

# Bayesian estimation of gene constraint from an evolutionary model with gene features

**Tony Zeng**

Stanford University <https://orcid.org/0000-0002-6509-9879>

**Jeffrey Spence**

Stanford University <https://orcid.org/0000-0002-3199-1447>

**Hakhamanesh Mostafavi**

Stanford University <https://orcid.org/0000-0002-1060-2844>

**Jonathan Pritchard** (✉ [pritch@stanford.edu](mailto:pritch@stanford.edu))

Stanford University <https://orcid.org/0000-0002-8828-5236>

---

## Article

### Keywords:

**Posted Date:** June 13th, 2023

**DOI:** <https://doi.org/10.21203/rs.3.rs-3012879/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

**Additional Declarations:** There is **NO** Competing Interest.

---

# Bayesian estimation of gene constraint from an evolutionary model with gene features

Tony Zeng<sup>1,\*†</sup>, Jeffrey P. Spence<sup>1,\*†</sup>, Hakhamanesh Mostafavi<sup>1</sup>, Jonathan K. Pritchard<sup>1,2,†</sup>

<sup>1</sup> Department of Genetics, Stanford University, Stanford CA

<sup>2</sup> Department of Biology, Stanford University, Stanford CA

\* Equal contribution

† Correspondence to:

tkzeng@stanford.edu, jspence@stanford.edu, pritch@stanford.edu

June 9, 2023

## Abstract

Measures of selective constraint on genes have been used for many applications including clinical interpretation of rare coding variants, disease gene discovery, and studies of genome evolution. However, widely-used metrics are severely underpowered at detecting constraint for the shortest ~25% of genes, potentially causing important pathogenic mutations to be overlooked. We developed a framework combining a population genetics model with machine learning on gene features to enable accurate inference of an interpretable constraint metric,  $s_{\text{het}}$ . Our estimates outperform existing metrics for prioritizing genes important for cell essentiality, human disease, and other phenotypes, especially for short genes. Our new estimates of selective constraint should have wide utility for characterizing genes relevant to human disease. Finally, our inference framework, GeneBayes, provides a flexible platform that can improve estimation of many gene-level properties, such as rare variant burden or gene expression differences.

# 1 Introduction

Identifying the genes important for disease and fitness is a central goal in human genetics. One particularly useful measure of importance is how much natural selection constrains a gene [1–4]. Constraint has been used to prioritize *de novo* and rare variants for clinical followup [5,6], predict the toxicity of drugs [7], link GWAS hits to genes [8], and characterize transcriptional regulation [9,10], among many other applications.

To estimate the amount of constraint on a gene, several metrics have been developed using loss-of-function variants (LOFs), such as protein truncating or splice disrupting variants. If a gene is important, then natural selection will act to remove LOFs from the population. Several metrics of gene importance have been developed based on this intuition to take advantage of large exome sequencing studies.

In one line of research, the number of observed unique LOFs is compared to the expected number under a model of no selective constraint. This approach has led to the widely-used metrics pLI [11] and LOEUF [12].

While pLI and LOEUF have proved useful for identifying genes intolerant to LOF mutations, they have important limitations [3]. First, they are uninterpretable in that they are only loosely related to the fitness consequences of LOFs. Their relationship with natural selection depends on the study’s sample size and other technical factors [3]. Second, they are not based on an explicit population genetics model so it is impossible to compare a given value of pLI or LOEUF to the strength of selection estimated for variants other than LOFs [3,4].

Another line of research has solved these issues of interpretability by estimating the fitness reduction for heterozygous carriers of an LOF in any given gene [1,2,4]. Throughout, we will adopt the notation of Cassa and colleagues and refer to this reduction in fitness as  $s_{\text{het}}$  [1,2], although the same population genetic quantity has been referred to as *hs* [4,13]. In [1], a deterministic approximation was used to estimate  $s_{\text{het}}$ , which was relaxed to incorporate the effects of genetic drift in [2]. This model was subsequently extended by Agarwal and colleagues to include the X chromosome and applied to a larger dataset, with a focus on the interpretability of  $s_{\text{het}}$  [4].

A major issue for most previous methods is that thousands of genes have few expected unique LOFs under neutrality, as they have short protein-coding sequences. For example, there are >5,000 genes that cannot be called as constrained by LOEUF, as they have too few expected unique LOFs to fall under the recommended LOEUF cutoff of 0.35 [14]. This problem is not limited to LOEUF, however, and all of these methods are severely underpowered to detect selection for this ~25% of genes.

Here, we present an approach that can accurately estimate  $s_{\text{het}}$  even for genes with few expected LOFs, while maintaining the interpretability of previous population-genetics based estimates [1,2,4].

Our approach has two main technical innovations. First, we use a novel population genetics model of LOF allele frequencies. Previous methods have either only modeled the number of unique LOFs, throwing away frequency information [11,12,15], or considered the sum of LOF frequencies across the gene [1,2,4], an approach that is not robust to misannotated LOFs. In contrast, we model the frequencies of individual LOF variants, allowing us to not only use the information

59 in such frequencies but also to model the possibility that any given LOF variant has been misan-  
60 notated, making our estimates more robust. Our approach uses new computational machinery,  
61 described in a companion paper [16], to accurately obtain the likelihood of observing an LOF at a  
62 given frequency without resorting to simulation [2,4] or deterministic approximations [1].

63 Second, our approach uses thousands of gene features, including gene expression patterns,  
64 protein structure information, and evolutionary constraint, to improve estimates for genes with  
65 few expected LOFs. By using these features, we can share information across similar genes. In-  
66 tuitively, this allows us to improve estimates for genes with few expected LOFs by leveraging  
67 information from genes with similar features that do have sufficient LOF data.

68 Adopting a similar approach, a recent preprint [15] used gene features in a deep learning  
69 model to improve estimation of constraint for genes with few expected LOFs, but did not use an  
70 explicit population genetics model, resulting in the same issues with interpretability faced by pLI  
71 and LOEUF.

72 We applied our method to a large exome sequencing cohort [12]. Our estimates of  $s_{\text{het}}$  are  
73 substantially more predictive than previous metrics at prioritizing essential and disease-associated  
74 genes. We also interrogated the relationship between gene features and natural selection, finding  
75 that evolutionary conservation, protein structure, and expression patterns are more predictive of  
76  $s_{\text{het}}$  than co-expression and protein-protein interaction networks. Expression patterns in the brain  
77 and expression patterns during development are particularly predictive of  $s_{\text{het}}$ . Finally, we use  
78  $s_{\text{het}}$  to highlight differences in selection on different categories of genes and consider  $s_{\text{het}}$  in the  
79 context of selection on variants beyond LOFs.

80 Our approach, GeneBayes, is extremely flexible and can be applied to improve estimation of  
81 numerous gene properties beyond  $s_{\text{het}}$ . Our implementation is available at [https://github.com/  
82 tkzeng/GeneBayes](https://github.com/tkzeng/GeneBayes).

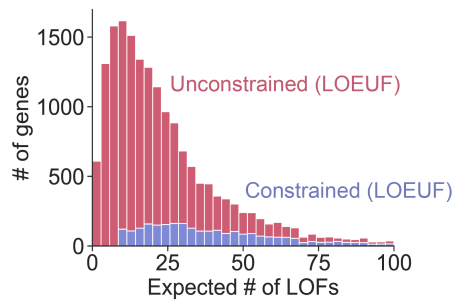
## 83 2 Results

### 84 2.1 Model Overview

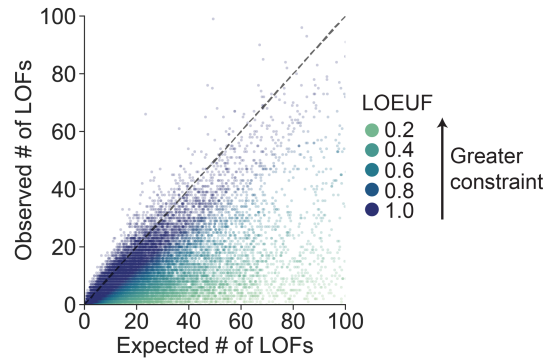
85 Using LOF data to infer gene constraint is challenging for genes with few expected LOFs, with  
86 metrics like LOEUF considering almost all such genes to be unconstrained (Figures 1A,B). We  
87 hypothesized that it would be possible to improve estimation using auxiliary information that  
88 may be predictive of LOF constraint, including gene expression patterns across tissues, protein  
89 structure, and evolutionary conservation. Intuitively, genes with similar features should have  
90 similar levels of constraint. By pooling information across groups of similar genes, constraint  
91 estimated for genes with sufficient LOF data may help improve estimation for underpowered  
92 genes.

93 However, while the frequencies of LOFs can be related to  $s_{\text{het}}$  through models from population  
94 genetics [1,2,4], we lack an understanding of how other gene features relate to constraint *a priori*.

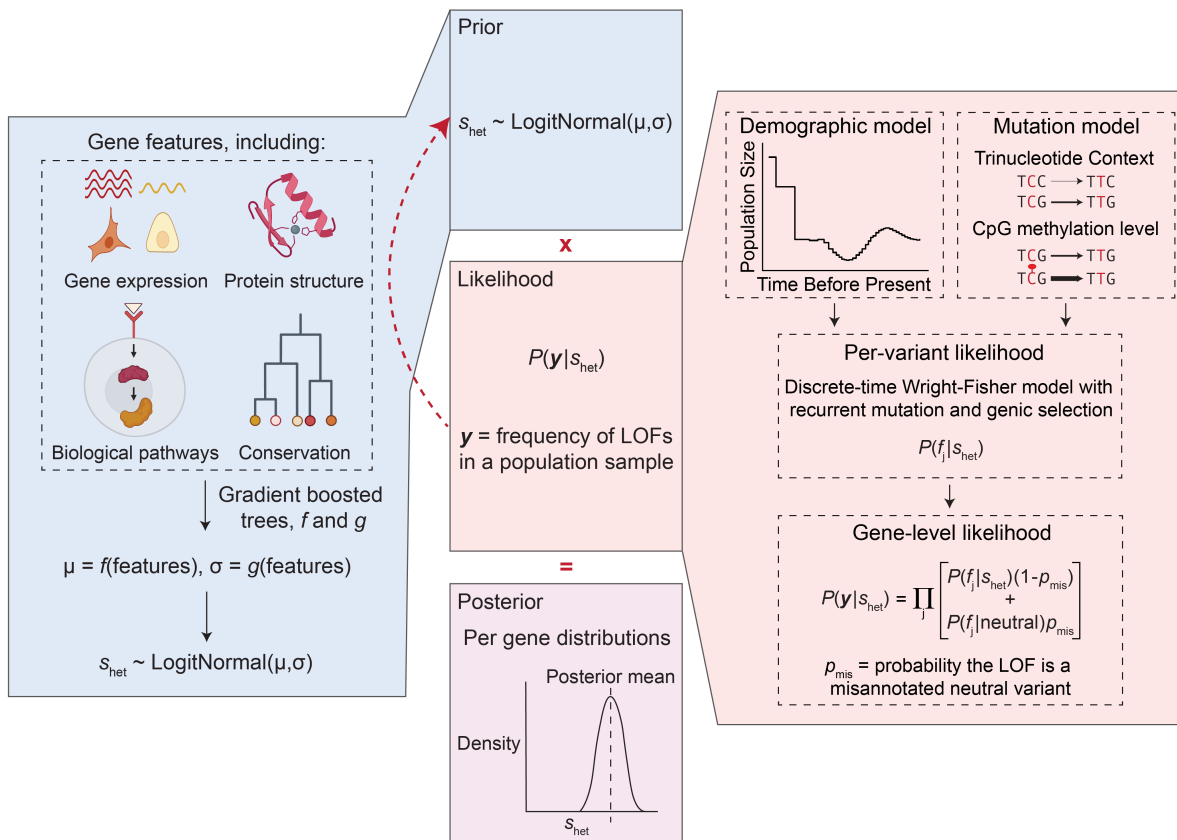
**A** Distribution of the number of expected LOFs



**B** LOEUF depends on the number of expected LOFs



**C** Schematic - estimating  $s_{\text{het}}$  using GeneBayes



**Figure 1: Limitations of LOEUF and schematic for inferring  $s_{\text{het}}$  using GeneBayes.** **A**) Stacked histogram of the expected number of unique LOFs per gene, where the distribution for genes considered unconstrained (respectively constrained) by LOEUF are colored in red (respectively blue). Genes with LOEUF < 0.35 are considered constrained, while all other genes are unconstrained (Methods). The plot is truncated on the x-axis at 100 expected LOFs. **B**) Scatterplot of the observed against the expected number of unique LOFs per gene. The dashed line denotes observed = expected. Each point is a gene, colored by its LOEUF score; genes with LOEUF > 1 are colored as LOEUF = 1. **C**) Schematic for estimating  $s_{\text{het}}$  using GeneBayes, highlighting the major components of the model: prior (blue boxes) and likelihood (red boxes). Parameters of the prior are learned by maximizing the likelihood (red arrow). Combining the prior and likelihood produces posteriors over  $s_{\text{het}}$  (purple box). See Methods for details.

95 To address this problem, we developed a flexible empirical Bayes framework, GeneBayes, that  
96 learns the relationship between gene features and  $s_{\text{het}}$  (Figure 1C). Our model consists of two main  
97 components. First, we model the prior on  $s_{\text{het}}$  for each gene as a function of its gene features (Fig-  
98 ure 1C, left). Specifically, we train gradient-boosted trees using NGBoost [17] to predict the param-  
99 eters of each gene’s prior distribution from its features. Our gene features include gene expression  
100 levels, Gene Ontology terms, conservation across species, neural network embeddings of pro-  
101 tein sequences, gene regulatory features, co-expression and protein-protein interaction features,  
102 sub-cellular localization, and intolerance to missense mutations (see Methods and Supplementary  
103 Note C for a full list).

104 Second, we use a model from population genetics to relate  $s_{\text{het}}$  to the observed LOF data (Fig-  
105 ure 1C, right). This model allows us to fit the gradient-boosted trees for the prior by maximizing  
106 the likelihood of the LOF data. Specifically, we use the discrete-time Wright Fisher model with  
107 genic selection, a standard model in population genetics that accounts for mutation and genetic  
108 drift [13, 18]. In our model,  $s_{\text{het}}$  is the reduction in fitness per copy of an LOF, and we infer  $s_{\text{het}}$   
109 while keeping the mutation rates and demography fixed to values taken from the literature (Sup-  
110 plementary Note B). Likelihoods are computed using new methods described in a companion  
111 paper [16].

112 Previous methods use either the number of *unique* LOFs or the sum of the frequencies of all  
113 LOFs in a gene, but we model the frequency of each individual LOF variant. We used LOF fre-  
114 quencies from the gnomAD consortium, which consists of exome sequences from  $\sim 125,000$  indi-  
115 viduals for 18,563 genes after filtering.

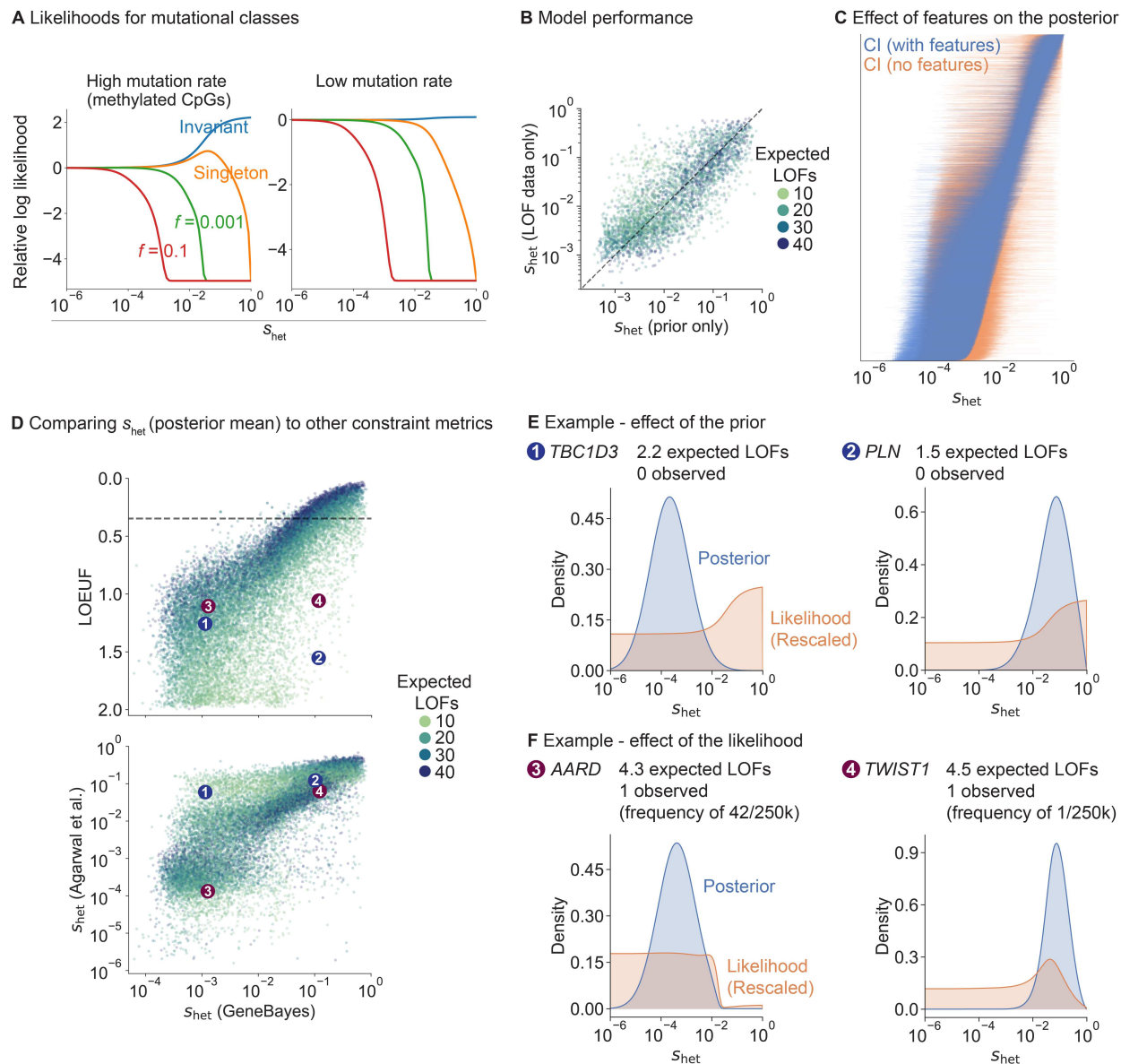
116 Combining these two components—the learned priors and the likelihood of the LOF data—we  
117 obtained posterior distributions over  $s_{\text{het}}$  for every gene. Throughout, we use the posterior mean  
118 value of  $s_{\text{het}}$  for each gene as a point estimate. See Methods for more details and Supplementary  
119 Table 2 for estimates of  $s_{\text{het}}$ .

## 120 2.2 Population genetics model and gene features both affect the estimation of $s_{\text{het}}$

121 First, we explored how LOF frequency and mutation rate relate to  $s_{\text{het}}$  in our population genet-  
122 ics model (Figure 2A). Invariant sites with high mutation rates are indicative of strong selection  
123 ( $s_{\text{het}} > 10^{-2}$ ), consistent with [19], while such sites with low mutation rates are consistent with  
124 essentially any value of  $s_{\text{het}}$  for the demographic model considered here. Regardless of mutation  
125 rate, singletons are consistent with most values of  $s_{\text{het}}$  but can rule out extremely strong selec-  
126 tion, and variants observed at a frequency of  $>10\%$  rule out even moderately strong selection  
127 ( $s_{\text{het}} > 10^{-3}$ ).

