

Supplemental material for:

Estimation of allele-specific fitness effects across human
protein-coding sequences and implications for disease

Yi-Fei Huang¹ and Adam Siepel¹

¹Simons Center for Quantitative Biology,
Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA

Details of the LASSIE model

Poisson Random Field model

LASSIE employs the Poisson Random Field (PRF) model to calculate the probability distribution of derived allele frequencies in the presence of natural selection, genetic drift, and new mutations (Sawyer and Hartl, 1992; Hartl et al., 1994; Williamson et al., 2005; Boyko et al., 2008; Evans et al., 2007). These calculations rely on an infinite sites assumption, which implies at most one mutation per site on the time scale of human population genetics (Kimura, 1969), thereby avoiding the complication of modeling multiple derived alleles per site. (Any sites with more than two alleles are removed from the input data; these sites are rare and only responsible for 0.3% segregating sites in the 1000 Genomes Project Yoruba data set.)

Let $f(y|S, \psi, \theta, t)$ be the probability density function for the derived allele frequency y at time t . Time is defined over the interval from $t = 0$, representing a deep ancestral population (prior to the emergence of any segregating polymorphisms), to the present, $t = t_{\text{current}}$. S represents the population-scaled selection coefficient, $S = 2N(0)s$, where $N(0)$ is the effective population at $t = 0$ and s is the genic selection coefficient associated with the derived mutation in question (see Evans et al. 2007). The genic selection model assumes that the fitnesses of heterozygous and homozygous carriers of the derived allele are equal to $1 + s$ and $1 + 2s$, respectively, relative to a fitness of 1 for homozygous carriers of the ancestral allele. ψ is a vector of demographic parameters that defines the relative effective population size at time t , which is denoted $\rho(t, \psi)$ (as detailed below). Finally, $\theta = 4N(0)\mu$ is the population-scaled mutation rate, where μ is the mutation rate per generation per nucleotide site.

To ease numerical computation, we additionally apply a transformation of the allele frequencies, $g(y|S, \psi, \theta, t) = y(1 - y)f(y|S, \psi, \theta, t)$. As demonstrated by Evans et al. (2007), $g(y|S, \psi, \theta, t)$ can therefore be calculated by solving the partial differential equation,

$$\frac{\partial}{\partial t}g(y|S, \psi, \theta, t) = \underbrace{-Sy(1-y)}_{\text{strength of selection}} \frac{\partial}{\partial y}g(y|S, \psi, \theta, t) + \underbrace{\frac{y(1-y)}{2\rho(t, \psi)}}_{\text{strength of drift}} \frac{\partial^2}{\partial y^2}g(y|S, \psi, \theta, t), \quad (1)$$

with boundary conditions $\lim_{y \downarrow 0} = \theta\rho(t, \psi)$ and $\lim_{y \uparrow 1} = 0$.

Equation 1 describes how allele frequencies change stochastically over time, in response to the population genetic parameters associated with natural selection (S), genetic drift (ψ), and mutation (θ). To calculate the distribution of allele frequencies in the modern population, $f(y|S, \psi, \theta, t_{\text{current}})$, LASSIE solves Equation 1 numerically using the Crank-Nicolson algorithm (Crank and Nicolson, 1947). More specifically, LASSIE discretizes the transformed allele frequency $g(y|S, \psi, \theta, t)$ into 1000 equal-size bins to form a discrete approximation of allele-frequencies. Similarly, it discretizes the scaled time t with a bin size of 0.0001. Then, it applies the Crank-Nicolson algorithm to solve Equation 1 iteratively, forward over time, to calculate $g(y|S, \psi, \theta, t_{\text{current}})$. Finally, the discretized density function of allele frequency in the modern population, $f(y|S, \psi, \theta, t_{\text{current}})$, is calculated using the inverse transformation, $f(y|S, \psi, \theta, t_{\text{current}}) = \frac{g(y|S, \psi, \theta, t_{\text{current}})}{y(1-y)}$.

