sup>16</sup>. This classification is consistent with the order of severity of the sequence variants estimated by Ensembl.

<b>Disease</b>	<b>N</b>
Coronary Artery Disease	474
Chronic Kidney Disease	424
Obesity	394
Hypertension	371
Type 2 diabetes	307
Osteoporosis	260
Atrial Fibrillation	259
Myocardial Infarction	246
Alzheimer's disease	214
Asthma	211
Osteoarthritis	209
Urinary Tract Infection	208
Systemic Lupus Erythematosus	207
Alcohol Dependence	199
Breast Cancer	196
Depression	171
Sleep Apnea	165
Autism Spectrum Disorders	160
Kidney Stones	160
Nicotine Dependence	158
Prostate Cancer	158
Gallstones	156
Schizophrenia	155
Attention Deficit Hyperactivity Disorder	154
Glaucoma	138
Mental Retardation	138
Colorectal Adenoma	134
Migraine	132
Age Related Macular Degeneration	127
Ischaemic Stroke	122
Psoriasis	120
Epilepsy	118
Tuberculosis	118
Basal Cell Carcinoma Of The Skin	115
Parkinson's Disease	108
Diverticular Disease	107
Cataract	106
Bronchitis	98
Tourette Syndrome	84
Heart Failure	81
Emphysema	80
Dyslexia	79
Hypothyroidism	79
Sick Sinus Syndrome	78
Peripheral Artery Disease	77
Rheumatoid Arthritis	73
Benign Prostatic Hyperplasia	70
Congenital Heart Disease	70
Hypertension In Pregnancy	69
Panic Disorder	69

**Supplementary Table 2.** The 50 most prevalent conditions among the 2,636 sequenced Icelanders.

<b>Demographic</b>	<b>Sequenced</b>	<b>Chip-typed</b>	<b>Relatives of chip-typed</b>
N	2,636	104,220	294,212
Female (%)	54.0	55.1	46.8
YOB <sup>a</sup> (SD)	1950 (23)	1953 (23)	1964 (36)
Alive <sup>b</sup> (%)	72.5	84.9	72.9
Age <sup>c</sup> (SD)	55 (20)	56 (70)	33 (20)
Lifespan <sup>d</sup> (SD)	79 (13)	80 (13)	59 (28)

**Supplementary Table 3.** The demographics of the sequenced, chip-typed and relatives of chip-typed individuals. <sup>a</sup>Year of birth. <sup>b</sup>Fraction currently alive. <sup>c</sup>Current age for the living. <sup>d</sup>Age at death.