128 To assess how informative gene features are about  $s_{\text{het}}$ , we trained our model on a subset  
129 of genes and evaluated the model on held-out genes (Figure 2B, Methods). We computed the  
130 Spearman correlation between  $s_{\text{het}}$  estimates from the prior and  $s_{\text{het}}$  estimated from the LOF data  
131 only. The correlation is high and comparable between train and test sets (Spearman  $\rho = 0.83$  and  
132  $0.78$  respectively), indicating the gene features alone are highly predictive of  $s_{\text{het}}$  and that this is  
133 not a consequence of overfitting.

134 To further characterize the impact of features on our estimates of  $s_{\text{het}}$ , we removed all features  
135 from our model and recalculated posterior distributions (Figure 2C). For most genes, posteriors



**Figure 2: Factors that contribute to our estimates of  $s_{het}$ .** **A)** Likelihood curves for different allele frequencies ( $f$ ) and mutation rates. **B)** Scatterplot of  $s_{het}$  estimated from LOF data (y-axis; posterior mean from a model without features) against the prior's predictions of  $s_{het}$  (x-axis; mean of learned prior). Dotted line denotes  $y = x$ . Each point is a gene, colored by the expected number of LOFs. **C)** Comparison of posterior distributions of  $s_{het}$  (95% Credible Intervals) from a model with (blue lines) and without (orange lines) gene features. Genes are ordered by their posterior mean in the model with gene features. **D)** Top: scatterplot of LOEUF (y-axis) and our  $s_{het}$  estimates (x-axis; posterior mean). Each point is a gene, colored by the expected number of LOFs. Bottom: scatterplot of  $s_{het}$  estimates from [4] (y-axis; posterior mode) and our  $s_{het}$  estimates (x-axis; posterior mean). Numbered points refer to genes in panels E and F. **E)** *TBC1D3* and *PLN* are two example genes where the gene features substantially affect the posterior. We plot their posterior distributions (blue) and likelihoods (orange; rescaled so that the area under the curve = 1). **F)** *AARD* and *TWIST1* are two example genes with the same LOEUF but different  $s_{het}$ . Posteriors and likelihoods are plotted as in panel E.

136 are substantially more concentrated when using gene features.

137 Next, we compared our estimates of  $s_{\text{het}}$  using GeneBayes to LOEUF and to selection coeffi-  
138 cients estimated by [4] (Figure 2D). To facilitate comparison, we use the posterior modes of  $s_{\text{het}}$   
139 reported in [4] as point estimates, but we note that [4] emphasizes the value of using full posterior  
140 distributions. While the correlation between our estimates is high for genes with sufficient LOFs  
141 (for genes with more LOFs than the median, Spearman  $\rho$  with LOEUF = 0.94;  $\rho$  with  $s_{\text{het}}$  from [4]  
142 = 0.88), it is lower for genes with few expected LOFs (for genes with fewer LOFs than the median,  
143 Spearman  $\rho$  with LOEUF = 0.71;  $\rho$  with  $s_{\text{het}}$  from [4] = 0.71).

144 We further explored the reduced correlations for genes with few expected LOFs. For example,  
145 *TBC1D3* and *PLN* have few expected LOFs, and their likelihoods are consistent with any level  
146 of constraint (Figure 2E). Due to the high degree of uncertainty, LOEUF considers both genes to  
147 be unconstrained, while the  $s_{\text{het}}$  point estimates from [4] err in the other direction and consider  
148 both genes to be constrained (Figure 2D). This uncertainty arises from use of the LOF data alone,  
149 and is captured by the wide posterior distributions for the  $s_{\text{het}}$  estimates from [4]. In contrast, by  
150 using gene features, our posterior distributions of  $s_{\text{het}}$  indicate that *PLN* is strongly constrained  
151 but *TBC1D3* is not, consistent with the observation that heterozygous LOFs in *PLN* cause severe  
152 cardiac dilation and heart failure [20].

153 In contrast to estimates of  $s_{\text{het}}$ , LOEUF further ignores information about allele frequencies by  
154 considering only the number of unique LOFs, resulting in a loss of information. For example,  
155 *AARD* and *TWIST1* have almost the same numbers of observed and expected unique LOFs, so  
156 LOEUF is similar for both (LOEUF = 1.1 and 1.06 respectively). However, while *TWIST1*'s ob-  
157 served LOF is present in only 1 of 246,192 alleles, *AARD*'s is  $\sim 40\times$  more frequent. Consequently,  
158 the likelihood rules out the possibility of strong constraint at *AARD* (Figure 2F), causing the two  
159 genes to differ in their estimated selection coefficients (Figure 2D).

160 In contrast, *TWIST1* has a posterior mean  $s_{\text{het}}$  of 0.11 when using gene features, indicating very  
161 strong selection. Consistent with this, *TWIST1* is a transcription factor critical for specification of  
162 the cranial mesoderm, and heterozygous LOFs in the gene are associated with Saethre-Chotzen  
163 syndrome, a disorder characterized by congenital skull and limb abnormalities [21,22].

164 Besides *PLN* and *TWIST1*, many genes are considered constrained by  $s_{\text{het}}$  but not by LOEUF,  
165 which is designed to be highly conservative. In Table 1, we list 15 examples with  $s_{\text{het}} > 0.1$   
166 and LOEUF  $> 0.5$ , selected based on their clinical significance and prominence in the literature  
167 (Methods). One notable example is a set of 16 ribosomal protein genes for which heterozygous  
168 disruption causes Diamond-Blackfan anemia—a rare genetic disorder characterized by an inabil-  
169 ity to produce red blood cells [23] (Supplementary Table 1). All are considered strongly con-  
170 strained by  $s_{\text{het}}$  (minimum  $s_{\text{het}} = 0.26$ ). In contrast, only 6 are considered constrained by LOEUF  
171 (LOEUF  $< 0.35$ ), as many of these genes have few expected unique LOFs.



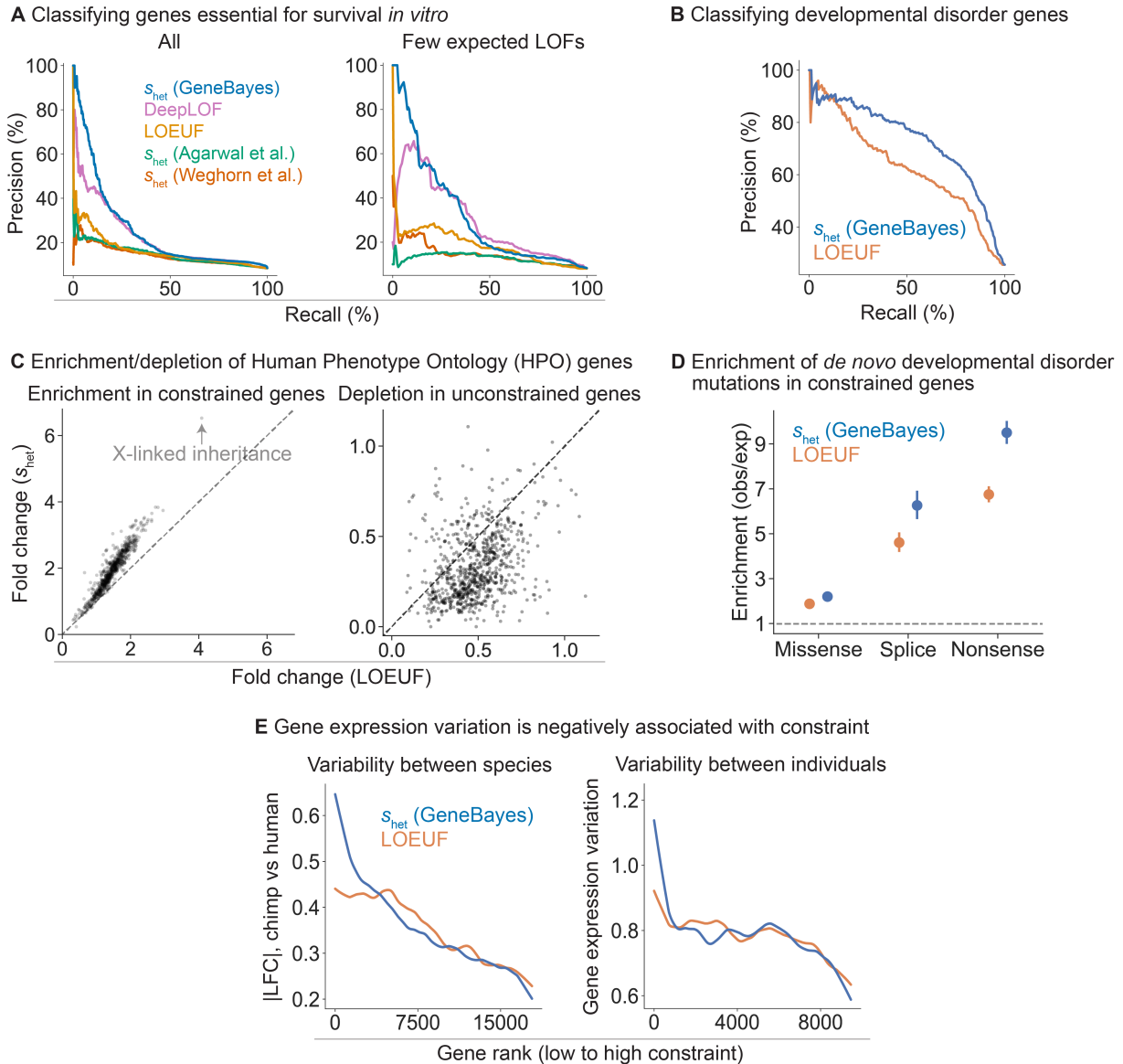
Gene	$s_{\text{het}}$	LOEUF	Obs.	Exp.	Condition and reference
<i>RPS15A</i> *	0.61	0.56	0	5.4	<i>Diamond-Blackfan anemia</i> : Red blood cell aplasia resulting in growth, craniofacial, and other congenital defects [23]
<i>DCX</i>	0.48	0.62	3	12.6	<i>Lissencephaly</i> : Migrational arrest of neurons resulting in mental retardation and seizures [24]
<i>SOX2</i>	0.33	0.57	1	8.3	<i>Syndromic microphthalmia</i> : Missing or small eyes from birth [25]
<i>NDP</i>	0.33	0.88	0	3.4	<i>Norrie disease</i> : Retinal dystrophy resulting in early childhood blindness, mental disorders, and deafness [26]
<i>EIF5A</i>	0.32	0.54	1	8.7	<i>Faundes-Banka syndrome</i> : Developmental delay, microcephaly, and facial dysmorphisms [27]
<i>CDKN1C</i>	0.27	0.53	0	5.7	<i>Beckwith-Wiedemann syndrome</i> : Pediatric overgrowth with predisposition to tumor development [28]
<i>TGIF1</i>	0.25	0.91	5	11.5	<i>Holoprosencephaly</i> : Structural malformation of the forebrain during development [29]
<i>SH2D1A</i>	0.23	0.96	1	4.9	<i>Lymphoproliferative syndrome</i> : Severe immune dysregulation due to improper lymphocyte apoptosis [30]
<i>CEBPA</i>	0.17	1.18	0	2.4	<i>Acute myeloid leukemia</i> : Blood and bone marrow cancer with rapid progression [31]
<i>GATA4</i>	0.15	0.53	3	14.7	<i>Atrial septal defect</i> : Congenital heart defect resulting in a hole between the atria [32]
<i>TIMP3</i>	0.13	0.53	2	11.8	<i>Sorsby fundus dystrophy</i> : Retinal dystrophy that causes loss of vision [33]
<i>FOXC2</i>	0.13	0.79	3	9.8	<i>Lymphedema-distichiasis syndrome</i> : Lymphedema of the limbs and double rows of eyelashes [34]
<i>IGF2</i>	0.12	1.13	3	6.8	<i>Silver-Russell syndrome</i> : Growth retardation, relative macrocephaly, and feeding difficulties [35]
<i>PLN</i>	0.12	1.56	0	1.5	<i>Dilated cardiomyopathy</i> : Enlarged heart chambers, decreased contractile function, and heart failure [20]
<i>TWIST1</i>	0.11	1.06	1	4.5	<i>Saethre-Chotzen syndrome</i> : Craniosynostosis, facial dysmorphism, and hand and foot abnormalities [21] [22]

Table 1: **OMIM genes constrained by  $s_{\text{het}}$  but not by LOEUF.** Mutations that disrupt the functions of these genes are associated with Mendelian diseases in the OMIM database [36]. Genes are ordered by  $s_{\text{het}}$  (posterior mean). Obs. and Exp. are the unique number of observed and expected LOFs respectively. \**RPS15A* is associated with *Diamond-Blackfan anemia* along with nine other genes considered constrained by  $s_{\text{het}}$  but not by LOEUF (Supplementary Table 1).

### 172 2.3 Utility of $s_{\text{het}}$ in prioritizing phenotypically important genes

173 To assess the accuracy of our  $s_{\text{het}}$  estimates and evaluate their ability to prioritize genes, we first  
174 used these estimates to classify genes essential for survival of human cells *in vitro*. Genome-wide  
175 CRISPR growth screens have measured the effects of gene knockouts on cell survival or prolif-  
176 eration, quantifying the *in vitro* importance of each gene for fitness [37, 38]. We find that our  
177 estimates of  $s_{\text{het}}$  outperform other constraint metrics at classifying essential genes (Figure 3A, left;  
178 bootstrap  $p < 2 \times 10^{-5}$  for pairwise differences in AUPRC between our estimates and other met-  
179 rics). The difference is largest for genes with few expected LOFs, where  $s_{\text{het}}$  (GeneBayes) retains  
180 similar precision and recall while other metrics lose performance (Figure 3A, right). In addition,  
181 our estimates of  $s_{\text{het}}$  outperform other metrics at classifying nonessential genes (Supplementary  
182 Figure 2A).

183 DeepLOF [15], the only other method that combines information from both LOF data and gene



**Figure 3: GeneBayes estimates of  $s_{het}$  perform well at identifying constrained and unconstrained genes.** **A)** Precision-recall curves comparing the performance of  $s_{het}$  against other methods in classifying essential genes (left: all genes, right: quartile of genes with the fewest expected unique LOFs). **B)** Precision-recall curves comparing the performance of  $s_{het}$  against LOEUF in classifying developmental disorder genes. **C)** Scatterplots showing the enrichment (respectively depletion) of the top 10% most (respectively least) constrained genes in HPO terms, with genes ranked by  $s_{het}$  (y-axis) or LOEUF (x-axis). **D)** Enrichment of *de novo* mutations in patients with developmental disorders, calculated as the observed number of mutations over the expected number under a null mutational model. We plot the enrichment of missense, splice, and nonsense variants in the 10% most constrained genes, ranked by  $s_{het}$  (blue) or LOEUF (orange). Bars represent 95% confidence intervals. **E)** Left: LOESS curve showing the relationship between constraint (gene rank, x-axis) and absolute log fold change in expression between chimp and human cortical cells (y-axis). Genes are ranked by  $s_{het}$  (blue) or LOEUF (orange) Right: LOESS curve showing the relationship between constraint (gene rank, x-axis) and gene expression variation (normalized standard deviation) in GTEx samples.

184 features, outperforms methods that rely exclusively on LOF data, highlighting the importance of

185 using auxiliary information. Yet, DeepLOF uses only the number of unique LOFs, discarding  
186 frequency information. As a result, it is outperformed by our method, indicating that careful  
187 modeling of LOF frequencies also contributes to the performance of our approach.

188 Next, we performed further comparisons of our estimates of  $s_{\text{het}}$  against LOEUF, as LOEUF  
189 and its predecessor pLI are extremely popular metrics of constraint. To evaluate the ability of  
190 these methods to prioritize disease genes, we first used  $s_{\text{het}}$  and LOEUF to classify curated devel-  
191 opmental disorder genes [39]. Here,  $s_{\text{het}}$  outperforms LOEUF (Figure 3B; bootstrap  $p = 2 \times 10^{-9}$   
192 for the difference in AUPRC) and performs favorably compared to additional constraint metrics  
193 (Supplementary Figure 2B).

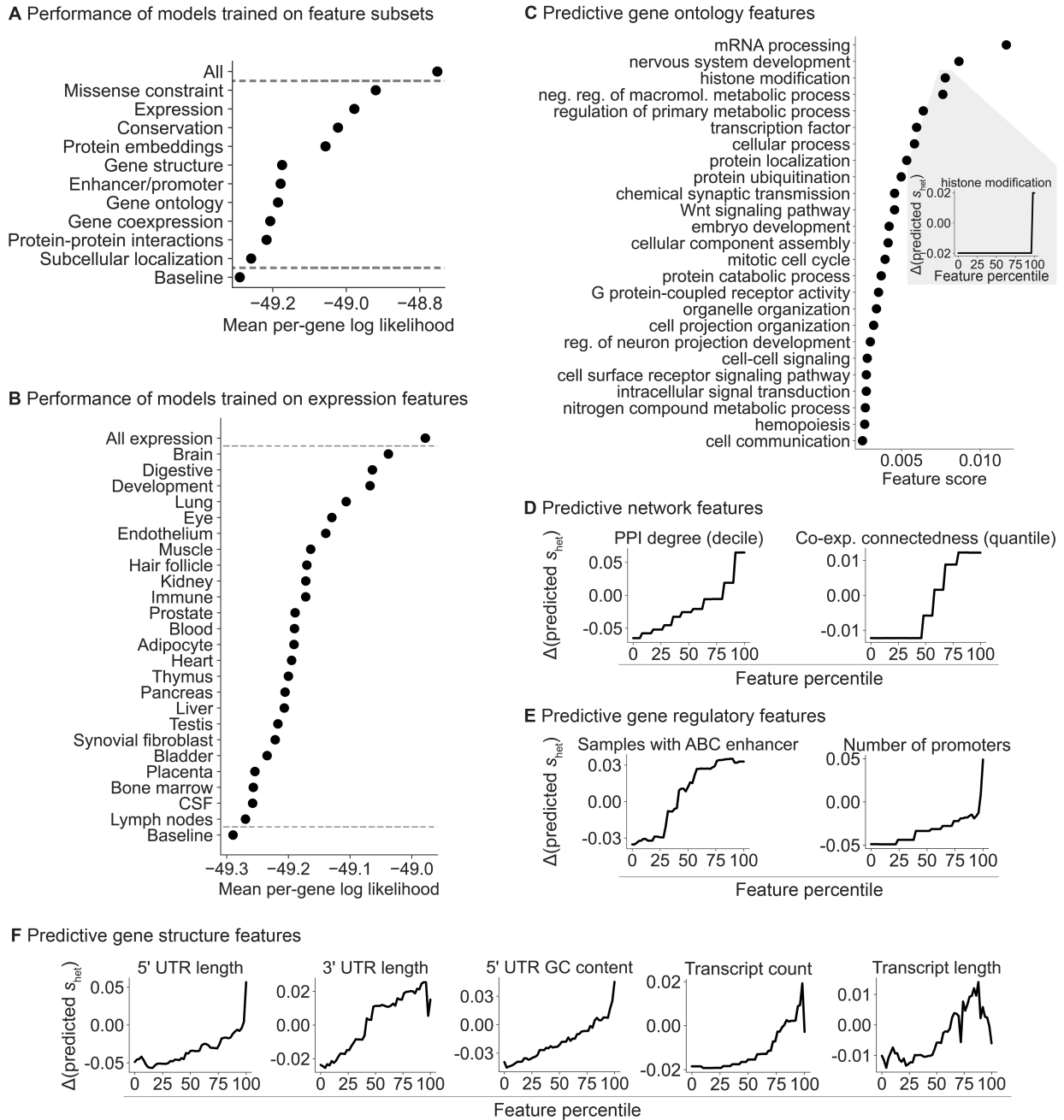
194 Next, we considered a broader range of phenotypic abnormalities annotated in the Human  
195 Phenotype Ontology (HPO) [40]. For each HPO term, we calculated the enrichment of the 10%  
196 most constrained genes and depletion of the 10% least constrained genes, ranked using  $s_{\text{het}}$  or  
197 LOEUF. Genes considered constrained by  $s_{\text{het}}$  are 1.9-fold enriched in HPO terms, compared to  
198 1.5-fold enrichment for genes considered constrained by LOEUF (Figure 3C, left). Additionally,  
199 genes considered unconstrained by  $s_{\text{het}}$  are 3.0-fold depleted in HPO terms, compared to 2.1-fold  
200 depletion for genes considered constrained by LOEUF (Figure 3C, right).