Sampling distribution of derived allele frequencies

The function $f(y|S, \psi, \theta, t_{\text{current}})$ represents the population-level distribution of derived allele frequencies, but in practice, we can only obtain a finite number of samples from the population. Therefore, it is essential to specify the sampling distribution of derived allele frequencies. Let M_i be the (haploid) sample size at site i . For example, the high-coverage Yoruba data set used in this study consists of 51 unrelated individuals, so $M_i = 102$ for all the autosomal sites without missing data. (Missing data is naturally handled in this framework by setting M_i equal to the number of alleles actually available at site i .) The probability of observing a polymorphic site i with m_i copies of the derived allele is simply given by the expectation of the binomial sampling distribution with respect to the density of continuous population-level allele frequencies (Evans et al., 2007),

$$Q(m_i|S, \psi, \theta) = \int_0^1 \underbrace{\binom{M_i}{m_i} y^{m_i} (1-y)^{M_i-m_i}}_{\text{binomial distribution}} \underbrace{f(y|S, \psi, \theta, t_{\text{current}})}_{\text{population-level allele frequency}} dy, \quad 1 \leq m_i \leq M_i - 1. \quad (2)$$

LASSIE employs Equation 4.1.18 in *Numerical Recipes in C* (Press et al., 1992) to numerically calculate the integral $Q(m_i|S, \psi, \theta)$, using the discretization scheme defined for the Crank-Nicolson algorithm (above).

Equation 2 assumes that the ancestral allele is known. In practice, we consider uncertainty in the reconstructed ancestral allele, as previously estimated in a phylogenetic analysis (Gronau

et al., 2013; Arbiza et al., 2013) (see Online Methods). In particular, let q_i and q'_i be the probabilities of the reconstructed ancestral allele being identical to the reference allele and alternative (non-reference) allele, respectively, and let m_i be the number of observed alternative alleles. Then the sampling distribution of segregating mutations can be calculated as,

$$P(m_i|S, \psi, \theta) = q_i Q(m_i|S, \psi, \theta) + q'_i Q(M_i - m_i|S, \psi, \theta). \quad (3)$$

Note that, for simplicity and speed, we abandon explicit handling of uncertainty in ancestral alleles when using the model in the context of the mixture density network (see below).

Demographic model

To control the effect of population expansions on the distribution of allele frequencies, we employ a three-epoch demographic model in which $\rho(t, \psi)$ is a step function with two change points. Previous studies suggest that a simple demographic model with two to three epochs is powerful enough to account for the impact of human expansions on the site frequency spectrum (Williamson et al., 2005; Boyko et al., 2008; Racimo and Schraiber, 2014). In our three-epoch model, the vector of demographic parameters is denoted by $\psi = (N_1, N_2, t_1, t_2)$, in which N_1 and N_2 represent the relative effective population sizes of the second and the third epochs, respectively, and t_1 and t_2 indicate the durations of the second and the third epochs, respectively.

We estimate the “neutral” parameters in the model (ψ and θ) from sites putatively free from selection (see Online Methods) by forcing $S = 0$ and maximizing the likelihood of the PRF model. Because there are millions of neutral sites in the Yoruba data set, however, we approximately estimate the neutral parameters in two steps. First, we estimate ψ using only the polymorphic sites in neutral regions (with $1 \leq m_i \leq M_i - 1$). Let $P'(m_i|S = 0, \psi, \theta)$ represent the sampling distribution for these polymorphic sites, with,

$$P'(m_i|S = 0, \psi, \theta) = \frac{P(m_i|S = 0, \psi, \theta)}{\sum_{n=1}^{M_i-1} P(n|S = 0, \psi, \theta)}. \quad (4)$$

ψ is estimated by maximizing a composite likelihood function for these sites alone, ignoring linkage between sites,

$$\hat{\psi} = \arg \max_{\psi} \prod_{i \in \mathcal{S}} P'(m_i|S = 0, \psi, \theta = c), \quad (5)$$

where \mathcal{S} is the set of segregating sites and c is an arbitrary small positive constant. Importantly, this estimator for $\hat{\psi}$ is invariant to the choice of c . This step provides a good approximation under the PRF model (which also ignores linkage), because the probability of the monomorphic sites depends primarily on θ and only very weakly on ψ .