Gene	Chr	Pos (build 36)	Type	Consequence	Info	N <sub>seq</sub>	Sanger genotype		
							Hom	Het	Fail
<i>AKR1A1</i>	1	45,804,882	Indel	frameshift	0.98	56	3	0	1
<i>LRR1Q3</i>	1	74,394,016	Indel	frameshift	0.98	22	3	0	0
<i>GPR113</i>	2	26,394,482	Indel	frameshift	0.98	18	3	0	0
<i>GCKR</i>	2	27,598,876	SNP	stop gained	1.00	23	4	0	0
<i>RETSAT</i>	2	85,431,608	Indel	frameshift	0.97	44	2	0	1
<i>NCKAP5</i>	2	133,199,764	Indel	frameshift	0.97	10	3	1	0
<i>RBM43</i>	2	151,817,528	Indel	frameshift	0.97	46	4	0	0
<i>LRTM1</i>	3	54,927,549	Indel	frameshift	1.00	36	2	0	0
<i>THPO</i>	3	185,573,262	Indel	frameshift	0.99	21	1	2	0
<i>SLC6A19</i>	5	1,270,062	SNP	splice donor	0.98	40	1	0	0
<i>IL31RA</i>	5	55,231,568	Indel	frameshift	0.99	38	4	0	1
<i>GPR98</i>	5	90,043,957	SNP	stop gained	0.98	22	4	0	0
<i>ERAP1</i>	5	96,147,326	Indel	frameshift	0.98	41	4	0	0
<i>C2</i>	6	32,010,044	Indel	splice donor	0.99	95	11	0	1
<i>CUL9</i>	6	43,282,034	Indel	splice acceptor	0.99	41	4	0	0
<i>AH11</i>	6	135,801,205	SNP	splice donor	0.94	10	1	1	0
<i>GNAT3</i>	7	79,946,241	SNP	stop gained	0.97	22	2	0	0
<i>PPP1R3A</i>	7	113,306,396	Indel	frameshift	0.97	30	3	0	0
<i>RP1L1</i>	8	10,507,566	Indel	frameshift	0.96	21	1	0	0
<i>CNGB3</i>	8	87,725,124	Indel	frameshift	1.00	28	4	0	1
<i>FANCG</i>	9	35,065,702	Indel	frameshift	1.00	7	1	0	0
<i>DOLK</i>	9	130,749,402	Indel	frameshift	0.94	50	2	0	0
<i>PNPLA7</i>	9	139,564,448	Indel	frameshift	0.97	45	1	0	0
<i>CCDC7</i>	10	32,896,804	Indel	frameshift	0.99	27	4	0	0
<i>HPS1</i>	10	100,168,005	SNP	splice acceptor	0.88	5	2	0	0
<i>OR51G1</i>	11	4,901,932	Indel	frameshift	0.98	43	1	0	0
<i>PYGM</i>	11	64,283,799	SNP	stop gained	0.99	19	1	0	0
<i>CTSW</i>	11	65,403,926	Indel	frameshift	0.99	22	4	0	0
<i>NLRX1</i>	11	118,555,861	Indel	frameshift	1.00	50	4	0	0
<i>KLRF1</i>	12	9,888,308	Indel	frameshift	1.00	18	3	0	0
<i>GPRC5D</i>	12	12,993,748	Indel	frameshift	1.00	39	2	0	0
<i>RAD9B</i>	12	109,434,981	Indel	frameshift	0.99	20	2	0	0
<i>ATP7B</i>	13	51,432,390	Indel	frameshift	0.99	31	3	0	0
<i>ABHD4</i>	14	22,145,205	Indel	frameshift	0.99	17	2	0	0
<i>ARID4A</i>	14	57,901,772	SNP	splice donor	1.00	16	2	0	0
<i>SNAPC1</i>	14	61,329,314	Indel	frameshift	0.98	51	7	0	0
<i>MLH3</i>	14	74,568,565	Indel	frameshift	0.98	21	2	0	0
<i>SPTBN5</i>	15	39,945,301	Indel	frameshift	0.97	24	2	0	0
<i>TGM5</i>	15	41,336,054	Indel	splice donor	0.98	10	2	0	0
<i>RLBP1</i>	15	87,554,642	SNP	stop gained	0.95	11	2	0	0
<i>SYNM</i>	15	97,490,739	Indel	frameshift	0.98	25	1	0	0
<i>NOD2</i>	16	49,321,279	Indel	frameshift	0.97	32	1	0	0
<i>KIF19</i>	17	69,852,027	Indel	frameshift	0.99	28	3	2	0
<i>PGLS</i>	19	17,492,874	Indel	frameshift	0.98	35	5	0	1
<i>TUBB1</i>	20	57,028,006	Indel	frameshift	1.00	16	1	0	0
<i>SLCO4A1</i>	20	60,766,798	Indel	frameshift	0.99	22	2	0	0
<i>TMPRSS3</i>	21	42,682,220	Indel	frameshift	0.95	10	2	0	0
Total	-	-	-	-	-	-	128	6	6

**Supplementary Table 5.** The results of Sanger sequencing individuals predicted to be homozygous for a set of rare LoF variants, enriched for indels over SNPs. Shown are the affected gene, chromosome, build 36 position, variant type (SNP or Indel), maximal

consequence predicted by VEP, imputation information, number of whole-genome sequenced carriers, and the number of Sanger sequenced individuals genotyped as homozygous, heterozygous (not matching the predicted genotype) and that could not be determined.