201 X-linked inheritance is one of the terms with the largest enrichment of constrained genes (6.6-  
202 fold enrichment for  $s_{\text{het}}$  and 4.2-fold enrichment for LOEUF). The ability of  $s_{\text{het}}$  to prioritize X-  
203 linked genes may prove particularly useful, as many disorders are enriched for X-chromosome  
204 genes [41] and the selection on losing a single copy of such genes is stronger on average [4].  
205 Yet, population-scale sequencing alone has less power to detect a given level of constraint on  
206 X-chromosome genes, as the number of X chromosomes in a cohort with males is smaller than the  
207 number of autosomes.

208 We next assessed if *de novo* disease-associated variants are enriched in constrained genes, simi-  
209 lar to the analyses in [4,5]. To this end, we used data from 31,058 trios to calculate for each gene the  
210 enrichment of *de novo* missense and LOF mutations in offspring with DDs relative to unaffected  
211 parents [5]. We found that for both classes of variants, enrichment is higher for genes considered  
212 constrained by  $s_{\text{het}}$ , with the highest enrichment observed for LOF variants (Figure 3D; enrich-  
213 ment of  $s_{\text{het}}$  and LOEUF respectively, for missense mutations = 2.2, 1.9; splice site mutations =  
214 6.3, 4.6; and nonsense mutations = 9.5, 6.7). Consistent with previous findings, the excess burden  
215 of *de novo* variants is predominantly in highly constrained genes (Supplementary Figure 2C, left).  
216 Notably, this difference in enrichment remains after removing known DD genes (Supplementary  
217 Figure 2C, right). Together, these results indicate that  $s_{\text{het}}$  not only improves identification of  
218 known disease genes but may also facilitate discovery of novel DD genes [5].

219 Finally, constraint can also be related to longer-term evolutionary processes that give rise to the  
220 variation among individuals or species, including variation in gene expression levels. We expect  
221 constrained genes to maintain expression levels closer to their optimal values across evolutionary  
222 time scales, as each LOF can be thought of as a  $\sim 50\%$  reduction in expression. Consistent with  
223 this expectation, we find that less constrained genes have larger absolute differences in expression  
224 between human and chimpanzee in cortical cells [42], with a stronger correlation for  $s_{\text{het}}$  than for  
225 LOEUF (Figure 3E). This pattern should also hold when considering the variation in expression  
226 within a species. We quantified variance using the normalized standard deviation of gene expres-  
227 sion levels estimated from RNA-seq samples in GTEx [43] and found that the variance decreases

228 with increased constraint, again with a stronger correlation for  $s_{\text{het}}$  (Figure 3E).



**Figure 4: Breakdown of the gene features important for  $s_{\text{het}}$  prediction.** **A)** Ordered from highest to lowest, plot of the mean per-gene log likelihood over the test genes for models separately trained on categories of features. “All” and “Baseline” include all and no features respectively. **B)** Plot of the mean per-gene log likelihood, as in panel A, for models separately trained on expression features grouped by tissue, cell type, or developmental stage. **C)** Ordered from highest to lowest, feature scores for individual gene ontology (GO) terms. Inset: lineplot showing the change in predicted  $s_{\text{het}}$  for a feature as the feature value is varied. **D)** Lineplot as in panel C (inset) for protein-protein interaction (PPI) and co-expression features, **E)** enhancer and promoter features, and **F)** gene structure features.

## 2.4 Interpreting the learned relationship between gene features and $s_{\text{het}}$

Our framework allows us to learn the relationship between gene features and  $s_{\text{het}}$  in a statistically principled way. In particular, by fitting a model with all of the features jointly, we can account for dependencies between the features. To interrogate the relationship between features and  $s_{\text{het}}$ , we divided our gene features into 10 distinct categories (Figure 4A) and trained a separate model per category using only the features in that category. We found that missense constraint, gene expression patterns, evolutionary conservation, and protein embeddings are the most informative categories.

Next, we further divided the expression features into 24 subgroups, representing tissues, cell types, and developmental stage (Table 6). Expression patterns in the brain, digestive system, and during development are the most predictive of constraint (Figure 4B). Notably, a study that matched Mendelian disorders to tissues through literature review found that a sizable plurality affect the brain [44]. Meanwhile, most of the top digestive expression features are also related to development (e.g., expression component loadings in a fetal digestive dataset [45]). The importance of developmental features is consistent with the severity of many developmental disorders and the expectation that selection is stronger on early-onset phenotypes [46], supported by the findings of [4].

To quantify the relationship between constraint and individual features, we changed the value of one feature at a time and used the variation in predicted  $s_{\text{het}}$  over the feature values as the score for each feature (Methods).

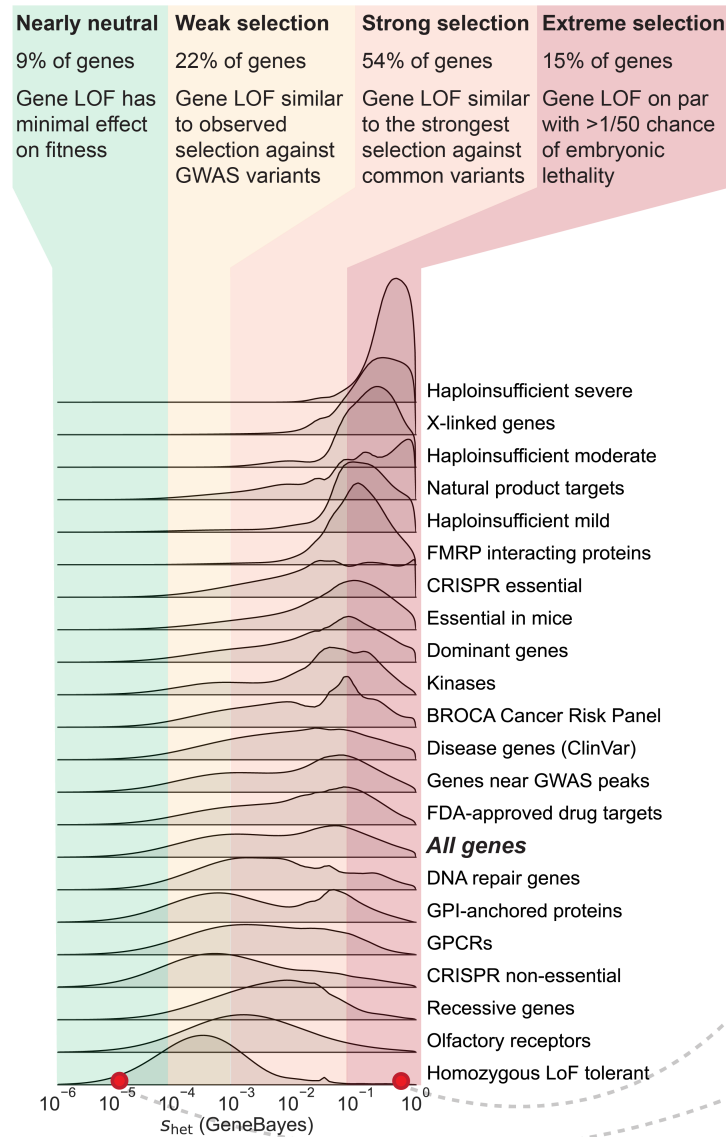
We first explored some of the individual Gene Ontology (GO) terms most predictive of constraint (Figure 4C). Consistent with the top expression features, the top GO features highlight developmental and brain-specific processes as important for selection.

Next, we analyzed network (Figure 4D), gene regulatory (Figure 4E), and gene structure (Figure 4F) features. Protein-protein interaction (PPI) and gene co-expression networks have highlighted “hub” genes involved in numerous cellular processes [47,48], while genes linked to GWAS variants have more complex enhancer landscapes [49]. Consistent with these studies, we find that connectedness in PPI and co-expression networks as well as enhancer and promoter count are positively associated with constraint (Figure 4D,E). In addition, gene structure affects gene function—for example, UTR length and GC content affect RNA stability, translation, and localization [50,51]—and likewise, several gene structure features are predictive of constraint (Figure 4F). Our results indicate that more complex genes—genes that are involved in more regulatory connections, that are more central to networks, and that have more complex gene structures—are generally more constrained.

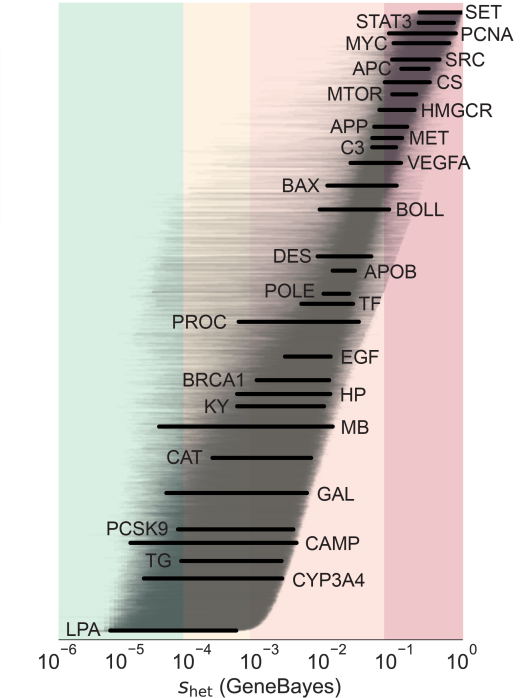
## 2.5 Contextualizing the strength of selection against gene loss-of-function

A major benefit of  $s_{\text{het}}$  over LOEUF and pLI is that  $s_{\text{het}}$  has a precise, intrinsic meaning in terms of fitness [1–4]. This facilitates comparison of  $s_{\text{het}}$  between genes, populations, species, and studies. For example,  $s_{\text{het}}$  can be compared to selection estimated from mutation accumulation or gene deletion experiments performed in model organisms [52,53]. More broadly, selection applies beyond LOFs. While we focused on estimating changes in fitness due to LOFs, consequences of

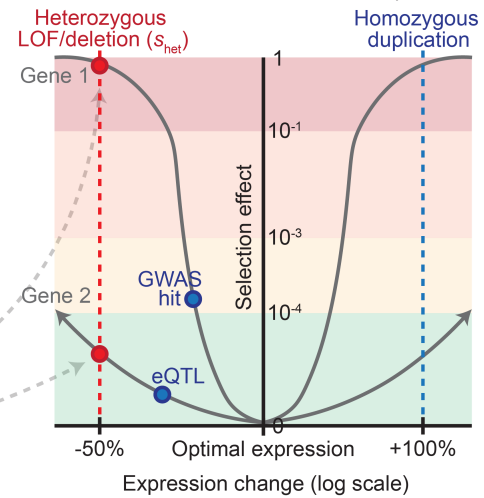
**A**  $s_{\text{het}}$  distributions for categories of genes



**B**  $s_{\text{het}}$  distributions for individual genes



**C** Model for selection as a function of expression



**Figure 5: Comparing selection on LOFs ( $s_{\text{het}}$ ) between genes and to selection on other variant types. **A)** Distributions of  $s_{\text{het}}$  for gene sets, calculated by averaging the posterior distributions for the genes in each gene set. Gene sets are sorted by the mean of their distributions. Colors represent four general selection regimes. **B)** Posterior distributions of  $s_{\text{het}}$  for individual genes, ordered by mean. Lines represent 95% credible intervals, with labeled genes represented by thick black lines. Colors represent the selection regimes in panel A. **C)** Schematic demonstrating the hypothesized relationship between changes in expression (x-axis,  $\log_2$  scale) and selection (y-axis) against these changes for two hypothetical genes, assuming stabilizing selection. The shapes of the curves are not estimated from real data. Background colors represent the selection regimes in panel A. The red points and line represent the effects of heterozygous LOFs and deletions on expression and selection, while the blue points and line represent the potential effects of other types of variants.**

269 non-coding, missense, and copy number variants can be understood through the same framework,  
270 as we expect such variants to also be under negative selection [19] due to ubiquitous stabilizing  
271 selection on traits [54]. Quantifying differences in the selection on variants will deepen our under-  
272 standing of the evolution and genetics of human traits (see Discussion).

273 To contextualize our  $s_{\text{het}}$  estimates, we compared the distributions of  $s_{\text{het}}$  for different gene sets  
274 (Figure 5A) and genes (Figure 5B), and analyzed them in terms of selection regimes. To define such  
275 regimes, we first conceptualized selection on variants as a function of their effects on expression  
276 (Figure 5C), where heterozygous LOFs reduce expression by  $\sim 50\%$  across all contexts relevant to  
277 selection. Under this framework, we can directly compare  $s_{\text{het}}$  to selection on other variant types—  
278 for the hypothetical genes in Figure 5C, a GWAS hit affecting Gene 1 has a stronger selective effect  
279 than a LOF affecting Gene 2, despite having a smaller effect on expression.

280 Next, we divided the range of possible  $s_{\text{het}}$  values into four regimes determined by theoretical  
281 considerations [55] and comparisons to other types of variants [56, 57]—nearly neutral (9% of  
282 genes), weak selection (22%), strong selection (54%), and extreme selection (15%). LOFs in nearly  
283 neutral genes ( $s_{\text{het}} < 10^{-4}$ ) have minimal effects on fitness—the frequency of such variants is  
284 dominated by genetic drift rather than selection [55]. Under the weak selection regime ( $s_{\text{het}}$  from  
285  $10^{-4}$  to  $10^{-3}$ ), gene LOFs have similar effects on fitness as typical GWAS hits, which usually have  
286 small or context-specific effects on gene expression or function [56]. Under the strong selection  
287 regime ( $s_{\text{het}}$  from  $10^{-3}$  to  $10^{-1}$ ), gene LOFs have fitness effects on par with the strongest selection  
288 coefficients measured for common variants, such as the selection estimated for adaptive mutations  
289 in *LCT* [57]. Finally, for genes in the extreme selection regime ( $s_{\text{het}} > 10^{-1}$ ), LOFs have an effect  
290 on fitness equivalent to a  $>2\%$  chance of embryonic lethality, indicating that such LOFs have an  
291 extreme effect on survival or reproduction.

292 Gene sets vary widely in their constraint. For example, genes known to be haploinsufficient  
293 for severe diseases are almost all under extreme selection. In contrast, genes that can tolerate  
294 homozygous LOFs are generally under weak selection. One notable example of such a gene is  
295 *LPA*—while high expression levels are associated with cardiovascular disease, low levels have  
296 minimal phenotypic consequences [58, 59], consistent with limited conservation in the sequence  
297 or gene expression of *LPA* across species and populations [60, 61]

298 Other gene sets have much broader distributions of  $s_{\text{het}}$  values. For example, manually curated  
299 recessive genes are under weak to strong selection, indicating that many such genes are either not  
300 fully recessive or have pleiotropic effects on other traits under selection. For example, homozy-  
301 gous LOFs in *PROC* can cause life-threatening congenital blood clotting [62], yet  $s_{\text{het}}$  for *PROC* is  
302 non-negligible (Figure 5B), consistent with observations that heterozygous LOFs can also increase  
303 blood clotting and cause deep vein thrombosis [63].

304 Similarly,  $s_{\text{het}}$  values for ClinVar disease genes [64] span the range from weak to extreme se-  
305 lection, with only moderate enrichment for greater constraint relative to all genes. Consistent  
306 with this, the effects of disease on fitness depend on disease severity, age-of-onset, and preva-  
307 lence throughout human history. For example, even though heterozygous loss of *BRCA1* greatly  
308 increases risk of breast and ovarian cancer [65], *BRCA1* is under strong rather than extreme se-  
309 lection. Possible partial explanations are that these cancers have an age-of-onset past reproduc-  
310 tive age and are less prevalent in males, or that *BRCA1* is subject to some form of antagonistic  
311 pleiotropy [14, 66].

312 **3 Discussion**

**Figure 6** Schematic of GeneBayes with example applications

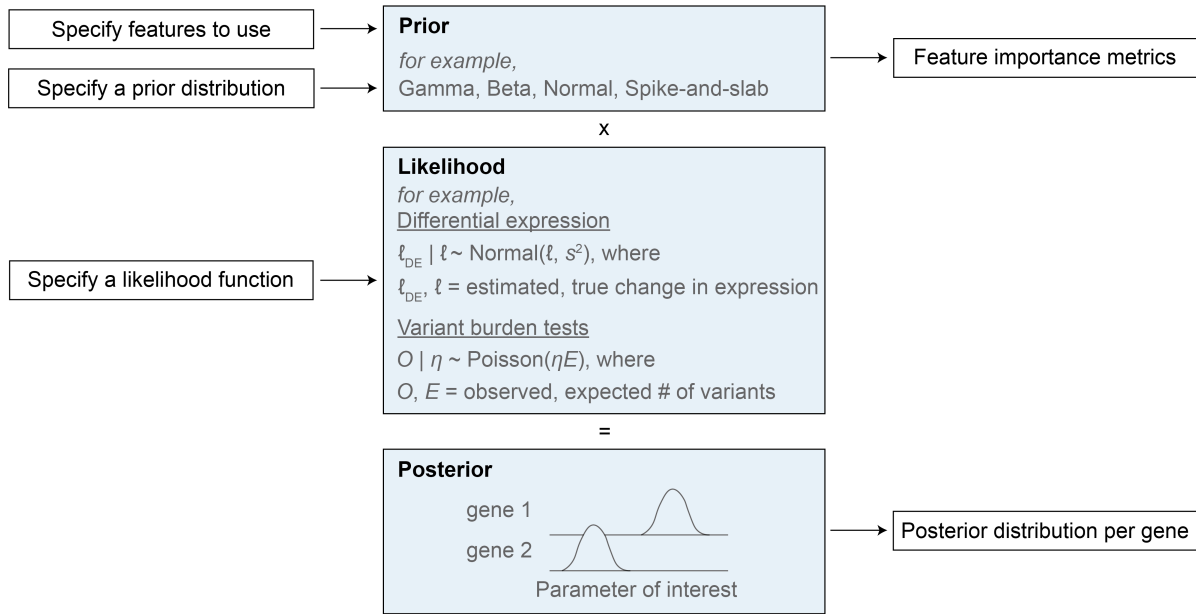


Figure 6: **GeneBayes is a flexible framework for estimating gene-level properties.** Schematic for how GeneBayes can be applied to estimate gene-level properties beyond  $s_{het}$ , showing the key inputs and outputs and two example applications. See Supplementary Note D for more details.

313 Here, we developed an empirical Bayes approach to accurately infer  $s_{het}$ , an interpretable met-  
 314 ric of gene constraint. Our approach uses powerful machine learning methods to leverage vast  
 315 amounts of functional and evolutionary information about each gene while coupling them to a  
 316 population genetics model.

317 There are two advantages of this approach. First, the additional data sources result in substan-  
 318 tially better performance than LOEUF across tasks, from classifying essential genes to identifying  
 319 pathogenic *de novo* mutations. These improvements are especially pronounced for the large frac-  
 320 tion of genes with few expected LOFs, where LOF data alone is underpowered for estimating  
 321 constraint.

322 Second, by inferring  $s_{het}$ , our estimates of constraint are interpretable in terms of fitness, and  
 323 we can directly compare the impact of a loss-of-function across genes, populations, species, and  
 324 studies.

325 As a selection coefficient,  $s_{het}$  can also be directly compared to other selection coefficients, even  
 326 for different types of variants [3,4]. In general, we believe genes are close to their optimal levels  
 327 of expression and experience stabilizing selection [54], in which case expression-altering variants  
 328 decrease fitness, with larger perturbations causing greater decreases (Figure 5C). Estimating the  
 329 fitness consequences of other types of expression-altering variants, such as duplications or eQTLs,  
 330 will allow us to map the relationship between genetic variation and fitness in detail, deepening  
 331 our understanding of the interplay of expression, complex traits, and fitness [10,56,67,68].