In the second step, we fix $\psi = \hat{\psi}$ and estimate θ using both monomorphic and polymorphic neutral sites. However, to reduce computational cost, we consider only a random sample of 5% of monomorphic and polymorphic sites. Because θ is a single scalar parameter and the data set is large, the downsampling procedure has a negligible impact on the accuracy of the estimated θ . Specifically, we calculate the probability of observing a monomorphic site if both monomorphic and polymorphic sites are included as,

$$P(m_i = 0 | S = 0, \hat{\psi}, \theta) = 1 - \sum_{m_i=1}^{M_i-1} P(m_i | S = 0, \hat{\psi}, \theta). \quad (6)$$

Then we estimate θ by maximizing the composite likelihood,

$$\hat{\theta} = \arg \max_{\theta} \prod_{i \in \mathcal{R}} P(m_i | S = 0, \hat{\psi}, \theta), \quad (7)$$

where \mathcal{R} is our subsample of monomorphic and polymorphic “neutral” sites.

Mixture model for inferring representative selection coefficients

Given the estimates of the neutral parameters, $\hat{\theta}$ and $\hat{\psi}$, we fit a three-component mixture model to represent the global distribution of fitness effects based on all polymorphic and monomorphic sites in coding regions. The first component in this model describes neutral evolution and its selection coefficient S_0 is fixed to 0 by definition. The second and the third components represent weak negative selection ($S_1 < 0$) and strong negative selection ($S_2 < S_1$), respectively. Let w_0 , w_1 , and w_2 represent the probabilities of these three mixture components, respectively. All parameters are estimated by maximum likelihood,

$$(\hat{w}_0, \hat{w}_1, \hat{w}_2, \hat{S}_1, \hat{S}_2) = \arg \max_{w_0, w_1, w_2, S_1, S_2} \prod_i \sum_{n=0}^2 w_n P(m_i | S_n, \hat{\psi}, \hat{\theta}) \quad (8)$$

subject to the linear constraint $\sum_{n=0}^2 w_n = 1$ as well as the constraints that $S_0 = 0$, $S_1 < 0$, and $S_2 < S_1$. In this likelihood function, $P(m_i \neq 0 | S_n, \hat{\psi}, \hat{\theta})$ is defined in Equation 3 and describes the

probability of observing a polymorphic site, and $P(m_i = 0|S_n, \hat{\psi}, \hat{\theta}) = 1 - \sum_{m_i=1}^{M_i-1} P(m_i|S_n, \hat{\psi}, \hat{\theta})$ describes the probability of observing a monomorphic site given a selection coefficient S_n . The estimated representative selection coefficients, $\hat{S}_0 = 0$, \hat{S}_1 , and \hat{S}_2 , are then fixed in the mixture density network described below. Finally, the unscaled selection coefficients \hat{s}_n that LASSIE reports are calculated using the equation $\hat{s}_n = \frac{2\hat{\mu}\hat{S}_n}{\hat{\theta}}$, where $\hat{\mu} = 1.26 \times 10^{-8}$ is the estimated human non-CpG mutation rate per site per generation (Rasmussen et al., 2014).

Mixture density network architecture

In the mixture density network (Bishop, 1994) for inferring allele-specific selection coefficients, we assume that the allele-specific weights of selection coefficients can be inferred from the vector of genomic features, \mathbf{X}^{ijk} , for a mutation event from ancestral allele j to derived allele k at site i . We denote \mathbb{P}^{ijk} as the vector of allele-specific probabilities of being under neutral evolution, weak selection, or strong selection. We assume that \mathbb{P}^{ijk} is determined by a mixture density network, which, in our experiments, either includes a single hidden layer or no hidden layers. If the mixture density network includes a hidden layer, the hidden layer is defined by,

$$\mathbf{H}^{ijk} = \text{dropout}(\text{ReLU}(\mathbf{X}^{ijk} \cdot \mathbf{W}_{\text{hidden}})) \quad (9)$$

where $\mathbf{W}_{\text{hidden}}$ denotes the weights and bias terms while **ReLU** and **dropout** denote the rectified linear layer (Nair and Hinton, 2010) and dropout layer (Srivastava et al., 2014), respectively. The rectified linear layer serves as an activation function and the dropout layer serves as a regularizer for preventing overfitting.