Gene	Chr	Pos (build 36)	Type	Consequence	Sanger genotype		
					Match	Mismatch	Fail
<i>RPF1</i>	1	84,728,000	SNP	splice donor	1	0	0
<i>COL24A1</i>	1	86,080,386	SNP	stop gained	0	0	1
<i>GBP1</i>	1	89,294,359	Indel	frameshift	1	0	0
<i>KIAA1107</i>	1	92,416,480	Indel	frameshift	1	0	0
<i>FNBP1L</i>	1	93,769,013	Indel	frameshift	1	0	0
<i>ARHGAP29</i>	1	94,469,582	Indel	frameshift	1	0	0
<i>FNDC7</i>	1	109,061,921	SNP	splice acceptor	1	0	0
<i>ADAMTSL4</i>	1	148,794,615	Indel	frameshift	1	0	0
<i>ASH1L</i>	1	153,584,271	SNP	stop gained	1	0	0
<i>SLC25A44</i>	1	154,444,412	Indel	frameshift	1	0	0
<i>USF1</i>	1	159,279,245	SNP	splice donor	1	0	0
<i>RC3H1</i>	1	172,196,936	SNP	stop gained	1	0	0
<i>JMJD4</i>	1	225,988,267	Indel	frameshift	1	0	0
<i>OBSCN</i>	1	226,513,846	Indel	frameshift	1	0	0
<i>HEATR1</i>	1	234,785,579	SNP	splice acceptor	1	0	0
<i>OR14A16</i>	1	246,045,218	Indel	frameshift	1	0	0
<i>PGBD2</i>	1	247,178,343	SNP	stop gained	1	0	0
<i>PXDN</i>	2	1,637,273	SNP	stop gained	1	0	0
<i>LAPTM4A</i>	2	20,100,645	SNP	stop gained	1	0	0
<i>GCKR</i>	2	27,595,227	Indel	frameshift	1	0	0
<i>CRIM1</i>	2	36,557,636	Indel	frameshift	1	0	0
<i>SRBD1</i>	2	45,469,964	Indel	frameshift	1	0	0
<i>VPS54</i>	2	64,052,875	Indel	frameshift	1	0	0
<i>GAD1</i>	2	171,408,810	SNP	stop gained	1	0	0
<i>CCDC141</i>	2	179,517,533	SNP	stop gained	1	0	0
<i>TMEM194B</i>	2	191,089,263	SNP	stop gained	1	0	0
<i>HSPD1</i>	2	198,067,124	SNP	splice donor	1	0	0
<i>RPE</i>	2	210,592,681	Indel	frameshift	1	0	0
<i>ANKZF1</i>	2	219,805,329	SNP	splice donor	1	0	0
<i>ANKZF1</i>	2	219,808,833	Indel	frameshift	1	0	0
<i>AGAP1</i>	2	236,622,528	SNP	stop gained	1	0	0
<i>UBE2F</i>	2	238,603,953	SNP	stop gained	1	0	0
<i>OR6B3</i>	2	240,633,259	Indel	frameshift	1	0	0
<i>ATP2B2</i>	3	10,345,775	Indel	frameshift	1	0	0
<i>KRBOX1</i>	3	42,957,805	SNP	stop gained	1	0	0
<i>SNRK</i>	3	43,359,929	Indel	frameshift	1	0	0
<i>RASSF1</i>	3	50,343,885	SNP	splice acceptor	1	0	0
<i>KBTBD8</i>	3	67,137,323	SNP	stop gained	1	0	0
<i>SLC35A5</i>	3	113,782,121	Indel	frameshift	1	0	0
<i>CEP63</i>	3	135,760,852	Indel	frameshift	1	0	0
<i>PLSCR4</i>	3	147,401,511	SNP	splice donor	1	0	0
<i>GNB4</i>	3	180,619,896	SNP	stop gained	1	0	0
<i>CLDN1</i>	3	191,513,354	SNP	splice donor	1	0	0
<i>ATP13A4</i>	3	194,703,123	SNP	splice acceptor	1	0	0
<i>LRRCL5</i>	3	195,562,901	Indel	frameshift	1	0	0
<i>ZNF721</i>	4	428,035	SNP	stop gained	1	0	0
<i>BOD1L1</i>	4	13,210,220	Indel	frameshift	1	0	0
<i>DHX15</i>	4	24,138,728	Indel	frameshift	1	0	0
<i>FRYL</i>	4	48,286,518	SNP	splice donor	1	0	0
<i>UGT2B11</i>	4	70,114,689	Indel	frameshift	0	0	1
<i>SEC31A</i>	4	83,959,376	SNP	stop gained	1	0	0
<i>PDE5A</i>	4	120,703,435	SNP	splice donor	1	0	0
<i>KIAA1109</i>	4	123,458,761	Indel	frameshift	1	0	0