332 A recent method, DeepLOF [15], uses a similar empirical Bayes approach, but by estimating  
333 constraint from the number of observed and expected unique LOFs, it inherits the same difficul-  
334 ties regarding interpretation as pLI and LOEUF, and loses information by not considering variant  
335 frequencies. On the other hand, another line of work [1,2], culminating in [4], solved the issues  
336 with interpretability by directly estimating  $s_{\text{het}}$ . Yet, by relying exclusively on LOFs, these esti-  
337 mates are underpowered for  $\sim 25\%$  of genes. Furthermore, by using the aggregate frequencies of  
338 all LOF variants, previous  $s_{\text{het}}$  estimates [1,2,4] are not robust to misannotated LOF variants. Our  
339 approach eliminates this tradeoff between power and interpretability present in existing metrics.

340 Our estimates of  $s_{\text{het}}$  will be useful for many applications. For example, by informing gene-  
341 level priors, LOEUF, pLI, and previous estimates of  $s_{\text{het}}$  have been used to increase the power of  
342 association studies based on rare or *de novo* mutations [5,6,69]. In such contexts, our  $s_{\text{het}}$  estimates  
343 can be used as a drop-in replacement. Additionally, extremely constrained and unconstrained  
344 genes may be interesting to study in their own right. Genes of unknown function with particularly  
345 high values of  $s_{\text{het}}$  should be prioritized for further study. Investigating highly constrained genes  
346 may give insights into the mechanisms by which cellular and organism-level phenotypes affect  
347 fitness [70].

348 While we primarily used the posterior means of  $s_{\text{het}}$  here, our approach provides the entire  
349 posterior distribution per gene, similar to [4]. In some applications, different aspects of the pos-  
350 terior may be more relevant than the mean. For example, when prioritizing rare variants for  
351 followup in a clinical setting, the posterior probability that  $s_{\text{het}}$  is high enough for the variant to  
352 severely reduce fitness may be more relevant.

353 As more exomes are sequenced, one might expect that we would be better able to more ac-  
354 curately estimate  $s_{\text{het}}$ . Yet, in a companion paper [16], we show that increasing the sample size  
355 used for estimating LOF frequencies will provide essentially no additional information for the  
356  $\sim 85\%$  of genes with the lowest values of  $s_{\text{het}}$ . This fundamental limit on how much we can learn  
357 about these genes from LOF data alone highlights the importance of approaches like ours that  
358 can leverage additional data types. By sharing information across genes, we can overcome this  
359 fundamental limit on how accurately we can estimate constraint.

360 Here we focused on estimating  $s_{\text{het}}$ , but our empirical Bayes framework, GeneBayes, can be  
361 used in any setting where one has a model that ties a gene-level parameter to gene-level observ-  
362 able data (Supplementary Note D). For example, GeneBayes can be used to find trait-associated  
363 genes using variants from case/control studies [71,72], or to improve power to find differen-  
364 tially expressed genes in RNA-seq experiments [73]. We provide a graphical overview of how  
365 GeneBayes can be applied more generally in Figure 6. Briefly, GeneBayes requires users to specify  
366 a likelihood model and the form of a prior distribution for their parameter of interest. Then, using  
367 empirical Bayes and a set of gene features, it improves power to estimate the parameter by flexibly  
368 sharing information across similar genes.

369 In summary, we developed a powerful framework for estimating a broadly applicable and  
370 readily interpretable metric of constraint,  $s_{\text{het}}$ . Our estimates provide a more informative ranking  
371 of gene importance than existing metrics, and our approach allows us to interrogate potential  
372 causes and consequences of natural selection.

### 373 **Data availability**

374 Posterior means and 95% credible intervals for  $s_{\text{het}}$  are available in Supplementary Table 2. Poste-  
375 rior densities for  $s_{\text{het}}$  are available in Supplementary Table 3. A description of the gene features is  
376 available in Supplementary Table 4. These supplementary tables are also available at [74], along  
377 with likelihoods for  $s_{\text{het}}$ , LOF variants with misannotation probabilities, and gene feature tables.

### 378 **Code availability**

379 GeneBayes and code for estimating  $s_{\text{het}}$  are available at <https://github.com/tkzeng/GeneBayes>.

### 380 **Acknowledgements**

381 We would like to thank Ipsita Agarwal, Molly Przeworski, Jesse Engreitz, and members of the  
382 Pritchard Lab for valuable feedback and discussions. This work was supported by NIH grants  
383 R01HG011432 and R01HG008140.

## References

- [1] Cassa CA, Weghorn D, Balick DJ, Jordan DM, Nusinow D, Samocha KE, et al. Estimating the selective effects of heterozygous protein-truncating variants from human exome data. *Nature Genetics*. 2017;49(5):806-10.
- [2] Weghorn D, Balick DJ, Cassa C, Kosmicki JA, Daly MJ, Beier DR, et al. Applicability of the Mutation–Selection Balance Model to Population Genetics of Heterozygous Protein-Truncating Variants in Humans. *Molecular Biology and Evolution*. 2019;36(8):1701-10.
- [3] Fuller ZL, Berg JJ, Mostafavi H, Sella G, Przeworski M. Measuring intolerance to mutation in human genetics. *Nature Genetics*. 2019;51(5):772-6.
- [4] Agarwal I, Fuller ZL, Myers SR, Przeworski M. Relating pathogenic loss-of function mutations in humans to their evolutionary fitness costs. *eLife*. 2023;12:e83172.
- [5] Kaplanis J, Samocha KE, Wiel L, Zhang Z, Arvai KJ, Eberhardt RY, et al. Evidence for 28 genetic disorders discovered by combining healthcare and research data. *Nature*. 2020;586(7831):757-62.
- [6] Fu JM, Satterstrom FK, Peng M, Brand H, Collins RL, Dong S, et al. Rare coding variation provides insight into the genetic architecture and phenotypic context of autism. *Nature Genetics*. 2022;54(9):1320-31.
- [7] Whiffin N, Armean IM, Kleinman A, Marshall JL, Minikel EV, Goodrich JK, et al. The effect of LRRK2 loss-of-function variants in humans. *Nature Medicine*. 2020;26(6):869-77.
- [8] Gazal S, Weissbrod O, Hormozdiari F, Dey KK, Nasser J, Jagadeesh KA, et al. Combining SNP-to-gene linking strategies to identify disease genes and assess disease omnigenicity. *Nature Genetics*. 2022;54(6):827-36.
- [9] Wang X, Goldstein DB. Enhancer domains predict gene pathogenicity and inform gene discovery in complex disease. *The American Journal of Human Genetics*. 2020;106(2):215-33.
- [10] Mostafavi H, Spence JP, Naqvi S, Pritchard JK. Limited overlap of eQTLs and GWAS hits due to systematic differences in discovery. *bioRxiv*. 2022:2022-05.
- [11] Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature*. 2016;536(7616):285-91.
- [12] Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*. 2020;581(7809):434-43.
- [13] Gillespie JH. *Population genetics: a concise guide*. JHU press; 2004.
- [14] Gudmundsson S, Singer-Berk M, Watts NA, Phu W, Goodrich JK, Solomonson M, et al. Variant interpretation using population databases: Lessons from gnomAD. *Human Mutation*. 2022;43(8):1012-30.

- 420 [15] LaPolice TM, Huang YF. A deep learning framework for predicting human essential genes  
421 from population and functional genomic data. *bioRxiv*. 2021:2021-12.
- 422 [16] Spence J, Zeng T, Mostafavi H, Pritchard J. Scaling the discrete-time Wright-Fisher model to  
423 biobank-scale datasets. *bioRxiv*. 2023.
- 424 [17] Duan T, Anand A, Ding DY, Thai KK, Basu S, Ng A, et al. Ngboost: Natural gradient  
425 boosting for probabilistic prediction. In: *International Conference on Machine Learning*.  
426 PMLR; 2020. p. 2690-700.
- 427 [18] Ewens WJ. *Mathematical population genetics: theoretical introduction*. vol. 27. Springer;  
428 2004.
- 429 [19] Agarwal I, Przeworski M. Mutation saturation for fitness effects at human CpG sites. *eLife*.  
430 2021;10:e71513.
- 431 [20] Haghghi K, Kolokathis F, Pater L, Lynch RA, Asahi M, Gramolini AO, et al. Human phos-  
432 pholamban null results in lethal dilated cardiomyopathy revealing a critical difference be-  
433 tween mouse and human. *The Journal of Clinical Investigation*. 2003;111(6):869-76.
- 434 [21] Howard TD, Paznekas WA, Green ED, Chiang LC, Ma N, Luna RIOD, et al. Mutations in  
435 TWIST, a basic helix–loop–helix transcription factor, in Saethre-Chatzen syndrome. *Nature*  
436 *Genetics*. 1997;15(1):36-41.
- 437 [22] Ghouzzi VE, Merrer ML, Perrin-Schmitt F, Lajeunie E, Benit P, Renier D, et al. Mutations of  
438 the TWIST gene in the Saethre-Chatzene syndrome. *Nature Genetics*. 1997;15(1):42-6.
- 439 [23] Da Costa L, Leblanc T, Mohandas N. Diamond-Blackfan anemia. *Blood*. 2020;136(11):1262-  
440 73.
- 441 [24] des Portes V, Pinard JM, Billuart P, Vinet MC, Koulakoff A, Carrié A, et al. A novel CNS gene  
442 required for neuronal migration and involved in X-linked subcortical laminar heterotopia  
443 and lissencephaly syndrome. *Cell*. 1998;92(1):51-61.
- 444 [25] Fantes J, Ragge NK, Lynch SA, McGill NI, Collin JRO, Howard-Peebles PN, et al. Mutations  
445 in SOX2 cause anophthalmia. *Nature Genetics*. 2003;33(4):462-3.
- 446 [26] Berger W, de Pol Dv, Warburg M, Gal A, Bleeker-Wagemakers L, de Silva H, et al. Mutations  
447 in the candidate gene for Norrie disease. *Human Molecular Genetics*. 1992;1(7):461-5.
- 448 [27] Faundes V, Jennings MD, Crilly S, Legraie S, Withers SE, Cuvertino S, et al. Impaired eIF5A  
449 function causes a Mendelian disorder that is partially rescued in model systems by spermi-  
450 dine. *Nature Communications*. 2021;12(1):833.
- 451 [28] Hatada I, Ohashi H, Fukushima Y, Kaneko Y, Inoue M, Komoto Y, et al. An imprinted gene  
452 p57 KIP2 is mutated in Beckwith–Wiedemann syndrome. *Nature Genetics*. 1996;14(2):171-3.
- 453 [29] Gripp KW, Wotton D, Edwards MC, Roessler E, Ades L, Meinecke P, et al. Mutations in TGIF  
454 cause holoprosencephaly and link NODAL signalling to human neural axis determination.  
455 *Nature Genetics*. 2000;25(2):205-8.

- 456 [30] Coffey AJ, Brooksbank RA, Brandau O, Oohashi T, Howell GR, Bye JM, et al. Host response  
457 to EBV infection in X-linked lymphoproliferative disease results from mutations in an SH2-  
458 domain encoding gene. *Nature Genetics*. 1998;20(2):129-35.
- 459 [31] Smith ML, Cavenagh JD, Lister TA, Fitzgibbon J. Mutation of CEBPA in familial acute  
460 myeloid leukemia. *New England Journal of Medicine*. 2004;351(23):2403-7.
- 461 [32] Garg V, Kathiriya IS, Barnes R, Schluterman MK, King IN, Butler CA, et al. GATA4 muta-  
462 tions cause human congenital heart defects and reveal an interaction with TBX5. *Nature*.  
463 2003;424(6947):443-7.
- 464 [33] Langton KP, McKie N, Curtis A, Goodship JA, Bond PM, Barker MD, et al. A novel tis-  
465 sue inhibitor of metalloproteinases-3 mutation reveals a common molecular phenotype in  
466 Sorsby's fundus dystrophy. *Journal of Biological Chemistry*. 2000;275(35):27027-31.
- 467 [34] Fang J, Dagenais SL, Erickson RP, Arlt MF, Glynn MW, Gorski JL, et al. Mutations in  
468 FOXC2 (MFH-1), a forkhead family transcription factor, are responsible for the hered-  
469 itary lymphedema-distichiasis syndrome. *The American Journal of Human Genetics*.  
470 2000;67(6):1382-8.
- 471 [35] Begemann M, Zirn B, Santen G, Wirthgen E, Soellner L, Büttel HM, et al. Paternally inherited  
472 IGF2 mutation and growth restriction. *New England Journal of Medicine*. 2015;373(4):349-  
473 56.
- 474 [36] Amberger JS, Bocchini CA, Schiettecatte F, Scott AF, Hamosh A. OMIM. org: Online  
475 Mendelian Inheritance in Man (OMIM®), an online catalog of human genes and genetic  
476 disorders. *Nucleic Acids Research*. 2015;43(D1):D789-98.
- 477 [37] Meyers RM, Bryan JG, McFarland JM, Weir BA, Sizemore AE, Xu H, et al. Computational  
478 correction of copy number effect improves specificity of CRISPR-Cas9 essentiality screens  
479 in cancer cells. *Nature Genetics*. 2017;49(12):1779-84.
- 480 [38] Ghandi M, Huang FW, Jané-Valbuena J, Kryukov GV, Lo CC, McDonald III ER, et al. Next-  
481 generation characterization of the cancer cell line encyclopedia. *Nature*. 2019;569(7757):503-  
482 8.
- 483 [39] Wright CF, Campbell P, Eberhardt RY, Aitken S, Perrett D, Brent S, et al. Genomic Diagnosis  
484 of Rare Pediatric Disease in the United Kingdom and Ireland. *New England Journal of*  
485 *Medicine*. 2023.
- 486 [40] Köhler S, Gargano M, Matentzoglou N, Carmody LC, Lewis-Smith D, Vasilevsky NA, et al.  
487 The human phenotype ontology in 2021. *Nucleic Acids Research*. 2021;49(D1):D1207-17.
- 488 [41] Leitão E, Schröder C, Parenti I, Dalle C, Rastetter A, Kühnel T, et al. Systematic analysis and  
489 prediction of genes associated with monogenic disorders on human chromosome X. *Nature*  
490 *Communications*. 2022;13(1):6570.
- 491 [42] Agoglia RM, Sun D, Birey F, Yoon SJ, Miura Y, Sabatini K, et al. Primate cell fusion disen-  
492 tangles gene regulatory divergence in neurodevelopment. *Nature*. 2021;592(7854):421-7.

- 493 [43] Consortium G. The GTEx Consortium atlas of genetic regulatory effects across human tis-  
494 sues. *Science*. 2020;369(6509):1318-30.
- 495 [44] Basha O, Argov CM, Artzy R, Zoabi Y, Hekselman I, Alfandari L, et al. Differential net-  
496 work analysis of multiple human tissue interactomes highlights tissue-selective processes  
497 and genetic disorder genes. *Bioinformatics*. 2020;36(9):2821-8.
- 498 [45] Gao S, Yan L, Wang R, Li J, Yong J, Zhou X, et al. Tracing the temporal-spatial transcriptome  
499 landscapes of the human fetal digestive tract using single-cell RNA-sequencing. *Nature Cell  
500 Biology*. 2018;20(6):721-34.
- 501 [46] Charlesworth B, et al. *Evolution in age-structured populations*. vol. 2. Cambridge Univer-  
502 sity Press Cambridge; 1994.
- 503 [47] Barrio-Hernandez I, Schwartzentruber J, Shrivastava A, Del-Toro N, Gonzalez A, Zhang Q,  
504 et al. Network expansion of genetic associations defines a pleiotropy map of human cell  
505 biology. *Nature Genetics*. 2023:1-10.
- 506 [48] Van Dam S, Vosa U, van der Graaf A, Franke L, de Magalhaes JP. Gene co-expression anal-  
507 ysis for functional classification and gene–disease predictions. *Briefings in Bioinformatics*.  
508 2018;19(4):575-92.
- 509 [49] Nasser J, Bergman DT, Fulco CP, Guckelberger P, Doughty BR, Patwardhan TA,  
510 et al. Genome-wide enhancer maps link risk variants to disease genes. *Nature*.  
511 2021;593(7858):238-43.
- 512 [50] Mayr C. Regulation by 3'-untranslated regions. *Annual Review of Genetics*. 2017;51:171-94.
- 513 [51] Leppek K, Das R, Barna M. Functional 5' UTR mRNA structures in eukaryotic translation  
514 regulation and how to find them. *Nature Reviews Molecular Cell Biology*. 2018;19(3):158-74.
- 515 [52] Agrawal AF, Whitlock MC. Inferences about the distribution of dominance drawn from  
516 yeast gene knockout data. *Genetics*. 2011;187(2):553-66.
- 517 [53] Mukai T, Chigusa SI, Mettler L, Crow JF. Mutation rate and dominance of genes affecting  
518 viability in *Drosophila melanogaster*. *Genetics*. 1972;72(2):335-55.
- 519 [54] Sella G, Barton NH. Thinking about the evolution of complex traits in the era of genome-  
520 wide association studies. *Annual Review of Genomics and Human Genetics*. 2019;20:461-93.
- 521 [55] Charlesworth B. Effective population size and patterns of molecular evolution and varia-  
522 tion. *Nature Reviews Genetics*. 2009;10(3):195-205.
- 523 [56] Simons YB, Mostafavi H, Smith CJ, Pritchard JK, Sella G. Simple scaling laws control the  
524 genetic architectures of human complex traits. *bioRxiv*. 2022:2022-10.
- 525 [57] Mathieson I, Terhorst J. Direct detection of natural selection in Bronze Age Britain. *Genome  
526 Research*. 2022;32(11-12):2057-67.
- 527 [58] Emdin CA, Khera AV, Natarajan P, Klarin D, Won HH, Peloso GM, et al. Phenotypic char-  
528 acterization of genetically lowered human lipoprotein (a) levels. *Journal of the American  
529 College of Cardiology*. 2016;68(25):2761-72.

- 530 [59] Langsted A, Nordestgaard BG, Kamstrup PR. Low lipoprotein (a) levels and risk of  
531 disease in a large, contemporary, general population study. *European Heart Journal*.  
532 2021;42(12):1147-56.
- 533 [60] Rausell A, Luo Y, Lopez M, Seeleuthner Y, Rapaport F, Favier A, et al. Common homozy-  
534 gosity for predicted loss-of-function variants reveals both redundant and advantageous  
535 effects of dispensable human genes. *Proceedings of the National Academy of Sciences*.  
536 2020;117(24):13626-36.
- 537 [61] Reyes-Soffer G, Ginsberg HN, Berglund L, Duell PB, Heffron SP, Kamstrup PR, et al.  
538 Lipoprotein (a): a genetically determined, causal, and prevalent risk factor for atheroscle-  
539 rotic cardiovascular disease: a scientific statement from the American Heart Association.  
540 *Arteriosclerosis, Thrombosis, and Vascular Biology*. 2022;42(1):e48-60.
- 541 [62] Millar DS, Johansen B, Berntorp E, Minford A, Bolton-Maggs P, Wensley R, et al. Molecular  
542 genetic analysis of severe protein C deficiency. *Human Genetics*. 2000;106:646-53.
- 543 [63] Romeo G, Hassan HJ, Staempfli S, Roncuzzi L, Cianetti L, Leonardi A, et al. Hereditary  
544 thrombophilia: identification of nonsense and missense mutations in the protein C gene.  
545 *Proceedings of the National Academy of Sciences*. 1987;84(9):2829-32.
- 546 [64] Landrum MJ, Chitipiralla S, Brown GR, Chen C, Gu B, Hart J, et al. ClinVar: improvements  
547 to accessing data. *Nucleic Acids Research*. 2020;48(D1):D835-44.
- 548 [65] Couch FJ, Nathanson KL, Offit K. Two decades after BRCA: setting paradigms in personal-  
549 ized cancer care and prevention. *Science*. 2014;343(6178):1466-70.
- 550 [66] Smith KR, Hanson HA, Hollingshaus MS. BRCA1 and BRCA2 mutations and female fertili-  
551 ty. *Current Opinion in Obstetrics & Gynecology*. 2013;25(3):207.
- 552 [67] O'Connor LJ, Schoech AP, Hormozdiari F, Gazal S, Patterson N, Price AL. Extreme poly-  
553 genicity of complex traits is explained by negative selection. *The American Journal of Hu-  
554 man Genetics*. 2019;105(3):456-76.
- 555 [68] Benton ML, Abraham A, LaBella AL, Abbot P, Rokas A, Capra JA. The influence of evolu-  
556 tionary history on human health and disease. *Nature Reviews Genetics*. 2021;22(5):269-83.
- 557 [69] Satterstrom FK, Kosmicki JA, Wang J, Breen MS, De Rubeis S, An JY, et al. Large-scale exome  
558 sequencing study implicates both developmental and functional changes in the neurobiol-  
559 ogy of autism. *Cell*. 2020;180(3):568-84.
- 560 [70] Gardner EJ, Neville MD, Samocha KE, Barclay K, Kolk M, Niemi ME, et al. Reduced  
561 reproductive success is associated with selective constraint on human genes. *Nature*.  
562 2022;603(7903):858-63.
- 563 [71] He X, Sanders SJ, Liu L, De Rubeis S, Lim ET, Sutcliffe JS, et al. Integrated model of de novo  
564 and inherited genetic variants yields greater power to identify risk genes. *PLoS Genetics*.  
565 2013;9(8):e1003671.
- 566 [72] Zhu X, Stephens M. Bayesian large-scale multiple regression with summary statistics from  
567 genome-wide association studies. *The Annals of Applied Statistics*. 2017;11(3):1561.

- 568 [73] Boyeau P, Regier J, Gayoso A, Jordan MI, Lopez R, Yosef N. An empirical Bayes method  
569 for differential expression analysis of single cells with deep generative models. *bioRxiv*.  
570 2022:2022-05.
- 571 [74] Zeng T, Spence JP, Mostafavi H, Pritchard JK. s\_het estimates from GeneBayes and other  
572 supplementary datasets. Zenodo; 2023. Available from: [https://doi.org/10.5281/  
573 zenodo.7939768](https://doi.org/10.5281/zenodo.7939768).
- 574 [75] Schiffels S, Durbin R. Inferring human population size and separation history from multiple  
575 genome sequences. *Nature Genetics*. 2014;46(8):919-25.
- 576 [76] Cummings BB, Karczewski KJ, Kosmicki JA, Seaby EG, Watts NA, Singer-Berk M, et al.  
577 Transcript expression-aware annotation improves rare variant interpretation. *Nature*.  
578 2020;581(7809):452-8.
- 579 [77] McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The ensembl variant  
580 effect predictor. *Genome Biology*. 2016;17(1):1-14.
- 581 [78] Frankish A, Carbonell-Sala S, Diekhans M, Jungreis I, Loveland JE, Mudge JM, et al. GEN-  
582 CODE: reference annotation for the human and mouse genomes in 2023. *Nucleic Acids  
583 Research*. 2023;51(D1):D942-9.
- 584 [79] Olson ND, Wagner J, McDaniel J, Stephens SH, Westreich ST, Prasanna AG, et al. Preci-  
585 sionFDA Truth Challenge V2: Calling variants from short and long reads in difficult-to-map  
586 regions. *Cell Genomics*. 2022;2(5):100129.
- 587 [80] Blake JA, Baldarelli R, Kadin JA, Richardson JE, Smith CL, Bult CJ. Mouse Genome Database  
588 (MGD): Knowledgebase for mouse–human comparative biology. *Nucleic Acids Research*.  
589 2021;49(D1):D981-7.
- 590 [81] Groza T, Gomez FL, Mashhadi HH, Muñoz-Fuentes V, Gunes O, Wilson R, et al. The Inter-  
591 national Mouse Phenotyping Consortium: comprehensive knockout phenotyping under-  
592 pinning the study of human disease. *Nucleic acids research*. 2023;51(D1):D1038-45.
- 593 [82] Hart T, Brown KR, Sircoulomb F, Rottapel R, Moffat J. Measuring error rates in genomic  
594 perturbation screens: gold standards for human functional genomics. *Molecular Systems  
595 Biology*. 2014;10(7):733.
- 596 [83] Blomen VA, Májek P, Jae LT, Bigenzahn JW, Nieuwenhuis J, Staring J, et al. Gene essentiality  
597 and synthetic lethality in haploid human cells. *Science*. 2015;350(6264):1092-6.
- 598 [84] Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, et al. A frame-  
599 work for the interpretation of de novo mutation in human disease. *Nature Genetics*.  
600 2014;46(9):944-50.
- 601 [85] Amari Si. Natural Gradient Works Efficiently in Learning. *Neural Computation*.  
602 1998;10(2):251-76.
- 603 [86] Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. Pytorch: An imperative  
604 style, high-performance deep learning library. *Advances in Neural Information Processing  
605 Systems*. 2019;32.



- 606 [87] Loshchilov I, Hutter F. Decoupled weight decay regularization. arXiv preprint  
607 arXiv:171105101. 2017.
- 608 [88] Gómez P, Toftevaag HH, Meoni G. torchquad: Numerical Integration in Arbitrary Dimen-  
609 sions with PyTorch. *Journal of Open Source Software*. 2021;6(64):3439.
- 610 [89] Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: *Proceedings of the 22nd*  
611 *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; 2016.  
612 p. 785-94.
- 613 [90] Sawyer SA, Hartl DL. Population genetics of polymorphism and divergence. *Genetics*.  
614 1992;132(4):1161-76.
- 615 [91] Harpak A, Bhaskar A, Pritchard JK. Mutation rate variation is a primary determinant of the  
616 distribution of allele frequencies in humans. *PLoS Genetics*. 2016;12(12):e1006489.
- 617 [92] Varin C, Reid N, Firth D. An overview of composite likelihood methods. *Statistica Sinica*.  
618 2011:5-42.
- 619 [93] Ramoni RB, Mulvihill JJ, Adams DR, Allard P, Ashley EA, Bernstein JA, et al. The undi-  
620 agnosed diseases network: accelerating discovery about health and disease. *The American*  
621 *Journal of Human Genetics*. 2017;100(2):185-92.
- 622 [94] Consortium GP, et al. An integrated map of genetic variation from 1,092 human genomes.  
623 *Nature*. 2012;491(7422):56.
- 624 [95] Lee Y, Nelder JA. Hierarchical generalized linear models. *Journal of the Royal Statistical*  
625 *Society: Series B (Methodological)*. 1996;58(4):619-56.
- 626 [96] Meng XL. Decoding the h-likelihood. *Statistical Science*. 2009;24(3):280-93.
- 627 [97] Weeks EM, Ulirsch JC, Cheng NY, Trippe BL, Fine RS, Miao J, et al. Leveraging poly-  
628 genic enrichments of gene features to predict genes underlying complex traits and diseases.  
629 medRxiv. 2020:2020-09.
- 630 [98] Boukas L, Bjornsson HT, Hansen KD. Promoter CpG density predicts downstream gene  
631 loss-of-function intolerance. *The American Journal of Human Genetics*. 2020;107(3):487-98.
- 632 [99] Pers TH, Karjalainen JM, Chan Y, Westra HJ, Wood AR, Yang J, et al. Biological interpreta-  
633 tion of genome-wide association studies using predicted gene functions. *Nature Communi-*  
634 *cations*. 2015;6(1):5890.
- 635 [100] The Gene Ontology resource: enriching a Gold mine. *Nucleic acids research*.  
636 2021;49(D1):D325-34.
- 637 [101] Raina P, Guinea R, Chatsirisupachai K, Lopes I, Farooq Z, Guinea C, et al. GeneFriends:  
638 gene co-expression databases and tools for humans and model organisms. *Nucleic Acids*  
639 *Research*. 2023;51(D1):D145-58.
- 640 [102] Consortium G, Ardlie KG, Deluca DS, Segrè AV, Sullivan TJ, Young TR, et al. The Genotype-  
641 Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science*.  
642 2015;348(6235):648-60.

- 643 [103] DGT RPC, Consortium F, et al. A promoter-level mammalian expression atlas. *Nature*.  
644 2014;507(7493):462-70.
- 645 [104] Fulco CP, Nasser J, Jones TR, Munson G, Bergman DT, Subramanian V, et al. Activity-by-  
646 contact model of enhancer–promoter regulation from thousands of CRISPR perturbations.  
647 *Nature Genetics*. 2019;51(12):1664-9.
- 648 [105] Roadmap EC, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, et al. Integrative analysis  
649 of 111 reference human epigenomes. *Nature*. 2015;518(7539):317-30.
- 650 [106] Liu Y, Sarkar A, Kheradpour P, Ernst J, Kellis M. Evidence of reduced recombination rate in  
651 human regulatory domains. *Genome Biology*. 2017;18(1):1-11.
- 652 [107] Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, et al. Evolutionarily  
653 conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Research*.  
654 2005;15(8):1034-50.
- 655 [108] Sullivan PF, Meadows JR, Gazal S, Phan BN, Li X, Genereux DP, et al. Leveraging base-  
656 pair mammalian constraint to understand genetic variation and human disease. *Science*.  
657 2023;380(6643):eabn2937.
- 658 [109] Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. Detection of nonneutral substitution rates  
659 on mammalian phylogenies. *Genome research*. 2010;20(1):110-21.
- 660 [110] Elnaggar A, Heinzinger M, Dallago C, Rehawi G, Wang Y, Jones L, et al. Prottrans: Toward  
661 understanding the language of life through self-supervised learning. *IEEE Transactions on*  
662 *Pattern Analysis and Machine Intelligence*. 2021;44(10):7112-27.
- 663 [111] Stärk H, Dallago C, Heinzinger M, Rost B. Light attention predicts protein location from the  
664 language of life. *Bioinformatics Advances*. 2021;1(1):vbab035.
- 665 [112] Huang YF. Unified inference of missense variant effects and gene constraints in the human  
666 genome. *PLoS Genetics*. 2020;16(7):e1008922.

## 667 4 Methods

### 668 Empirical Bayes overview

669 Many genes have few observed loss-of-function variants, making it challenging to infer constraint  
670 without additional information. Bayesian approaches that specify a prior distribution for each  
671 gene can provide such information to improve constraint estimates, but specifying prior distri-  
672 butions is challenging as we have limited prior knowledge about the selection coefficients  $s_{\text{het}}$ .  
673 Empirical Bayes procedures allow us to learn a prior distribution for each gene by combining  
674 information across genes.

675 To use the information contained in the gene features, we learn a mapping from a gene’s fea-  
676 tures to a prior specific for that gene. We parameterize this mapping using gradient-boosted trees,  
677 as implemented in NGBoost [17]. Intuitively, this approach learns a notion of “similarity” between  
678 genes based on their features, and then shares information across similar genes to learn how  $s_{\text{het}}$   
679 relates to the gene features. This approach has two major benefits. First, by sharing information  
680 between similar genes, it can dramatically improve the accuracy of the predicted  $s_{\text{het}}$  values, par-  
681 ticularly for genes with few expected LOFs. Second, by leveraging the LOF data, this approach  
682 allows us to learn about how the various gene features relate to fitness, which cannot be modeled  
683 from first principles.

684 For a more in-depth description of our approach along with mathematical and implementation  
685 details, see Supplementary Note A.

### 686 Population genetic likelihood

687 To model how  $s_{\text{het}}$  relates to the frequency of individual LOF variants, we used the discrete-time  
688 Wright-Fisher model, with an approximation of diploid selection with additive fitness effects. We  
689 used a composite likelihood approach, assuming independence across individual LOF variants to  
690 obtain gene-level likelihoods. Within this composite likelihood, we model each individual variant  
691 as either having a selection coefficient of  $s_{\text{het}}$  with probability  $1 - p_{\text{miss}}$ , or having a selection  
692 coefficient of 0 with probability  $p_{\text{miss}}$ . That is,  $p_{\text{miss}}$  acts as the prior probability that a given variant  
693 is misannotated, and we assume that misannotated variants evolve neutrally regardless of the  
694 strength of selection on the gene. All likelihoods were computed using new machinery developed  
695 in a companion paper [16].

696 Our model depends on a number of parameters—a demographic model of past population  
697 sizes, mutation rates for each site, and the probability of misannotation. The demographic model  
698 is taken from the literature [75] with modifications as described in [4]. The mutation rates account  
699 for trinucleotide context as well as methylation status at CpGs [12]. Finally, we estimated the  
700 probability of misannotation from the data.

701 For additional technical details and intuition see Supplementary Note B.

702 **Curation of LOF variants**

703 We obtained annotations for the consequences of all possible single nucleotide changes to the  
 704 hg19 reference genome from [76]. The effects of variants on protein function were predicted us-  
 705 ing Variant Effect Predictor (VEP) version 85 [77] using GENCODE v19 gene annotations [78] as  
 706 a reference. We defined a variant as a LOF if it was predicted by VEP to be a splice acceptor,  
 707 splice donor, or stop gain variant. In addition, predicted LOFs were further annotated using LOF-  
 708 TEE [12], which implements a series of filters to identify variants that may be misannotated (for  
 709 example, LOFTTE considers predicted LOFs near the ends of transcripts as likely misannotations).  
 710 For our analyses, we only kept predicted LOFs labelled as High Confidence by LOFTTE, which  
 711 are LOFs that passed all of LOFTTE’s filters.

712 Next, we considered potential criteria for further filtering LOFs: cutoffs for the median ex-  
 713 ome sequencing read depth, cutoffs for the mean pext (proportion expressed across transcripts)  
 714 score [76], whether to exclude variants that fall in segmental duplications or regions with low  
 715 mappability [79], and whether to exclude variants flagged by LOFTTE as potentially problematic  
 716 but that passed LOFTTE’s primary filters.

717 We trained models with these filters one at a time and in combination, and chose the model  
 718 that had the best AUPRC in classifying essential from nonessential genes in mice. The filters  
 719 we evaluated and chose for the final model are reported in Table 2. Since we used mouse gene  
 720 essentiality data to choose the filters, we do not further evaluate  $s_{\text{het}}$  on these data.

721 We considered genes to be essential in mice if they are heterozygous lethal, as determined  
 722 by [12] using data from heterozygous knockouts reported in Mouse Genome Informatics [80].  
 723 We classify genes as nonessential if they are reported as “Viable with No Phenotype” by the  
 724 International Mouse Phenotyping Consortium [81] (annotations downloaded on 12/08/22 from  
 725 <https://www.ebi.ac.uk/mi/impc/essential-genes-search/>).

Filtering criterion	Tested values	Best value
Cutoff for sequencing read depth (median across exomes)	5×, 10×, 20×	20×
Cutoff for mean pext across tissues	0.05, 0.1	0.05
Filter if variant falls in a segmental duplication or low mappability region	True, False	False
Filter if variant is flagged as potentially problematic	True, False	True

Table 2: Filtering criteria for LOF curation

726 Finally, we annotated each variant with its frequency in the gnomAD v2.1.1 exomes [12], a  
 727 dataset of 125,748 uniformly-analyzed exomes that were largely curated from case–control stud-  
 728 ies of common adult-onset diseases. gnomAD provides precomputed allele frequencies for all  
 729 variants that they call.

730 For potential LOFs that are not segregating, gnomAD does not release the number of indi-  
 731 viduals that were genotyped at those positions. For these sites, we used the median number of  
 732 genotyped individuals at the positions for which gnomAD does provide this information. We  
 733 performed this separately on the autosomes and X chromosome.

734 Data sources for the variant annotations, filters, and frequencies, as well as additional infor-  
 735 mation used to compute likelihoods are listed in Table 3.

Resource	Link
Annotations for possible LOFs	<a href="gs://gnomad-public/papers/2019-tx-annotation/pre_computed/all.possible.snvs.tx_annotated.GTEx.v7.021520.tsv">gs://gnomad-public/papers/2019-tx-annotation/pre_computed/all.possible.snvs.tx_annotated.GTEx.v7.021520.tsv</a>
Mean methylation for CpG sites	<a href="gs://gcp-public-data--gnomad/resources/methylation">gs://gcp-public-data--gnomad/resources/methylation</a>
Exome sequencing coverage	<a href="gs://gcp-public-data--gnomad/release/2.1/coverage/exomes/gnomad.exomes.coverage.summary.tsv.bgz">gs://gcp-public-data--gnomad/release/2.1/coverage/exomes/gnomad.exomes.coverage.summary.tsv.bgz</a>
Variant frequencies	<a href="gs://gcp-public-data--gnomad/release/2.1.1/vcf/exomes/gnomad.exomes.r2.1.1.sites.vcf.bgz">gs://gcp-public-data--gnomad/release/2.1.1/vcf/exomes/gnomad.exomes.r2.1.1.sites.vcf.bgz</a>
Low mappability and segmental duplications	<a href="https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/genome-stratifications/v3.1/GRCh37/Union/GRCh37_allowmapandsegdupregions.bed.gz">https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/genome-stratifications/v3.1/GRCh37/Union/GRCh37_allowmapandsegdupregions.bed.gz</a>

Table 3: Sources for LOF data

### 736 Feature processing and selection

737 We compiled 10 types of gene features from several sources:

- 738 1. Gene structure (e.g., number of transcripts, number of exons, GC content)
- 739 2. Gene expression across tissues and cell lines
- 740 3. Biological pathways and Gene Ontology terms
- 741 4. Protein-protein interaction networks
- 742 5. Co-expression networks
- 743 6. Gene regulatory landscape (e.g., number and properties of enhancers and promoters)
- 744 7. Conservation across species
- 745 8. Protein embeddings
- 746 9. Subcellular localization
- 747 10. Missense constraint

748 Additionally, we included an indicator variable that is 1 if the gene is on the non-pseudoautosomal  
749 region of the X chromosome and 0 otherwise.

750 For a description of the features within each category and where we acquired them, see Sup-  
751 plementary Note C.

### 752 Training and validation

753 We fine-tuned a set of hyperparameters for our full empirical Bayes approach, using the best hy-  
754 perparameters from an initial feature selection step (described in Supplementary Note C) as a  
755 starting point. To minimize overfitting, we split the genes into three sets—a training set (chromo-  
756 somes 7-22, X), a validation set for hyperparameter tuning (chromosomes 2, 4, 6), and a test set

757 to evaluate overfitting (chromosomes 1, 3, 5). During each training iteration, one or more trees  
 758 were added to the model to fit the natural gradient of the loss on the training set. We stopped  
 759 model training once the loss on the validation set did not improve for 10 iterations in a row (or the  
 760 maximum number of iterations, 1,000, was reached). Using this approach, we performed a grid  
 761 search over the hyperparameters listed in Table 4 and used the combination that minimized the  
 762 validation loss.

Parameter(s)	Tested values	Best value
Learning rate	0.0125, 0.05, 0.2	0.0125
Maximum tree depth ( <code>max_depth</code> )	3, 4, 5	3
Data subsampling ratio ( <code>subsample</code> )	0.6, 0.8, 1	0.8
Minimum weight of a leaf node ( <code>min_child_weight</code> )	1, 2, 4	1
L1 regularization ( <code>alpha</code> )	0, 1, 2	2
L2 regularization ( <code>lambda</code> )	1, 2, 4	1
Number of trees to fit per iteration ( <code>n_estimators</code> )	1, 2, 4	4

Table 4: Parameters for fitting the gradient-boosted trees

763 For Figure 2B, we reported results from the best model learned using the training set. For all  
 764 other results, we trained a model on all genes using the hyperparameters and number of training  
 765 iterations learned during this hyperparameter fine-tuning step.

## 766 Choosing genes for Table 1

767 To identify genes that are considered constrained by  $s_{\text{het}}$  but not by LOEUF, we filtered for genes  
 768 with  $s_{\text{het}} > 0.1$  (top  $\sim 17\%$  most constrained genes, analogous to the recommended LOEUF cutoff  
 769 of 0.35 [14], which corresponds to the top  $\sim 16\%$  of genes) and  $\text{LOEUF} > 0.5$  (least constrained  
 770  $\sim 73\%$  of genes). Of these, we identified genes where heterozygous or hemizygous mutations that  
 771 decrease the amount of functional protein (e.g. LOF mutations) are associated with Mendelian  
 772 disorders in the Online Mendelian Inheritance in Man (OMIM) database [36]. We chose genes for  
 773 Table 1 primarily based on their prominence in the existing literature.