<i>SPRY1</i>	4	124,543,023	Indel	frameshift	1	0	0
<i>ANKRD31</i>	5	74,449,746	Indel	frameshift	1	0	0
<i>ARSK</i>	5	94,944,390	Indel	frameshift	1	0	0
<i>SEPT8</i>	5	132,126,120	Indel	frameshift	1	0	0
<i>PCDHB7</i>	5	140,534,842	SNP	stop gained	0	1	0
<i>PCDHGA1</i>	5	140,690,906	Indel	frameshift	1	0	0
<i>TCERG1</i>	5	145,823,310	Indel	frameshift	1	0	0
<i>RIPK1</i>	6	3,034,821	Indel	frameshift	1	0	0
<i>TJAP1</i>	6	43,581,406	Indel	frameshift	1	0	0
<i>TDRD6</i>	6	46,766,887	Indel	frameshift	1	0	0
<i>PKHD1</i>	6	52,038,772	Indel	frameshift	0	0	1
<i>DOPEY1</i>	6	83,923,681	Indel	frameshift	1	0	0
<i>COQ3</i>	6	99,924,209	Indel	frameshift	1	0	0
<i>METTL24</i>	6	110,750,769	SNP	splice acceptor	1	0	0
<i>MED23</i>	6	131,965,093	SNP	stop gained	1	0	0
<i>MED23</i>	6	131,967,064	Indel	frameshift	0	0	1
<i>UTRN</i>	6	145,144,830	Indel	frameshift	1	0	0
<i>INTS1</i>	7	1,478,549	SNP	splice acceptor	1	0	0
<i>TRIM50</i>	7	72,376,322	SNP	splice donor	0	1	0
<i>ABCB1</i>	7	87,012,099	Indel	frameshift	1	0	0
<i>ZKSCAN5</i>	7	98,967,205	Indel	frameshift	1	0	0
<i>PUS7</i>	7	104,930,102	SNP	splice donor	1	0	0
<i>HBP1</i>	7	106,627,976	Indel	frameshift	1	0	0
<i>TUSC3</i>	8	15,575,717	SNP	splice donor	1	0	0
<i>TNFRSF10B</i>	8	22,937,664	SNP	stop gained	1	0	0
<i>EBF2</i>	8	25,774,592	Indel	frameshift	1	0	0
<i>ESCO2</i>	8	27,702,411	Indel	frameshift	1	0	0
<i>SDCBP</i>	8	59,654,887	SNP	stop gained	1	0	0
<i>ASPH</i>	8	62,659,055	Indel	splice donor	1	0	0
<i>GDAP1</i>	8	75,438,881	Indel	frameshift	1	0	0
<i>DCAF13</i>	8	104,508,657	SNP	splice donor	1	0	0
<i>FOCAD</i>	9	20,897,164	Indel	frameshift	1	0	0
<i>IFNW1</i>	9	21,131,250	Indel	frameshift	1	0	0
<i>POLR1E</i>	9	37,476,805	SNP	splice donor	1	0	0
<i>CNTNAP3</i>	9	39,168,221	Indel	frameshift	0	0	1
<i>RMI1</i>	9	85,805,947	Indel	frameshift	1	0	0
<i>PHF2</i>	9	95,465,694	Indel	frameshift	1	0	0
<i>RNF20</i>	9	103,363,232	SNP	stop gained	1	0	0
<i>EPB41L4B</i>	9	111,060,351	Indel	frameshift	1	0	0
<i>RC3H2</i>	9	124,653,538	SNP	splice acceptor	1	0	0
<i>CAMSAP1</i>	9	137,881,802	Indel	splice donor	1	0	0
<i>APBB1IP</i>	10	26,842,478	Indel	frameshift	1	0	0
<i>BAMBI</i>	10	29,010,935	SNP	stop gained	1	0	0
<i>KIAA1462</i>	10	30,358,117	SNP	stop gained	1	0	0
<i>C10orf68</i>	10	33,057,932	SNP	splice donor	1	0	0
<i>GJD4</i>	10	35,936,511	SNP	splice acceptor	1	0	0
<i>PCDH15</i>	10	55,238,694	SNP	stop gained	1	0	0
<i>PCDH15</i>	10	56,094,027	SNP	initiator codon	1	0	0
<i>CCAR1</i>	10	70,202,759	Indel	frameshift	1	0	0
<i>MMP26</i>	11	4,967,664	Indel	frameshift	1	0	0
<i>TRIM48</i>	11	54,792,401	SNP	splice acceptor	1	0	0
<i>UBE2L6</i>	11	57,084,385	SNP	splice donor	1	0	0
<i>PATL1</i>	11	59,182,941	SNP	stop gained	1	0	0
<i>TMEM135</i>	11	86,708,068	SNP	splice donor	1	0	0
<i>PPP2R1B</i>	11	111,117,994	Indel	frameshift	1	0	0