## 774 Evaluation on additional datasets

### 775 Definition of human essential and nonessential genes

776 We obtained data from 1,085 CRISPR knockout screens quantifying the effects of genes on cell  
 777 survival or proliferation from the DepMap portal (22Q2 release) [37, 38]. Scores from each screen  
 778 are normalized such that nonessential genes identified by [82] have a median score of 0 and that  
 779 common essential genes identified by [82, 83] have a median score of  $-1$ .

780 In classifying essential genes (Figure 3A), we define a gene as essential if its score is  $< -1$   
 781 in at least 25% of screens, and as *not* essential if its score is  $> -1$  in all screens. In classifying  
 782 nonessential genes, we define a gene as nonessential if it has a minimal effect on growth in most  
 783 cell lines (score  $> -0.25$  and  $< 0.25$  in at least 99% of screens), and as *not* nonessential if its score  
 784 is  $< 0$  in all screens.

## 785 **Definition of developmental disorder genes**

786 Through the Deciphering Developmental Disorders (DDD) study [39], clinicians have annotated  
787 a subset of genes with the strength and nature of their association with developmental disorder  
788 s. We classify genes as developmental disorder genes if they are annotated by the DDD study  
789 with `confidence_category = definitive` and `allelic_requirement = monoallelic_autosomal,`  
790 `monoallelic_X_hem` (hemizygous), or `monoallelic_X_het` (heterozygous).

791 We classify genes as not associated with developmental disorders if they are annotated by the  
792 DDD study, do not meet the above criteria, and are not annotated with `confidence_category =`  
793 `strong` or `moderate` and `allelic_requirement = monoallelic_autosomal, monoallelic_X_hem,`  
794 or `monoallelic_X_het`.

795 We downloaded genes with DDD annotations from [https://www.deciphergenomics.org/ddd/](https://www.deciphergenomics.org/ddd/ddgenes)  
796 `ddgenes` on 05/06/2023.

## 797 **Enrichment/depletion of Human Phenotype Ontology (HPO) genes**

798 The Human Phenotype Ontology (HPO) provides a structured organization of phenotypic abnormal-  
799 ities and the genes associated with them, with each HPO term corresponding to a phenotypic  
800 abnormality. We calculated the enrichment of constrained genes in each HPO term with at least  
801 200 genes as the ratio (fraction of HPO genes under constraint)/(fraction of background genes  
802 under constraint). We defined genes under constraint to be the decile of genes considered most  
803 constrained by  $s_{\text{het}}$  or LOEUF. To choose background genes, we sampled from the set of all genes  
804 to match each HPO term's distribution of expected unique LOFs. Similarly, we calculated the de-  
805 pletion of unconstrained genes in each HPO term as the ratio (fraction of HPO genes not under  
806 constraint)/(fraction of background genes not under constraint), where we define genes not under  
807 constraint to be the decile of genes considered least constrained by  $s_{\text{het}}$  or LOEUF.

808 We downloaded HPO phenotype-to-gene annotations from [http://purl.obolibrary.org/](http://purl.obolibrary.org/obo/hp/hpoa/phenotype_to_genes.txt)  
809 `obo/hp/hpoa/phenotype_to_genes.txt` on 01/27/2023.

## 810 **Enrichment of *de novo* mutations in developmental disorder patients**

811 We used the enrichment metric developed by [5] in their analysis of *de novo* mutations (DNMs)  
812 identified from exome sequencing of 31,058 developmental disorder patients and their unaffected  
813 parents. Enrichment of DNMs in developmental disorder patients was calculated as the ratio  
814 of observed DNMs in patients over the expected number under a null mutational model that  
815 accounts for the study sample size and triplet mutation rate at the mutation sites [84].

816 For Figure 3D, we calculated the enrichment of DNMs in constrained genes, defined as the  
817 decile of genes considered most constrained by  $s_{\text{het}}$  or LOEUF. For Supplementary Figure 2C, we  
818 calculated the enrichment of DNMs in constrained genes with and without known associations  
819 with development disorders. We defined a gene as having a known association if it is anno-  
820 tated by the DDD study (see Methods section "Definition of developmental disorder genes") with  
821 `confidence_category = definitive` or `strong` and `allelic_requirement = monoallelic_autosomal,`  
822 `monoallelic_X_hem` (hemizygous), or `monoallelic_X_het` (heterozygous).

823 For each set of genes, we computed the mean enrichment over sites and 95% Poisson confi-  
824 dence intervals for the mean using the code provided by [5].

### 825 **Expression variability across species**

826 To understand the variability in expression between humans and other species, we focused on  
827 gene expression differences between human and chimpanzee as estimated from RNA sequencing  
828 of an *in vitro* model of the developing cerebral cortex for each species [42]. As a metric of vari-  
829 ability between the two species, we used the absolute log-fold change (LFC) in gene expression  
830 between human and chimpanzee cortical spheroids, which was calculated from samples collected  
831 at several time points throughout differentiation of the spheroids. LFC estimates were obtained  
832 from Supplementary Table 9 of [42].

833 To visualize the relationship between constraint and absolute LFC, we plotted a LOESS curve  
834 between the constraint on a gene (gene rank from least to most constrained using either  $s_{\text{het}}$  or  
835 LOEUF as the constraint metric) and the absolute LFC for the gene. Curves were calculated using  
836 the LOWESS function from the `statsmodels` package with parameters `frac = 0.15` and `delta = 10`.

### 837 **Expression variability across individuals**

838 We used the coefficient of variance (CV) as a metric for gene expression variability across individ-  
839 uals, defined as  $CV = \sigma_i / \mu_i$  where  $\sigma_i$  and  $\mu_i$  are the standard deviation and mean of the expression  
840 level of gene  $i$  respectively. Here, expression is in units of Transcripts Per Million. We calculated  
841 CV using 17,398 RNA-seq samples in the GTEx v8 release [43], with data from 838 donors and 52  
842 tissues/cell lines.

843 Another potential metric for gene expression variability is the standard deviation for a gene,  $\sigma_i$ .  
844 However, as the mean expression for a gene,  $\mu_i$ , is strongly correlated with  $\sigma_i$  (Spearman  $\rho = 0.73$   
845 in GTEx), the relation between  $\sigma_i$  and  $s_{\text{het}}^{(i)}$  may be confounded by the relation between  $\mu_i$  and  $s_{\text{het}}^{(i)}$ .  
846 In contrast, we found that CV is only slightly correlated with  $\mu_i$  (Spearman  $\rho = -0.06$  in GTEx).

847 LOESS curves were computed as in “Expression variability across species.”

### 848 **Feature interpretation**

#### 849 **Training models on feature subsets**

850 We grouped features into categories (see Supplementary Table 4 for the features in each category),  
851 and trained a model for each category to predict  $s_{\text{het}}$  from the corresponding features. For each  
852 model, we tuned hyperparameters over a subset of the values we considered for the full model  
853 (Table 5), and chose the combination of hyperparameters that minimized the loss over genes in  
854 the validation set. As a baseline, we trained a model with no features, such that all genes have a  
855 shared prior distribution that is learned from the LOF data—this model is analogous to a standard  
856 empirical Bayes model.



Parameter(s)	Tested values
Learning rate	0.0125, 0.05
Maximum tree depth ( <code>max_depth</code> )	3
Data subsampling ratio ( <code>subsample</code> )	0.8, 1
Minimum weight of a leaf node ( <code>min_child_weight</code> )	1
L1 regularization ( <code>alpha</code> )	0, 1, 2
L2 regularization ( <code>lambda</code> )	1
Number of trees to fit per iteration ( <code>n_estimators</code> )	1, 2, 4

Table 5: Parameters for feature subsets

## 857 Definition of expression feature subsets

858 We grouped gene expression features into 24 categories representing tissues, cell types, and de-  
859 velopmental stage using terms present in the feature names (Table 6).

Category	Terms in the feature (not case sensitive)
Brain	brain, nerve, microglia, hippocampus
Digestive	digestive, gut, gutendoderm, intestine, colon, ileum
Development	development, gastrulation, embryo
Lung	lung, airway
Eye	eye, retina
Endothelium	endothelium
Muscle	muscle
Hair follicle	hairfollicle
Kidney	kidney
Immune	immune, monocytes, nk, tcell, pbmc
Prostate	prostate
Blood	blood, heme, fetalblood
Adipocyte	adipocyte
Heart	heart, aorta
Thymus	thymus
Pancreas	pancreas, islets, pancreasductal
Liver	liver
Testis	testis
Synovial fibroblast	synovialfibroblast
Bladder	bladder
Placenta	placenta
Bone marrow	bonemarrow
CSF	csf
Lymph nodes	lymphnodes

Table 6: Terms used to define tissues for expression features

## 860 Scoring individual features

861 To score individual gene features, we varied the value of one feature at a time and calculated  
862 the variance in predicted  $s_{\text{het}}$  as a feature score. In more detail, we fixed each feature to val-  
863 ues spanning the range of observed values for that feature (0th, 2nd, ..., 98th, and 100th per-  
864 centile), such that all genes shared the same feature value. Then, for each of these 51 feature  
865 values, we averaged the  $s_{\text{het}}$  values predicted by the learned priors over all genes, where the

866 predicted  $s_{\text{het}}$  for each gene is the mean of its prior. We denote this averaged prediction by  
867  $s_{\text{het}}^{(f)}\{p\}$  for some feature  $f$  and percentile  $p$ . Finally, we define the score for feature  $f$  as  $\text{score}_f =$   
868  $\text{sd}(s_{\text{het}}^{(f)}\{0\}, s_{\text{het}}^{(f)}\{2\}, \dots, s_{\text{het}}^{(f)}\{98\}, s_{\text{het}}^{(f)}\{100\})$ , where  $\text{sd}$  is a function computing the sample standard  
869 deviation. In other words, a feature with a high score is one for which varying its value causes  
870 high variance in the predicted  $s_{\text{het}}$ .

871 For the lineplots in Figures 4C-4F, we scale the predictions  $s_{\text{het}}^{(f)}\{p\}$  for each feature  $f$  by sub-  
872 tracting  $(s_{\text{het}}^{(f)}\{0\} + s_{\text{het}}^{(f)}\{100\})/2$  from each prediction.

### 873 **Pruning features before computing feature scores**

874 While investigating the effects of features on predicted  $s_{\text{het}}$ , we found that including highly corre-  
875 lated features in the model could produce unintuitive results, such as opposite correlations with  
876  $s_{\text{het}}$  for highly similar features. Therefore, for Figures 4C-4F, we first pruned the set of features  
877 to minimize pairwise correlations between the remaining features. To do this, we randomly kept  
878 one feature in each group of correlated features, where such a group is defined as a set of features  
879 where each feature in the set has an absolute Spearman  $\rho > 0.7$  to some other feature in the set.

880 For Figures 4C-4F, we trained models on the relevant features in this pruned set (gene ontol-  
881 ogy, network, gene regulatory, and gene structure features for Figures 4C, 4D, 4E, and 4F respec-  
882 tively). After feature pruning, we found the directions of effect for the features were consistent  
883 with their marginal directions of effect.

## 884 Supplementary Material

### 885 A Empirical Bayes with NGBoost

#### 886 Empirical Bayes overview

In the simplest version of empirical Bayes, we specify the form of the prior distribution and assume that that prior is shared across all genes—for example, for gene  $i$  we might assume the prior distribution is  $s_{\text{het}}^{(i)} \sim \text{LogitNormal}(\mu, \sigma)$  with density  $p_{\mu, \sigma}(s_{\text{het}}^{(i)})$ , where the  $\text{LogitNormal}(\mu, \sigma)$  distribution is defined such that  $\text{logit}(s_{\text{het}}^{(i)}) = \log(s_{\text{het}}^{(i)} / (1 - s_{\text{het}}^{(i)}))$  is normally distributed with mean  $\mu$  and variance  $\sigma^2$ . We can then estimate  $\mu$  and  $\sigma$  using the observed LOF data for each gene,  $\mathbf{y}_1, \dots, \mathbf{y}_M$ , by maximizing the marginal likelihood:

$$\prod_{i=1}^M \int_0^1 p(\mathbf{y}_i | s_{\text{het}}^{(i)}) p_{\mu, \sigma}(s_{\text{het}}^{(i)}) ds_{\text{het}}^{(i)}. \quad (1)$$

Next, we can compute the posterior distribution of  $s_{\text{het}}^{(i)}$  for each gene,

$$p(s_{\text{het}}^{(i)} | \mathbf{y}_i) = \frac{p(\mathbf{y}_i | s_{\text{het}}^{(i)}) p_{\mu, \sigma}(s_{\text{het}}^{(i)})}{\int_0^1 p(\mathbf{y}_i | s_{\text{het}}^{(i)}) p_{\mu, \sigma}(s_{\text{het}}^{(i)}) ds_{\text{het}}^{(i)}}. \quad (2)$$

887 However, rather than learning the parameters for the prior from only the LOF data, we can also  
888 use gene features to learn gene-specific prior parameters,  $\mu_i$  and  $\sigma_i$ . To do this, we used a machine  
889 learning approach, NGBoost, to learn functions  $f$  and  $g$  such that  $\mu_i = f(\mathbf{x}_i)$  and  $\sigma_i = g(\mathbf{x}_i)$ , where  
890  $\mathbf{x}_i$  is a vector of gene features associated with gene  $i$ . In the next few sections, we will describe  
891 how we learned  $f$  and  $g$ .

#### 892 NGBoost

893 NGBoost (Natural Gradient Boosting) is an approach for training gradient boosted trees to predict  
894 the parameters of a probability distribution [17]. Gradient boosted trees are a type of machine  
895 learning model typically used to predict outcomes  $y$ , from features  $X$ , producing point estimates  
896 such as predictions of  $\mathbb{E}[y | X]$ ; in contrast, NGBoost uses gradient boosted trees to predict  $p(y |$   
897  $X = \mathbf{x})$  by learning parameters of  $p(y | X = \mathbf{x})$  as functions of  $\mathbf{x}$ —in other words, NGBoost allows  
898 us to learn the full distribution of  $y$  conditioned on observing the features  $\mathbf{x}$ .

Specifically, for gene  $i$ , we assume the prior distribution is  $s_{\text{het}}^{(i)} \sim \text{LogitNormal}(\mu_i, \sigma_i)$ , with density  $p_{\mu_i, \sigma_i}(s_{\text{het}}^{(i)})$ .  $\mu_i = f(\mathbf{x}_i)$  and  $\sigma_i = g(\mathbf{x}_i)$  are functions of the vector of gene features  $\mathbf{x}_i$ , where  $f$  and  $g$  are parameterized as gradient-boosted trees. We chose this distribution as previous work has suggested that  $s_{\text{het}}^{(i)}$  is distributed on a logarithmic scale [1, 2, 4], yet,  $s_{\text{het}}^{(i)}$  is also bounded between 0 and 1. Both of these properties are enforced by the LogitNormal distribution. In Supplementary Note B, we develop a population genetic likelihood  $p(\mathbf{y}_i | s_{\text{het}}^{(i)})$ , where  $\mathbf{y}_i$  is a vector that represents the observed frequencies of each possible loss of function variant for the gene.

Then, with  $M$  genes in the training set, the score that NGBoost maximizes during training is:

$$\sum_{i=1}^M S(\mathbf{y}_i; \mu_i, \sigma_i) = \sum_{i=1}^M \log p(\mathbf{y}_i) = \sum_{i=1}^M \log \left( \int_0^1 p(\mathbf{y}_i | s_{\text{het}}^{(i)}) p_{\mu_i, \sigma_i}(s_{\text{het}}^{(i)}) ds_{\text{het}}^{(i)} \right). \quad (3)$$

To do this, NGBoost first initializes the parameters of  $f$  and  $g$  such that all genes have the same prior distribution. Next, NGBoost adopts a gradient descent approach to maximize the score function: for each iteration until training ends, NGBoost first computes the natural gradient of gene  $i$ 's score with respect to the parameters  $\mu_i$  and  $\sigma_i$  of  $p_{\mu_i, \sigma_i}(s_{\text{het}}^{(i)})$ , where the natural gradient of  $S = S(\mathbf{y}_i; \mu_i, \sigma_i)$ , is defined as:

$$\tilde{\nabla} S \propto \mathcal{I}_{\mu_i, \sigma_i}^{-1} \nabla_{\mu_i, \sigma_i} S \quad (4)$$

where

$$\mathcal{I}_{\mu_i, \sigma_i} = \mathbb{E}_{s_{\text{het}}^{(i)} \sim p_{\mu_i, \sigma_i}} \left[ \left( \nabla_{\mu_i, \sigma_i} \log p_{\mu_i, \sigma_i}(s_{\text{het}}^{(i)}) \right) \left( \nabla_{\mu_i, \sigma_i} \log p_{\mu_i, \sigma_i}(s_{\text{het}}^{(i)}) \right)^T \right] \quad (5)$$

999 is the Fisher Information Matrix for  $p_{\mu_i, \sigma_i}(s_{\text{het}}^{(i)})$  and  $\nabla_{\mu_i, \sigma_i}$  represents differentiation with respect to  
 900  $\mu_i$  and  $\sigma_i$ . Natural gradients take into account the underlying ‘‘information geometry’’ of the space  
 901 of distributions in a way that standard gradients do not [85]. As an example, changing the variance  
 902 of a Normal distribution from 0.1 to 0.2 is much more dramatic than changing the variance from  
 903 10.1 to 10.2. After computing the natural gradient, NGBoost fits a decision tree to each dimension  
 904 of the natural gradient, updating  $\mu_i$  and  $\sigma_i$  in the direction that most steeply increases the gene’s  
 905 score. While gradient-boosting algorithms (including NGBoost, by default) typically fit a single  
 906 decision tree at each iteration, we allow NGBoost to fit one or more trees, which performs slightly  
 907 better in practice (see ‘‘Training and Validation’’ in Methods).

908 Below, we summarize the training algorithm. Let  $\mu_i^{(t)}, \sigma_i^{(t)}$  denote the parameters of the prior at  
 909 training iteration  $t$ .

910 1. Initialize parameters for all genes,  $i = 1, \dots, M$ :

911  $\mu_i^{(0)}, \sigma_i^{(0)} = \operatorname{argmax}_{\mu, \sigma} \sum_{i=1}^M S(\mathbf{y}_i; \mu, \sigma)$

912 2. For iterations  $t = 1, \dots, T$ :

913 (a) For each gene, calculate natural gradients of the score:

914  $\tilde{\nabla} S(\mathbf{y}_i; \mu_i^{(t)}, \sigma_i^{(t)})$ , whose two components we denote as  $\tilde{\nabla} S_{\mu}$  and  $\tilde{\nabla} S_{\sigma}$

915 (b) Fit decision trees  $f^{(t)}$  and  $g^{(t)}$  on the natural gradients:

916  $f^{(t)} = \operatorname{fit} \left( \left\{ \mathbf{x}_i, \tilde{\nabla} S_{\mu_i} \right\}_{i=1}^M \right)$

917  $g^{(t)} = \operatorname{fit} \left( \left\{ \mathbf{x}_i, \tilde{\nabla} S_{\sigma_i} \right\}_{i=1}^M \right)$

918 (c) Update the parameters for each gene, where  $\eta$  is a learning rate that is chosen by the  
 919 user as a hyperparameter

920  $\mu_i^{(t)} = \mu_i^{(t-1)} - \eta f^{(t)}(\mathbf{x}_i)$

921  $\sigma_i^{(t)} = \sigma_i^{(t-1)} - \eta g^{(t)}(\mathbf{x}_i)$

Once training is complete, we obtain a learned prior with parameters  $\mu_i^{(T)}, \sigma_i^{(T)}$ , and can compute the posterior distribution of  $s_{\text{het}}$

$$p\left(s_{\text{het}}^{(i)} \mid \mathbf{y}_i\right) = \frac{p\left(\mathbf{y}_i \mid s_{\text{het}}^{(i)}\right) p_{\mu_i^{(T)}, \sigma_i^{(T)}}\left(s_{\text{het}}^{(i)}\right)}{p\left(\mathbf{y}_i\right)} \quad (6)$$

as well as the mean of this distribution

$$\mathbb{E}\left[s_{\text{het}}^{(i)} \mid \mathbf{y}_i\right] = \int_0^1 s_{\text{het}}^{(i)} p\left(s_{\text{het}}^{(i)} \mid \mathbf{y}_i\right) ds_{\text{het}}^{(i)} \quad (7)$$

922 To compute 95% Credible Intervals, we compute the CDF of the posterior distribution using  
 923 Pytorch’s `cumulative_trapezoid` function [86]. Then, the 95% Credible Interval per gene is de-  
 924 fined as  $[\text{lb}^{(i)}, \text{ub}^{(i)}]$  such that  $P(s_{\text{het}}^{(i)} < \text{lb}^{(i)}) = 0.025$  and  $P(s_{\text{het}}^{(i)} < \text{ub}^{(i)}) = 0.975$ .