<i>ARHGEF12</i>	11	119,805,652	SNP	stop gained	1	0	0
<i>PRB4</i>	12	11,352,511	SNP	stop gained	0	0	1
<i>HEBP1</i>	12	13,019,510	SNP	stop lost	1	0	0
<i>LMBR1L</i>	12	47,778,054	Indel	frameshift	1	0	0
<i>MYL6B</i>	12	54,835,239	Indel	frameshift	1	0	0
<i>STAB2</i>	12	102,680,237	Indel	frameshift	1	0	0
<i>CCDC63</i>	12	109,826,980	SNP	splice donor	1	0	0
<i>COQ5</i>	12	119,426,177	Indel	frameshift	1	0	0
<i>ZNF140</i>	12	132,193,175	SNP	stop gained	1	0	0
<i>TRPC4</i>	13	37,164,429	SNP	stop gained	1	0	0
<i>PHF11</i>	13	48,990,187	Indel	frameshift	1	0	0
<i>PHF11</i>	13	48,998,596	Indel	splice donor	1	0	0
<i>NOP9</i>	14	23,842,872	Indel	frameshift	1	0	0
<i>BAZ1A</i>	14	34,301,136	Indel	frameshift	1	0	0
<i>RALGAP1</i>	14	35,217,012	SNP	stop gained	1	0	0
<i>ABHD12B</i>	14	50,438,332	Indel	frameshift	1	0	0
<i>CCDC175</i>	14	59,081,653	SNP	stop gained	1	0	0
<i>SYNE2</i>	14	63,618,026	Indel	frameshift	1	0	0
<i>SEL1L</i>	14	81,013,275	SNP	stop gained	1	0	0
<i>ATG2B</i>	14	95,852,687	SNP	stop gained	1	0	0
<i>ANKDD1A</i>	15	63,036,499	SNP	stop gained	1	0	0
<i>SLC24A1</i>	15	63,703,604	Indel	frameshift	1	0	0
<i>AKAP13</i>	15	83,990,144	Indel	frameshift	1	0	0
<i>SLX4</i>	16	3,580,804	SNP	stop gained	1	0	0
<i>DOC2A</i>	16	29,925,332	SNP	splice acceptor	1	0	0
<i>NOD2</i>	16	49,321,314	SNP	splice donor	1	0	0
<i>CES3</i>	16	65,555,260	Indel	frameshift	1	0	0
<i>FA2H</i>	16	73,307,795	SNP	stop gained	1	0	0
<i>DNAH2</i>	17	7,645,627	SNP	stop gained	1	0	0
<i>KRBA2</i>	17	8,213,920	Indel	frameshift	1	0	0
<i>ABCC3</i>	17	46,107,722	Indel	frameshift	1	0	0
<i>STXBP4</i>	17	50,479,512	Indel	splice donor	1	0	0
<i>FBF1</i>	17	71,430,875	Indel	frameshift	1	0	0
<i>SMCHD1</i>	18	2,740,487	Indel	frameshift	1	0	0
<i>PPP4R1</i>	18	9,573,114	SNP	splice donor	1	0	0
<i>AFG3L2</i>	18	12,367,037	SNP	stop gained	0	0	1
<i>TRAPPC8</i>	18	27,765,477	Indel	frameshift	1	0	0
<i>SERPINB13</i>	18	59,406,864	SNP	splice acceptor	1	0	0
<i>R3HDM4</i>	19	851,108	Indel	frameshift	1	0	0
<i>DENND1C</i>	19	6,419,658	SNP	splice acceptor	1	0	0
<i>C19orf45</i>	19	7,476,323	SNP	splice donor	1	0	0
<i>MYO1F</i>	19	8,501,250	SNP	splice acceptor	1	0	0
<i>MUC16</i>	19	8,934,495	SNP	stop gained	1	0	0
<i>JAK3</i>	19	17,814,916	Indel	frameshift	1	0	0
<i>ZNF793</i>	19	42,705,271	Indel	splice acceptor	1	0	0
<i>FCGBP</i>	19	45,049,601	SNP	splice acceptor	1	0	0
<i>PSG11</i>	19	48,220,940	SNP	stop gained	0	1	0
<i>TMEM150B</i>	19	60,524,154	SNP	stop gained	1	0	0
<i>CDS2</i>	20	5,113,502	SNP	splice acceptor	1	0	0
<i>DZANK1</i>	20	18,383,966	Indel	frameshift	1	0	0
<i>KIF3B</i>	20	30,361,887	Indel	frameshift	1	0	0
<i>LSS</i>	21	46,435,572	Indel	frameshift	1	0	0
<i>GSTT1</i>	22	22,706,821	SNP	splice donor	1	0	0
<i>SLC5A4</i>	22	30,950,354	Indel	frameshift	1	0	0
Total	-	-	-	-	152	3	7

**Supplementary Table 6.** The results of Sanger sequencing individuals who were observed to be the only carrier of a LoF variant among the 2,636 whole-genome sequenced Icelanders. Shown are the affected gene, chromosome, build 36 position, variant type (SNP or Indel), maximal consequence predicted by VEP and indicators of whether the Sangers sequencing was successful and matched the whole-genome sequencing.

	<b>All frequencies</b>	<b>MAF &lt; 2%</b>	<b>MAF &lt; 0.5%</b>
Number of genes	402		
Number of genes with LoF variants	117	79	67
Number of LoF variants	140	94	77
Completely knocked out genes	-	34	20
Completely knocked out individuals	-	251	43
Variants contributing to knockout	-	39	21

**Supplementary Table 7.** The number of sensory perception of smell genes affected by LoF variants and individuals that have these genes completely knocked out by LoF variants (GO:0007608).

	<b>Sequenced and phased</b>		<b>Imputed</b>	
	<b>Homozygote</b>	<b>Heterozygote</b>	<b>Homozygote</b>	<b>Heterozygote</b>
<b>MAF &lt; 0.5%</b>	0.018	9.4	0.016	8.5
<b>0.5% &lt; MAF &lt; 2%</b>	0.066	9.5	0.056	8.9
<b>MAF &gt; 2%</b>	21.0	111.1	20.5	110.0
<b>Total</b>	21.1	130.0	20.6	127.4

**Supplementary Table 8.** Average LoF minor allele genotype counts per individual stratified on minor allele frequency (MAF).

Tissue	N genes	Median coding size	Knocked out genes			
			MAF < 2 %		MAF < 0.5 %	
			N	Percentage (95% CI)	N	Percentage (95% CI)
brain	318	1,407	4	1.3 (0.0-2.5)	3	0.9 (0.0-2.0)
placenta	71	1,008	4	5.6 (0.1-11.2)	4	5.6 (0.1-11.1)
skin	104	1,078	6	5.8 (1.2-10.4)	6	5.8 (1.2-10.4)
kidney	62	1,413	4	6.5 (0.2-12.7)	2	3.2 (0.0-7.7)
pancreas	44	906	3	6.8 (0.0-14.5)	2	4.5 (0.0-11.0)
liver	171	1,377	13	7.6 (3.5-11.7)	12	7.0 (3.1-11.0)
esophagus	61	1,207	5	8.2 (1.1-15.3)	4	6.6 (0.1-13.1)
bone marrow	82	768	7	8.5 (2.2-14.9)	1	1.2 (0.0-3.6)
salivary gland	46	627	4	8.7 (0.1-17.2)	3	6.5 (0.0-13.9)
testis	885	1,156	102	11.5 (9.4-13.7)	72	8.1 (6.3-9.9)
heart	98	1,260	12	12.2 (5.4-19.1)	9	9.2 (3.2-15.1)

**Supplementary Table 9.** Counts of genes completely knocked out by rare LoF variants for genes with human tissue-specific expression<sup>28</sup>. Counts are shown for tissues that have at least 40 (autosomal) tissue-specific genes.

			<b>Knocked out genes</b>								
			<b>MAF &lt; 2%</b>				<b>MAF &lt; 0.5%</b>				
			<b>RefSeq genes</b>		<b>Genes linked to conditions under a recessive mode of inheritance<sup>26</sup></b>		<b>RefSeq genes</b>		<b>Genes linked to conditions under a recessive mode of inheritance<sup>26</sup></b>		
					<b>Severe diseases<sup>27</sup></b>				<b>Severe diseases<sup>27</sup></b>		
<b>Mouse phenotype</b>	<b>N</b>	<b>Median coding size</b>	<b>N</b>	<b>Percentage (95% CI)</b>	<b>Other conditions<sup>c</sup></b>	<b>N</b>	<b>Percentage (95% CI)</b>	<b>Other conditions<sup>c</sup></b>	<b>N</b>	<b>Percentage (95% CI)</b>	<b>Other conditions<sup>c</sup></b>
craniofacial	1,005	1,647	23	2.3 (1.4,3.2)	5	3	1.7 (0.9,2.5)	3	17	1.7 (0.9,2.5)	3
pigmentation	362	1,608	9	2.5 (0.9,4.1)	4	2	2.2 (0.7,3.7)	4	8	2.2 (0.7,3.7)	4
embryogenesis	1,530	1,575	43	2.8 (2.0,3.6)	5	4	2.2 (1.4,2.9)	4	33	2.2 (1.4,2.9)	4
digestive/alimentary	1,142	1,554	35	3.1 (2.1,4.1)	3	4	2.3 (1.4,3.1)	3	26	2.3 (1.4,3.1)	3
skeleton	1,483	1,509	49	3.3 (2.4,4.2)	8	1	2.5 (1.7,3.3)	6	37	2.5 (1.7,3.3)	6
mortality/aging	3,897	1,554	140	3.6 (3.0,4.2)	15	13	2.4 (1.9,2.8)	9	92	2.4 (1.9,2.8)	9
adipose tissue	647	1,509	24	3.7 (2.2,5.2)	6	2	2.9 (1.6,4.3)	6	19	2.9 (1.6,4.3)	6
growth/size/body	3,142	1,560	116	3.7 (3.1,4.3)	18	12	2.4 (1.9,2.9)	12	76	2.4 (1.9,2.9)	12
muscle	1,224	1,569	45	3.7 (2.6,4.7)	5	2	2.4 (1.5,3.2)	3	29	2.4 (1.5,3.2)	3
cardiovascular system	2,004	1,567	77	3.8 (3.0,4.7)	13	5	2.7 (2.0,3.4)	11	55	2.7 (2.0,3.4)	11
nervous system	2,814	1,599	108	3.8 (3.1,4.5)	17	11	2.8 (2.2,3.4)	12	78	2.8 (2.2,3.4)	12
integument	1,519	1,581	60	3.9 (3.0,4.9)	7	4	3.0 (2.1,3.8)	7	45	3.0 (2.1,3.8)	7
vision/eye	1,187	1,566	46	3.9 (2.8,5.0)	12	7	3.2 (2.2,4.2)	11	38	3.2 (2.2,4.2)	11
endocrine/exocrine gland	1,474	1,467	59	4.0 (3.0,5.0)	8	5	3.1 (2.3,4.0)	7	46	3.1 (2.3,4.0)	7
cellular	2,792	1,527	114	4.1 (3.4,4.8)	10	8	3.0 (2.4,3.7)	8	85	3.0 (2.4,3.7)	8
respiratory system	1,149	1,513	48	4.2 (3.0,5.3)	4	5	2.7 (1.8,3.6)	3	31	2.7 (1.8,3.6)	3
behavior/neurological	2,462	1,527	107	4.3 (3.6,5.1)	13	11	3.1 (2.4,3.7)	10	76	3.1 (2.4,3.7)	10
limbs/digits/tail	776	1,611	34	4.4 (2.9,5.8)	5	0	3.0 (1.8,4.1)	3	23	3.0 (1.8,4.1)	3
liver/biliary system	992	1,515	46	4.6 (3.3,6.0)	4	5	3.4 (2.3,4.6)	4	34	3.4 (2.3,4.6)	4
hearing/vestibular/ear	529	1,593	25	4.7 (2.9,6.5)	11	2	3.6 (2.0,5.2)	8	19	3.6 (2.0,5.2)	8
hematopoietic system	2,458	1,446	119	4.8 (4.0,5.7)	14	7	3.3 (2.6,3.9)	10	80	3.3 (2.6,3.9)	10
tumorigenesis	773	1,519	37	4.8 (3.3,6.3)	3	2	3.4 (2.1,4.6)	3	26	3.4 (2.1,4.6)	3
reproductive system	1,579	1,512	79	5.0 (3.9,6.1)	11	5	3.9 (3.0,4.9)	9	62	3.9 (3.0,4.9)	9
immune system	2,651	1,450	136	5.1 (4.3,5.9)	13	7	3.4 (2.7,4.1)	11	90	3.4 (2.7,4.1)	11
renal/urinary system	1,004	1,573	53	5.3 (3.9,6.7)	11	7	4.1 (2.9,5.3)	9	41	4.1 (2.9,5.3)	9
homeostasis/metabolism	3,422	1,467	186	5.4 (4.7,6.2)	22	13	3.8 (3.2,4.4)	18	129	3.8 (3.2,4.4)	18