## 925 NGBBoost— implementation details

926 To initialize parameters (step 1 in the training algorithm), we perform gradient descent with the  
 927 AdamW optimizer [87] implemented in PyTorch [86] with a learning rate of  $5 \times 10^{-4}$  and other-  
 928 wise default settings. We initialize the optimization at  $\mu = -5$  and  $\sigma = 0.5$ .

929 To compute the integrals in the score calculation, we use the `torchquad` package for numerical  
 930 integration [88], which allows us to use PyTorch’s automatic differentiation system to compute  
 931 gradients. We perform integration using Boole’s rule, integrating from  $5 \times 10^{-8}$  to  $1 - 5 \times 10^{-8}$   
 932 with  $10^6$  sample points.

933 The Fisher Information Matrix is approximated using a Monte Carlo approach: we sample  $s_{\text{het}}$   
 934 from the prior 1,000 times, compute the gradient for each sample, and approximate the expectation  
 935 using the sample mean.

936 To flexibly fit decision trees at each training iteration, we use the XGBoost package, a library  
 937 used for fitting standard gradient boosted trees [89]. In comparison to the default NGBBoost learner,  
 938 XGBoost supports missing features and allows for adjustment of numerous hyperparameters (see  
 939 “Training and Validation” in Methods). In contrast to typical applications of XGBoost, we only  
 940 allow a few (1-4) trees to be fit at each training iteration, as we are using XGBoost within a training  
 941 loop rather than as a standalone approach for model fitting.

942 All distributions were implemented using PyTorch, and training was conducted with GPU  
 943 support when available, with `tree_method = "gpu_hist"` for the XGBoost learners.

## 944 B Population Genetics Model

### 945 Overview of model

946 Some of the most commonly used measures of gene constraint (pLI [11], LOEUF [12]) are framed  
947 in terms of the number of unique LOFs observed in gene,  $O$ , relative to the number expected  
948 under a null model,  $E$ . While operationalizing constraint as some function of  $O$  and  $E$  captures the  
949 intuition that seeing fewer LOFs than expected is evidence that a gene is conserved, the numerical  
950 values of pLI and LOEUF are difficult to interpret. In practice this means that such measures  
951 can be useful for ranking which genes are important, but it makes it difficult to contextualize  
952 these results in terms of other types of variants, such as missense or noncoding variants, or copy  
953 number variants. Previous approaches have pioneered using a population genetics model in this  
954 context to obtain interpretable estimates, albeit with different technical details that we discuss  
955 below [1,2,4].

In order to obtain a more interpretable measure of constraint, we formalize constraint as the strength of natural selection acting against gene loss-of-function in a population genetics model. That is, we can ask how much fitness is reduced on average for an individual with one or two non-functional copies of a gene relative to individuals with two functional copies, following previous work [1,2,4]. To tie this concept of constraint to observed allele frequency data, we use a slightly simplified version of the discrete-time Wright Fisher model. This model contains mutation, selection, and genetic drift, and assumes that there are only two alleles and that the population is panmictic, monoecious, and has non-overlapping generations. While all of these assumptions are violated in humans (there are four nucleotides, population structure, two sexes, and overlapping generations), the model still provides a good approximation to allele frequency dynamics through time. If the allele frequency in generation  $k$  is  $f_k$ , then we model the allele frequency in the next generation via binomial sampling:

$$2N_{k+1}f_{k+1} \sim \text{Binomial}(2N_{k+1}, p(f_k)), \quad (8)$$

956 where  $N_{k+1}$  is the number of diploid individuals in generation  $k + 1$ , with

$$p(f_k) := \frac{(1 - s_{\text{het}})\tilde{f}_k(1 - \tilde{f}_k) + (1 - s_{\text{hom}})\tilde{f}_k^2}{(1 - \tilde{f}_k)^2 + 2(1 - s_{\text{het}})\tilde{f}_k(1 - \tilde{f}_k) + (1 - s_{\text{hom}})\tilde{f}_k^2},$$

957 where  $\tilde{f}_k = f_k(1 - \mu_{1 \rightarrow 0}) + \mu_{0 \rightarrow 1}(1 - f_k)$  is the allele frequency after alleles change from non-  
958 LOF to LOF at rate  $\mu_{0 \rightarrow 1}$  and from LOF to non-LOF at rate  $\mu_{1 \rightarrow 0}$ . The function  $p(\cdot)$  arises from  
959 considering bidirectional mutation and approximating a model of diploid selection where the  
960 relative reproductive success of individuals with 0, 1, or 2 copies of the LOF are 1,  $1 - s_{\text{het}}$ , and  $1 -$   
961  $s_{\text{hom}}$  respectively [13]. In practice, most LOF variants are extremely rare, and so it is exceedingly  
962 unlikely to find individuals homozygous for the LOF. This makes estimating  $s_{\text{hom}}$  as a separate  
963 parameter very difficult, and so we instead assume that  $s_{\text{hom}} = \min\{2s_{\text{het}}, 1\}$ . This is equivalent  
964 to assuming genic selection (i.e., additive fitness effects) with the constraint that an individual's  
965 relative fitness cannot be lower than 0.

966 Equation 8 fully specifies the model except for an initial condition. That is, we need to know  
967 what the distribution of frequencies is in generation 0. One mathematically appealing choice

968 would be to assume that the population is at equilibrium at time 0, but this seemingly straight-  
 969 forward choice results in nonsensical conclusions. To see why, if the mutation rates are low and  
 970 selection is negligible, then at equilibrium, with extremely high probability the population will  
 971 either be in a state where the frequency of the LOF allele is very close to zero or in a state where  
 972 the frequency of the LOF allele is very close to one. If the mutation rates between the two alleles  
 973 are close to equal, then these two cases happen roughly equally often. That is, we would expect  
 974 there to be a  $\sim 50\%$  chance that the population is fixed or nearly fixed for the LOF mutation. If  
 975 there are multiple independently evolving sites at which an LOF could arise (or if there are many  
 976 more ways to mutate to an LOF state than a non-LOF state), then the chance that any of these sites  
 977 is fixed or nearly fixed for an LOF rapidly approaches 100%. Under this equilibrium assumption,  
 978 we thus reach the absurd conclusion that the mere act of observing a gene that is functional in a  
 979 majority of the population is overwhelming evidence that the gene is strongly selected for. An-  
 980 other way of viewing this is that in reality we can only observe genes that are functional in an  
 981 appreciable fraction of the population, and so we should somehow be conditioning on this event,  
 982 whereas the equilibrium assumption looks at a given randomly chosen stretch of DNA and asks  
 983 whether it could be a gene given some set of mutations. Indeed, any randomly chosen stretch of  
 984 DNA could be made a gene through a series of mutations, but for any given stretch it would be  
 985 extremely unlikely to be a functional gene, and the equilibrium assumption exactly captures how  
 986 rare this would be.

987 We instead use the equilibrium of another process as the initial condition, which avoids these  
 988 conceptual pitfalls. We assume the distribution of frequencies at generation 0 is the equilibrium  
 989 conditioned on the LOF allele never reaching fixation in the population. We then compute the like-  
 990 lihood of observing a given present-day frequency while continuing to condition on non-fixation  
 991 of the LOF allele. This assumption implies that no matter the current frequency of the LOF vari-  
 992 ant, we know that at some point in the past the population was fixed for the functional version of  
 993 the gene, and the LOF variant can thus be thought of as being “derived” and the non-LOF variant  
 994 “ancestral”. In the limit of infinitely low (but non-zero) mutation rates, this assumption become  
 995 equivalent to the commonly assumed “infinite sites” model commonly used to compute frequency  
 996 in population genetics [90]. In contrast to the infinite sites model, where the probability that any  
 997 given site is segregating must be 0, our model allows us to compute the probability that a given  
 998 site is segregating. Furthermore, we can easily model recurrent mutation which can be important  
 999 for sites with large mutation rates (such as CpGs) and large sample sizes [91], whereas under the  
 1000 infinite sites model each mutation necessarily happens at a unique position in the genome, ruling  
 1001 out the possibility of recurrent mutation. Below we will write  $p_{DTWF}(y | s_{het})$  for the probability  
 1002 mass function computed using this procedure, with “DTWF” representing Discrete-Time Wright-  
 1003 Fisher, and  $y$  being an observed LOF allele frequency.

1004 Equation 8 is easy to describe and simulate under, and a very similar model has been used  
 1005 in an approximate Bayesian computation approach to estimate  $s_{het}$  [4]. While simulation is easy,  
 1006 computing likelihoods under this model is difficult for large sample sizes, and unfortunately we  
 1007 need explicit likelihoods in our empirical Bayes approach. In recent work [16], we have developed  
 1008 an efficient method for computing likelihoods under this model. The key idea is that the above  
 1009 dynamics can be written as

$$\mathbf{v}_{k+1} = \mathbf{M}_k^T \mathbf{v}_k$$

1010 where  $\mathbf{v}_k$  is a vector of dimension  $2N + 1$  where entry  $i$  is the probability that there are  $i$  haploids

1011 that have the LOF allele in generation  $k$ , and  $\mathbf{M}_k$  is a matrix where row  $i$  is the the probability mass  
1012 function of the Binomial distribution in Equation 8 given that the allele frequency in generation  
1013  $k$  is  $i/2N_k$ . This formulation makes clear that we can obtain the likelihood of observing a given  
1014 frequency at present given some initial distribution by performing a series of matrix-vector multi-  
1015 plications. Naively this would be prohibitively slow as  $\mathbf{M}_k$  can be as large as  $10^7 \times 10^7$ , but in [16]  
1016 we show that  $\mathbf{M}_k$  is approximately highly structured — it is both approximately extremely sparse  
1017 and approximately extremely low rank. Combining these insights we can perform matrix-vector  
1018 multiplication that is provably accurate while reducing the runtime for matrix-vector multiplica-  
1019 tion from  $O(N_k^2)$  to  $O(N_k)$ . Similar insights can be used to speed up the computation of equilibria,  
1020 which we discuss in detail in [16]. Furthermore, as discussed above, we actually want to com-  
1021 pute likelihoods conditioned on non-fixation of the LOF allele, but that is as simple as setting the  
1022 column of  $\mathbf{M}_k$  corresponding to fixation to 0, and then renormalizing  $\mathbf{v}$ . We precompute these  
1023 likelihoods for each possible pair of mutation rates (to and from the LOF allele) across a range of  
1024  $s_{\text{het}}$  values (100 log-linearly spaced points between  $10^{-8}$  and 1, as well as 0). We describe how we  
1025 set the mutation rates and the population sizes implicit in  $\mathbf{M}_k$  below.

## 1026 Modeling misannotation of LOFs

1027 Under the likelihood described above, and as seen in Figure 2A, positions where a LOF variant  
1028 could occur, but no LOF alleles are observed are slight evidence in favor of selection, while high  
1029 frequency variants are extremely strong evidence against selection. Meanwhile, we suspect that  
1030 many variants that are annotated as causing LOF actually have little to no effect on the gene prod-  
1031 uct due to some form of misannotation. If these misannotated variants evolve effectively neutrally,  
1032 they can reach high frequencies and cause us to artifactually infer artificially low levels of selec-  
1033 tion. These misannotated variants can be particularly problematic for approaches that combine  
1034 frequencies across all LOFs within a gene to obtain an aggregate gene-level LOF frequency [1,2,4].

1035 LOEUF [12] and pLI [11] avoid this problem by throwing away all frequency information  
1036 except for whether an LOF is segregating or not. While this approach is more robust, the ignored  
1037 frequency information is extremely useful for estimating the strength of selection. For example,  
1038 consider a gene where we expect to see 5 unique LOFs under neutrality and we see 3 segregating  
1039 LOFs. This might seem like weak or negligible constraint ( $O/E = 0.6$ ), but if those 3 sites are all  
1040 highly mutable and the variants at those sites are each only present in a single individual, then it  
1041 is plausible that this gene is quite constrained.

1042 To take full advantage of the information in the LOF frequencies while remaining robust to  
1043 misannotation, we take a composite likelihood approach [92], closely related to the Poisson ran-  
1044 dom field assumption commonly used in population genetics [90]. We approximate gene-level  
1045 likelihoods as a product of variant level likelihoods

$$p^{(i)}(\mathbf{y}^{(i)} | s_{\text{het}}^{(i)}) \approx \prod_{j=1}^{J_i} p_{\text{variant}}(\mathbf{y}_j^{(i)} | s_{\text{het}}^{(i)}),$$

1046 where  $\mathbf{y}^{(i)}$  is a vector of the observed allele frequencies at each possible LOF site in gene  $i$ , and  
1047  $s_{\text{het}}^{(i)}$  is the selection coefficient for having a heterozygous loss-of-function of gene  $i$ . Under this  
1048 formulation, we can easily model misannotation by assuming that each LOF independently has



1049 some probability of being misannotated,  $p_{\text{miss}}$ , and that misannotated variants evolve neutrally:

$$p_{\text{variant}} \left( \mathbf{y}_j^{(i)} \mid s_{\text{het}}^{(i)} \right) = (1 - p_{\text{miss}}) p_{\text{DTWF}} \left( \mathbf{y}_j^{(i)} \mid s_{\text{het}}^{(i)} \right) + p_{\text{miss}} p_{\text{DTWF}} \left( \mathbf{y}_j^{(i)} \mid 0 \right).$$

1050 Using this formulation, we can take full advantage of the rich information included in the exact  
1051 sample frequencies of each LOF variant, while still being robust to occasional misannotation. In  
1052 practice, we precompute  $p_{\text{variant}}$  using a grid of  $p_{\text{miss}}$  values, and then to obtain the likelihood at  
1053 arbitrary values of  $s_{\text{het}}$  and  $p_{\text{miss}}$  we linearly interpolate in log-likelihood space. Below, we discuss  
1054 our approach for setting  $p_{\text{miss}}$ .

1055 Given a probability of misannoation, we can then calculate a posterior probability that any  
1056 given variant has been misannotated. We include a table of these misannotation probabilities for  
1057 all possible LOFs in Supplementary Table XXX.

1058 As an example of the importance of correcting for misannotation, we consider the case of the  
1059 gene PPFIA3 (ENSG00000177380). This gene has a LOEUF score of 0.12 and so appears very  
1060 constrained, but in an early version of our model where we did not incorporate variant mis-  
1061 annotation, we inferred a posterior mean value of  $s_{\text{het}}$  of  $\sim 2 \times 10^{-4}$ , which is right at the bor-  
1062 der of being nearly neutral. Inspecting the LOF data for this gene, we find that all potential  
1063 LOFs are either not observed or observed in a single individual, except for a single splice donor-  
1064 disrupting variant at 16% frequency. There are no obvious signs indicating that this variant is  
1065 misannotated (e.g., in terms of coverage or mappability). If we model misannotation, however,  
1066 we find that this variant is likely misannotated (posterior probability of misannotation  $> 99.999\%$ ),  
1067 and as a result we estimate extremely strong selection against gene loss-of-function (posterior  
1068 mean  $s_{\text{het}}$  of  $\sim 0.234$ ). Indeed, a single autosomal dominant missense variant in this gene is  
1069 suspected to have caused a number of severe symptoms including developmental delay, intel-  
1070 lectual disability, seizures, and macrocephaly in an Undiagnosed Diseases Network participant  
1071 (<https://undiagnosed.hms.harvard.edu/participants/participant-159/>) [93].

## 1072 Modeling the X chromosome

1073 We must slightly modify our model when applying it to the X chromosome. Because males only  
1074 have one copy of the X chromosome, there are only 3/4 as many X chromosomes as autosomes  
1075 (assuming an approximately equal sex ratio). As a result, when dealing with the X chromosome  
1076 we scale all population sizes to 3/4 of the size used for the autosomes (rounded to the nearest  
1077 integer). We also need to slightly modify the expected frequency in the next generation. We as-  
1078 sume haploid selection in males with strength  $s_{\text{hom}}$ , and diploid selection in females with selection  
1079 coefficients  $s_{\text{het}}$  and  $s_{\text{hom}}$  for individuals heterozygous and homozygous for the LOF variant re-  
1080 spectively. This selection results in modified allele frequencies in the pool of males and females,  
1081 and the we assume that each chromosome in the next generation has 1/3 probability of coming  
1082 from a male, and 2/3 probability of coming from a female. This means that the expected fre-  
1083 quency in the next generation is 1/3 times the post-selection frequency in males plus 2/3 times  
1084 the post-selection frequency in females. Variants within the pseudoautosomal regions on the X  
1085 are modeled identically to variants on the autosomes. Agarwal and colleagues also considered  
1086 selection on the X in the context of LOF variants, with a model similar to that described here [4].

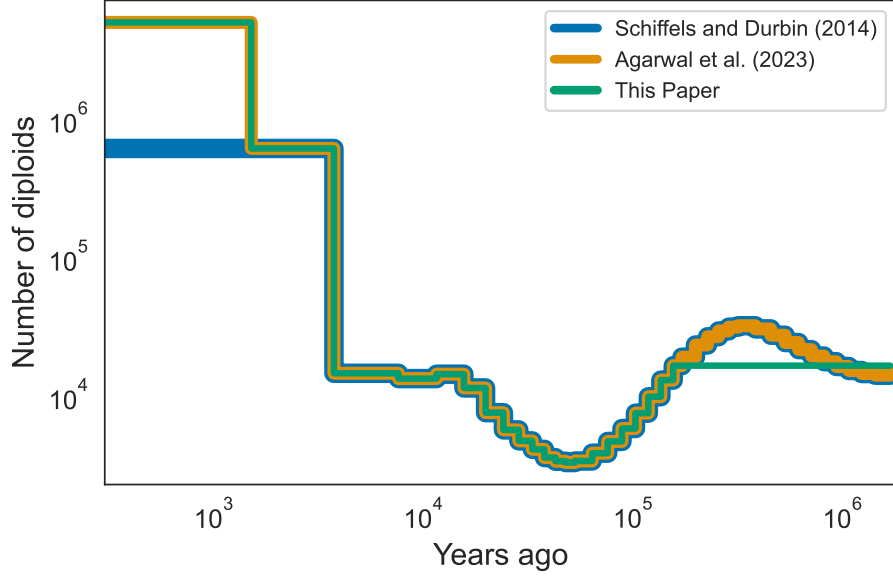
## 1087 Model parameters

1088 Our model has three key parameters — the mutation rate, the demographic model (i.e., population  
1089 sizes through time), and the probability that different variants are misannotated.

1090 We obtained mutation rates from gnomAD [12, Supplemental Dataset 10], which take into ac-  
1091 count trinucleotide context and methylation level (for CpG to TpG mutations). In our population  
1092 genetics model, we assume that there are only two alleles (a functional allele and an LOF allele),  
1093 whereas in reality there are four nucleotides. We approximate the rate of mutating from the func-  
1094 tional allele to the LOF allele as being the sum of the mutation rates from the reference nucleotide  
1095 to any nucleotide that might result in LOF. For example, if the reference allele is A, and either a  
1096 C or a T would result in LOF, then we say that the rate at which the functional allele mutates to  
1097 the LOF allele is the rate at which A mutates to C in this context plus the rate at which A mutates  
1098 to T in this context. For the rate of back mutation from the LOF allele to the functional allele, we  
1099 compute a weighted average of the rates of each possible LOF nucleotide back-mutating to any  
1100 possible non-LOF nucleotide, weighed by the probability that the original non-LOF nucleotide  
1101 mutated to that particular LOF nucleotide. Continuing our previous example, suppose A mutates  
1102 to C at rate  $1 \times 10^{-8}$  and A mutates to T at a rate  $1.5 \times 10^{-8}$ . Then conditioned on there having  
1103 been a single mutation resulting in a LOF variant, there is a  $1/2.5 = 0.4$  chance that the LOF is C  
1104 and 0.6 chance that the LOF is T. We then compute the back mutation rate as 0.4 times the rate at  
1105 which C mutates to A in this context plus the rate at which C mutates to G in this context (since  
1106 both A and G do not result in LOF) plus 0.6 times the rate at which T mutates to A in this con-  
1107 text plus the rate at which T mutates to G in this context. Implicitly this scheme assumes that the  
1108 flanking nucleotides in the trinucleotide context do not change, and we further assume that all  
1109 mutations resulting in CpGs result in unmethylated CpGs.

1110 For the population sizes in each generation, we used the “CEU” model inferred in [75] using  
1111 the 1000 Genomes Project data [94]. This model was also used in [4]. Population sizes under this  
1112 model are relatively constant before 5156 generations ago (approximately 155 thousand years ago)  
1113 and the effects of strong selection are relatively insensitive to all but the most recent population  
1114 sizes, so for a computational speedup we assumed that the population size was constant prior  
1115 to 5156 generations ago. Recently, [4] found that this CEU model underestimates the number  
1116 of low frequency variants and that changing the population size to 5,000,000 for the most recent  
1117 50 generations provides a better fit to the data. We used both demographic models and found  
1118 qualitatively similar results, with slightly better fit provided by the modified model, so we used  
1119 that demographic model for all subsequent analyses. In both cases, we modified the most ancient  
1120 population sizes, which are relatively constant, to be actually constant to speed up likelihood  
1121 calculations. The demographic models are presented in Supplementary Figure 1.

1122 The only remaining model parameter is  $p_{\text{miss}}$  the probability that any given LOF is misan-  
1123 notated. Throughout we focus on LOFs that either introduce early stop codons, disrupt splice  
1124 donors, or disrupts splice acceptors. Given that predicting which variants have these different  
1125 consequences involves different bioinformatic challenges, we inferred separate misannotation  
1126 probabilities  $p_{\text{miss}}^c$  for  $c \in \{\text{stop codon, splice donor, splice acceptor}\}$ . Below we write  $p_{\text{miss}}$  for the  
1127 collection of these three misannotation parameters. To get a rough estimate of these parameters  
1128 and avoid excessive computational burden, we took an h-likelihood approach [95,96]. That is, we  
1129 jointly maximized the likelihood across all genes with respect to their selective constraints as well



Supplementary Figure 1: CEU Demography inferred by Schiffels and Durbin [75], modified by Agarwal and colleagues [4], and further modified for this paper.

1130 as the the three misannotation probabilities that are shared across all genes:

$$\max_{p_{\text{miss}}, s_{\text{het}}^{(1)}, \dots, s_{\text{het}}^{(M)}} \sum_{i=1}^M \log p \left( \mathbf{y}^{(i)} \mid s_{\text{het}}^{(i)}, p_{\text{miss}} \right).$$

1131 This approach of just using the maximum likelihood estimates of  $s_{\text{het}}$  for each gene contrasts with  
 1132 the standard empirical Bayes approach, which would involve marginalizing out the unknown  $s_{\text{het}}$   
 1133 values. Yet, this marginalization step depends on the prior on  $s_{\text{het}}$ , which we learn via our NGBBoost  
 1134 framework. As a result, we would need to repeatedly run our NGBBoost framework as an inner loop  
 1135 to perform the standard empirical Bayes approach on  $p_{\text{miss}}$ . For our application, these values are  
 1136 nuisance parameters, and the results are relatively insensitive to their exact values so we opted for  
 1137 this simpler h-likelihood approach. Ultimately, we estimate that the probability of misannotation  
 1138 is 0.7%, 6.1%, and 8.4% for stop codons, splice donors, and splice acceptors respectively.

## 1139 C Feature processing and selection

1140 We compiled 10 types of gene features from several sources:

- 1141 1. *Gene structure*. Gene structure features were derived from GENCODE gene annotations (Re-  
1142 lease 39) [78]. Such features include the number of transcripts and, for the primary transcript  
1143 of each gene (the transcript tagged `Ensembl_canonical`), the number of exons as well as the  
1144 length and GC content of the transcript, total coding region, 5' UTR, and 3' UTR.
- 1145 2. *Gene expression*. We used gene features from 77 bulk and single-cell RNA-seq datasets, pro-  
1146 cessed and derived in [97]. These datasets can be grouped into 24 categories representing  
1147 tissues, cell types, and developmental stage (Table 6). For each dataset, features were de-  
1148 rived separately from all data and from individual cell clusters (for example, gene loadings  
1149 on principal components). In addition, features were derived from comparisons between  
1150 clusters (for example, t-statistics for differential expression). Finally, we include a metric,  $\tau$ ,  
1151 that summarizes the tissue-specificity of gene expression [98].
- 1152 3. *Biological pathways and Gene Ontology terms*. First, we included previously curated biological  
1153 pathway features [97,99]. In addition, to include GO terms that capture additional known  
1154 relationships between genes, we downloaded Biological Pathway (BP), Molecular Function  
1155 (MF), and Cellular Component (CC) terms [100] with at least 10 member genes using the  
1156 procedure described in [10]. Features for each gene were encoded as binary indicators of the  
1157 gene's membership in the pathways and GO terms.
- 1158 4. *Connectedness in protein-protein interaction (PPI) networks*. We included previously computed  
1159 measures of the connectedness of protein products of genes in PPI networks [10]. Connect-  
1160 edness was calculated as the number of interactions per protein weighted by the interaction  
1161 confidence scores.
- 1162 5. *Co-expression*. First, we included previously computed measures of the connectedness of  
1163 genes in co-expression networks [10], where connectedness measures the relative number  
1164 of neighbors of each gene in the network, averaged over tissues. Next, for each gene, we  
1165 derived features representing its co-expression with other genes (i.e. correlation in their ex-  
1166 pression levels across samples). To do this, we downloaded from the GeneFriends database  
1167 a co-expression network derived from GTEx RNA-seq samples [101,102], calculated the vari-  
1168 ance in the co-expression for each gene, and kept the 6,000 most variable genes. Then, we  
1169 included the co-expression with each of these 6,000 genes as a feature.
- 1170 6. *Gene regulatory landscape*. Gene regulatory features include the counts and properties of the  
1171 enhancers and promoters that regulate each gene. First, we included the number of pro-  
1172 moters per gene estimated by the FANTOM consortium using Cap Analysis of Gene Ex-  
1173 pression [10, 103]. Next, for each gene, we calculated the number, summed length, and  
1174 summed score of enhancer-to-gene links predicted using the Activity-By-Contact (ABC) ap-  
1175 proach [49,104], where an enhancer is considered linked to a gene if its ABC score is  $\geq 0.015$ .  
1176 We computed separate features for each of 131 biosamples. We also included features de-  
1177 rived by aggregating over all biosamples for both ABC enhancers and predicted enhancers

1178 from the Roadmap Epigenomics Consortium [10, 105, 106]—these feature include the num-  
1179 ber of biosamples with an active enhancer element, the total number of enhancer elements,  
1180 the total number of enhancer elements after taking merging enhancer domains, the total  
1181 length of the merged domains, and the average total enhancer length in an active cell type.  
1182 Finally, we included the enhancer-domain score for each gene [9] as a feature.

1183 7. *Conservation across species.* For each gene, we calculated the mean and 95th percentile phast-  
1184 Cons scores over the gene’s exons for multiple alignments of 7, 17, 20, 30, and 100 verte-  
1185 brate species to the human genome [107]. We downloaded phastCons Scores from [https://](https://hgdownload.soe.ucsc.edu/goldenPath/hg38/)  
1186 [hgdownload.soe.ucsc.edu/goldenPath/hg38/](https://hgdownload.soe.ucsc.edu/goldenPath/hg38/). In addition, we included the fraction of  
1187 coding sequence (CDS) or exons constrained across 240 mammals or 43 primates sequenced  
1188 in the Zoonomia project [108], with constraint determined by the per-base phyloP [109] or  
1189 phastCons score. Zoonomia data were downloaded from [https://figshare.com/articles/](https://figshare.com/articles/dataset/geneMatrix/13335548)  
1190 [dataset/geneMatrix/13335548](https://figshare.com/articles/dataset/geneMatrix/13335548).

1191 8. *Protein embedding features.* We included as features the embeddings learned by an autoen-  
1192 coder (ProtT5) trained on protein sequences [110]. Embeddings were downloaded from  
1193 <https://zenodo.org/record/5047020>. The embedding for each protein is a fixed-size vec-  
1194 tor that captures some of the protein’s biophysical and functional properties. For each gene  
1195 with more than one protein product, we averaged the embeddings of the proteins for that  
1196 gene.

1197 9. *Subcellular localization.* We included as features the subcellular localization of each pro-  
1198 tein and whether the protein is membrane-bound or soluble, as predicted by deep neu-  
1199 ral networks trained on the ProtT5 protein embeddings [110, 111]. Possible subcellular  
1200 classes included nucleus, cytoplasm, extracellular space, mitochondrion, cell membrane,  
1201 endoplasmatic reticulum, plastid, Golgi apparatus, lysosome or vacuole, and peroxisome.  
1202 Predictions were one-hot encoded, and for each gene with more than one protein product,  
1203 we summed the predictions for the gene’s proteins. Predictions were downloaded from  
1204 <https://zenodo.org/record/5047020>.

1205 10. *Missense constraint.* We included a measure of each gene’s average intolerance to missense  
1206 variants (UNEECON-G score) [112]. UNEECON-G scores incorporate variant-level features  
1207 to account for differences in the effects of missense variants on gene function.

1208 In addition to these 10 groups of features, we included a binary indicator for whether the  
1209 gene is located on the X chromosome. Genes in the pseudoautosomal regions were categorized as  
1210 autosomal.

1211 After compiling these features (total of 65,383), we performed feature selection to minimize  
1212 the practical complexity of training on such a large feature set and the complexity of the resulting  
1213 model. First, we removed features with zero variance and features where the Spearman corre-  
1214 lation of the feature values with  $O/E$  (the ratio of observed over expected unique LOF variants,  
1215 computed using gnomAD data) was less than 0.1 or had a nominal p-value  $\geq 0.05$ . Next, we per-  
1216 formed simultaneous feature selection and an initial round of hyperparameter tuning using the  
1217 shap-hypetune package, which uses Bayesian optimization to identify a set of features and hyper-  
1218 parameters that minimize the loss of a machine learning model fit on the training data. Specifically,  
1219 we fit gradient-boosted trees using XGBoost to predict  $O/E$  from the gene features; we chose to

1220 perform feature selection using XGBoost rather than NGBoost as training XGBoost models is sub-  
1221 stantially faster, and because we expect features/hyperparameters that perform well for XGBoost  
1222 to also perform well for NGBoost. For each set of hyperparameters, shap-hypetune performs back-  
1223 ward step-wise selection by removing the  $k$  least influential features (we chose  $k = 1000$  and  
1224 calculated influence using SHAP scores) at each step. Finally, we performed further feature se-  
1225 lection using shap-hypetune by fixing the hyperparameters and performing backward step-wise  
1226 selection with  $k = 50$ . Ultimately, we included 1,248 features in the model.

## 1227 D Estimating additional gene properties using GeneBayes

1228 GeneBayes is a flexible framework that can be used to infer other gene-level properties of interest  
1229 beyond  $s_{\text{het}}$ . In Figure 6, we presented a schematic of the key components of GeneBayes that users  
1230 should specify, which we describe in more detail now.

1231 First, users should specify the gene features to use as predictors. We expect the gene features  
1232 we use for  $s_{\text{het}}$  estimation to work well for other applications, but GeneBayes supports any choice  
1233 of features. In particular, GeneBayes can handle categorical and continuous features without fea-  
1234 ture scaling, as well as features with missing values.

1235 Next, users should specify the form of the prior distribution. GeneBayes supports the distri-  
1236 butions defined by the `distributions` package of PyTorch. GeneBayes also supports custom dis-  
1237 tributions, as long as they implement the methods used by GeneBayes (i.e. `log_prob` and `sample`)  
1238 and are differentiable within the PyTorch framework.

1239 Finally, users need to specify a likelihood function that relates their gene property of interest to  
1240 observed data. The likelihood can be specified in terms of a PyTorch distribution, or as a custom  
1241 function.

1242 After model training, GeneBayes outputs a per-gene posterior mean and 95% credible interval  
1243 for the property of interest. For each parameter in the prior, GeneBayes also outputs a metric for  
1244 each feature that represents the contribution of the feature to predictions of the parameter.

1245 In the next section, we describe in more detail the two example applications that we outlined  
1246 in Figure 6.

### 1247 Example applications

#### 1248 Differential expression

1249 In this example, users have estimates of log-fold changes in gene expression between conditions  
1250 and their standard errors from a differential expression workflow, and would like to estimate log-  
1251 fold changes with greater power (e.g. for lowly-expressed genes with noisy estimates).

1252 **Likelihood** We define  $\ell_{\text{DE}}^{(i)}$  and  $\ell_i$  as the estimated and true log-fold change in expression respec-  
1253 tively for gene  $i$ , and  $s_i$  as the standard error for the estimate. Then, we define the likelihood for  $\ell_i$   
1254 as

$$\ell_{\text{DE}}^{(i)} | \ell_i \sim \text{Normal}(\ell_i, s_i^2).$$

1255 **Prior** We describe two potential priors that one may choose to try. The first is a normal prior  
1256 with parameters  $\mu_i$  and  $\sigma_i$ :

$$\ell_i \sim \text{Normal}(\mu_i, \sigma_i^2).$$

1257 The second is a spike-and-slab prior with parameters  $\pi_i$ ,  $\mu_i$ , and  $\sigma_i$ , which assumes that gene  $i$

1258 only has a  $\pi_i$  probability of being differentially expressed:

$$z_i \sim \text{Bernoulli}(\pi_i)$$
$$\ell_i | z_i \sim \begin{cases} 0, & \text{if } z_i = 0 \\ \text{Normal}(\mu_i, \sigma_i^2), & \text{if } z_i = 1 \end{cases}$$

### 1259 **Variant burden tests**

1260 In this example, users have sequencing data from patients with a disease or (if calling *de novo*  
1261 mutations) sequencing data from family trios, and would like to identify genes with excess muta-  
1262 tional burden in patients (e.g. an excess of missense or LOF variants). One approach is to infer the  
1263 relative risk for each gene (denoted as  $\gamma_i$  for gene  $i$ ), defined as the expected ratio of the number  
1264 of variants in patients to the number of variants in healthy individuals.

1265 **Likelihood** Let  $E_i$  be the number of variants we expect to observe for gene  $i$  given the study  
1266 sample size and sequence-dependent mutation rates (e.g. expected counts obtained using the  
1267 mutational model developed by [84]). Next, let  $O_i$  be the number of variants observed in patients  
1268 for gene  $i$ . Then, we define the likelihood for  $\eta_i$  as

$$O_i | \eta_i \sim \text{Poisson}(\eta_i E_i).$$

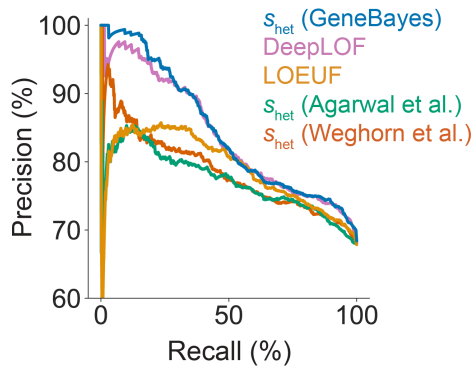
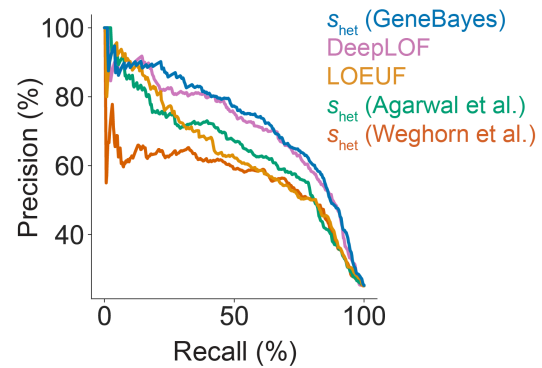
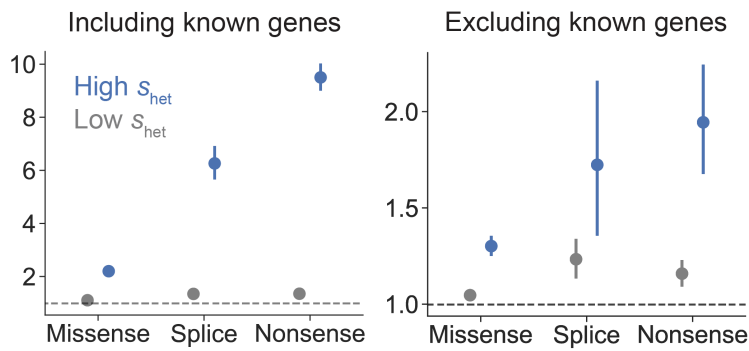
1269 **Prior** Because  $\eta_i$  is non-negative, one may want to choose a gamma prior with parameters  $\alpha_i$   
1270 and  $\beta_i$ :

$$\eta_i \sim \text{Gamma}(\alpha_i, \beta_i).$$



Gene	$s_{\text{het}}$	LOEUF
<i>RPL11</i>	0.75	0.3
<i>RPL18</i>	0.72	0.28
<i>RPL5</i>	0.71	0.17
<i>RPL35A</i>	0.67	0.41
<i>RPL15</i>	0.61	0.27
<i>RPL26</i>	0.61	0.38
<i>RPS15A</i>	0.61	0.56
<i>RPS7</i>	0.60	0.31
<i>RPS10</i>	0.60	0.27
<i>RPS26</i>	0.58	0.48
<i>RPL27</i>	0.56	0.48
<i>RPS24</i>	0.48	0.59
<i>RPS29</i>	0.40	1.2
<i>RPS27</i>	0.31	0.64
<i>RPS28</i>	0.26	0.8
<i>RPL35</i>	0.25	0.72

Supplementary Table 1: LOEUF and  $s_{\text{het}}$  for ribosomal proteins associated with Diamond-Blackfan anemia

**A** Classifying genes nonessential for survival *in vitro***B** Classifying developmental disorder genes**C** Enrichment of *de novo* developmental disorder mutations in constrained genes

Supplementary Figure 2: **Additional validation analyses.** **A)** Precision-recall curves comparing the performance of  $s_{het}$  estimates from GeneBayes against other constraint metrics in classifying non-essential genes. **B)** Precision-recall curves comparing the performance of  $s_{het}$  against other constraint metrics in classifying developmental disorder genes. **C)** Enrichment of *de novo* mutations in patients with developmental disorders, calculated as the observed number of mutations over the expected number under a null mutational model. We plot the enrichment of missense, splice, and nonsense variants in the 10% of genes considered most constrained by  $s_{het}$  (blue) and in all other genes (gray), including (left) and excluding (right) known developmental disorder genes. Bars represent 95% confidence intervals.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryTable2.txt](#)
- [SupplementaryTable3.tsv.zip](#)
- [SupplementaryTable4.xlsx](#)