other	270	1,638	17	6.3 (3.3,9.3)	3	2	9	3.3 (1.2,5.5)	2	2
taste/olfaction	120	1,522	11	9.2 (3.7,14.6)	0	0	10	8.3 (3.2,13.5)	0	0
Any mouse phenotype <sup>a</sup>	7,174	1,485	361	5.0 (4.6,5.4)	16	44	262	3.7 (3.3,4.0)	13	32
			(2,240)		(47)	(237)	(583)		(27)	(58)

**Supplementary Table 10.** Counts of gene and individual knock-outs with rare variants (DAF < 2%) for human orthologs of mouse genes with phenotypes in the Mouse Genome Informatics (MGI) database. <sup>a</sup>Number of individuals with a complete knockout in a human ortholog of a gene with mouse phenotype is shown in parenthesis.



Class	N markers	MAF (%)			N het. parent	N minor transm.	Minor transm. frequency (SE)	N both parents het.	N double minor transm.	Double minor transm. frequency (SE)
		Thresh	Med	Mean						
LoF	2,741	<0.5	0.11	0.15	566,454	282,556	0.49882 (0.00068)	932	208	0.2232 (0.0151)
LoF	3,059	<1	0.12	0.20	867,035	432,649	0.49900 (0.00055)	2,205	504	0.2286 (0.0098)
LoF	3,235	<2	0.13	0.27	1,204,213	601,849	0.49979 (0.00046)	4,860	1,149	0.2364 (0.0064)
MODERATE	50,079	<0.5	0.12	0.16	11,146,218	5,570,282	0.49975 (0.00017)	19,383	4,772	0.2462 (0.0038)
MODERATE	57,943	<1	0.15	0.23	18,809,718	9,404,723	0.49999 (0.00013)	51,641	12,760	0.2471 (0.0024)
MODERATE	63,601	<2	0.17	0.34	29,899,791	14,951,803	0.50006 (0.00010)	140,233	34,860	0.2486 (0.0014)
Intergenic	3,849,601	<0.5	0.12	0.16	930,068,954	464,961,218	0.49992 (0.00007)	1,677,095	414,847	0.2474 (0.0019)
Integenic	4,571,969	<1	0.15	0.25	1,681,403,646	840,702,580	0.50000 (0.00006)	4,914,196	1,220,935	0.2485 (0.0012)
Intergenic	5,193,617	<2	0.19	0.39	2,986,623,050	1,493,458,184	0.50005 (0.00005)	15,552,151	3,878,138	0.2494 (0.0007)
<i>LoF:</i>										
RVIS 0-20% <sup>29</sup>	429	<2	0.11	0.20	115,940	57,783	0.49839 (0.00150)	290	52	0.1793 (0.0266)
RVIS 20-40%	380	<2	0.11	0.22	121,432	60,692	0.49980 (0.00146)	464	96	0.2069 (0.0214)
RVIS 40-60%	388	<2	0.13	0.28	155,663	77,726	0.49932 (0.00127)	752	211	0.2806 (0.0156)
RVIS 60-80%	529	<2	0.14	0.26	211,012	105,381	0.49941 (0.00112)	770	162	0.2104 (0.0167)
RVIS 80-100%	1,131	<2	0.15	0.31	487,101	243,595	0.50009 (0.00072)	2,184	542	0.2482 (0.0095)
Essential 0-20% <sup>30</sup>	150	<2	0.13	0.21	47,533	23,550	0.49545 (0.00231)	150	26	0.1733 (0.0376)
Essential 20-40%	214	<2	0.15	0.27	73,165	36,851	0.50367 (0.00184)	285	54	0.1895 (0.0274)
Essential 40-60%	220	<2	0.13	0.30	90,284	45,206	0.50071 (0.00169)	369	95	0.2575 (0.0226)
Essential 60-80%	215	<2	0.15	0.30	84,794	42,258	0.49836 (0.00174)	346	87	0.2514 (0.0241)
Essential 80-100%	203	<2	0.12	0.29	82,759	41,411	0.50038 (0.00175)	357	91	0.2549 (0.0221)
Severe recessive <sup>27</sup>	95	<2	0.13	0.21	29,112	14,607	0.50175 (0.00348)	72	14	0.1944 (0.0572)

**Supplementary Table 11.** Transmission from single heterozygous parents and two heterozygous parents. For the Residual Variation Intolerance Score (RVIS), essential, recessive and Kingsmore classes only LoF variants were considered. Genes were divided into five groups by their RVIS (available for 18,329 genes) and the essential score (available for 7,114 genes) percentiles.

## Supplementary References

1. Gudbjartsson, D.F. Large-scale whole-genome sequencing of the Icelandic population. *Nature Genetics* **In press**(2015).
2. Kong, A. *et al.* Detection of sharing by descent, long-range phasing and haplotype imputation. *Nat. Genet.* **40**, 1068-1075 (2008).
3. Kong, A. *et al.* Parental origin of sequence variants associated with complex diseases. *Nature* **462**, 868-874 (2009).
4. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).
5. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-303 (2010).
6. DePristo, M.A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet* **43**, 491-8 (2011).
7. Kong, A. *et al.* Detection of sharing by descent, long-range phasing and haplotype imputation. *Nat Genet* **40**, 1068-75 (2008).
8. Marchini, J., Howie, B., Myers, S., McVean, G. & Donnelly, P. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat Genet* **39**, 906-13 (2007).
9. Kong, A. *et al.* Fine-scale recombination rate differences between sexes, populations and individuals. *Nature* **467**, 1099-103 (2010).
10. Gudbjartsson, D. The Sequence of Icelanders. *Submitted to Nature Genetics* (2014).
11. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* **27**, 573-80 (1999).
12. Kent, W.J. *et al.* The human genome browser at UCSC. *Genome Res* **12**, 996-1006 (2002).
13. McLaren, W. *et al.* Deriving the consequences of genomic variants with the Ensembl API and SNP Effect Predictor. *Bioinformatics* **26**, 2069-70 (2010).
14. Flicek, P. *et al.* Ensembl 2012. *Nucleic Acids Res* **40**, D84-90 (2012).
15. Pruitt, K.D., Tatusova, T., Brown, G.R. & Maglott, D.R. NCBI Reference Sequences (RefSeq): current status, new features and genome annotation policy. *Nucleic Acids Res* **40**, D130-5 (2012).
16. Eilbeck, K. *et al.* The Sequence Ontology: a tool for the unification of genome annotations. *Genome Biol* **6**, R44 (2005).
17. Sherry, S.T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res* **29**, 308-11 (2001).
18. Nielsen, R., Paul, J.S., Albrechtsen, A. & Song, Y.S. Genotype and SNP calling from next-generation sequencing data. *Nat Rev Genet* **12**, 443-51 (2011).
19. Efron, B. & Tibshirani, R. *An introduction to the bootstrap*, xvi, 436 p. (Chapman & Hall, New York, 1993).
20. Trapnell, C. *et al.* Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* **7**, 562-78 (2012).
21. Trapnell, C., Pachter, L. & Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105-11 (2009).
22. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**, 511-5 (2010).
23. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-9 (2009).
24. Fagerberg, L. *et al.* Analysis of the human tissue-specific expression by genome-wide integration of transcriptomics and antibody-based proteomics. *Mol Cell Proteomics* **13**, 397-406 (2014).
25. Blake, J.A., Bult, C.J., Eppig, J.T., Kadin, J.A. & Richardson, J.E. The Mouse Genome Database: integration of and access to knowledge about the laboratory mouse. *Nucleic Acids Res* **42**, D810-7 (2014).

26. Kohler, S. *et al.* The Human Phenotype Ontology project: linking molecular biology and disease through phenotype data. *Nucleic Acids Res* **42**, D966-74 (2014).
27. Saunders, C.J. *et al.* Rapid whole-genome sequencing for genetic disease diagnosis in neonatal intensive care units. *Sci Transl Med* **4**, 154ra135 (2012).
28. Tennessen, J.A. *et al.* Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *Science* **337**, 64-9 (2012).
29. Petrovski, S., Wang, Q., Heinzen, E.L., Allen, A.S. & Goldstein, D.B. Genic intolerance to functional variation and the interpretation of personal genomes. *PLoS Genet* **9**, e1003709 (2013).
30. Wang, T., Wei, J.J., Sabatini, D.M. & Lander, E.S. Genetic screens in human cells using the CRISPR-Cas9 system. *Science* **343**, 80-4 (2014).