At the top layer of the mixture density network, LASSIE first calculates an affine transformation of the hidden features and then employs a softmax function to transform the affine transformation to a normalized probability vector, which defines the allele-specific weights of selection components. This probability vector is defined by,

$$\mathbb{P}^{ijk} = \text{softmax}(\mathbf{H}^{ijk} \cdot \mathbf{W}_{\text{output}}), \quad (10)$$

where $\mathbf{W}_{\text{output}}$ denotes the weights and biases associated with the hidden feature vector \mathbf{H}^{ijk} , and **softmax** denotes the softmax layer. If no hidden layer is included, Equation 9 is simply replaced by the identity transformation, $\mathbf{H}^{ijk} \equiv \mathbf{X}^{ijk}$.

Objective function of the mixture density network

To estimate the parameters $\mathbf{W}_{\text{output}}$ and $\mathbf{W}_{\text{hidden}}$ in the mixture density network, we need to design a loss function that captures the discrepancy between the predictions from the mixture density network and the observed polymorphism patterns. We define this loss function as the negative logarithmic value of a likelihood function derived from the PRF model. While uncertainty in ancestral alleles can be considered (see above), here we simply assume that the ancestral allele is known for simplicity and speed. Accordingly, only sites with unambiguous ancestral alleles ($q_i > 0.98$ or $q'_i > 0.98$) are used in the training of the mixture density network. The scaled mutation rate per site is fixed at $\hat{\theta}$ and we assume an equal mutation rate to each alternative allele, so the scaled mutation rate for each possible derived allele is equal to $\frac{\hat{\theta}}{3}$. Under these assumptions, the probability of observing a mutation from ancestral allele j to derived allele k at a polymorphic site i can be calculated as,

$$\mathcal{L}^i = \frac{1}{3} \sum_{n=0}^2 \mathbb{P}_n^{ijk} Q(m_i | \hat{S}_n, \hat{\psi}, \hat{\theta}), \quad (11)$$

where $Q(m_i | \hat{S}_n, \hat{\psi}, \hat{\theta})$ is defined in Equation 2 and \mathbb{P}_n^{ijk} is the probability that the mutation belongs to selection component n . The probability of observing a monomorphic site is then equal to one minus the total probability of observing a mutation at this site,

$$\begin{aligned} \mathcal{L}^i &= 1 - \sum_{k \in \{A,G,C,T\} \setminus \{j\}} \sum_{n=0}^2 \sum_{m=1}^{M-1} \frac{1}{3} \mathbb{P}_n^{ijk} Q(m | \hat{S}_n, \hat{\psi}, \hat{\theta}) \\ &= \frac{1}{3} \sum_{k \in \{A,G,C,T\} \setminus \{j\}} \sum_{n=0}^2 \mathbb{P}_n^{ijk} - \frac{1}{3} \sum_{k \in \{A,G,C,T\} \setminus \{j\}} \sum_{n=0}^2 \sum_{m=1}^{M-1} \mathbb{P}_n^{ijk} Q(m | \hat{S}_n, \hat{\psi}, \hat{\theta}) \\ &= \frac{1}{3} \sum_{k \in \{A,G,C,T\} \setminus \{j\}} \sum_{n=0}^2 \mathbb{P}_n^{ijk} \left[1 - \sum_{m=1}^{M-1} Q(m | \hat{S}_n, \hat{\psi}, \hat{\theta}) \right] \\ &= \frac{1}{3} \sum_{k \in \{A,G,C,T\} \setminus \{j\}} \sum_{n=0}^2 \mathbb{P}_n^{ijk} Q(m = 0 | \hat{S}_n, \hat{\psi}, \hat{\theta}), \end{aligned} \quad (12)$$

where $Q(m = 0 | \hat{S}_n, \hat{\psi}, \hat{\theta}) = 1 - \sum_{m=1}^{M-1} Q(m | \hat{S}_n, \hat{\psi}, \hat{\theta})$. It is worth noting that $Q(m | \hat{S}_n, \hat{\psi}, \hat{\theta})$ is not dependent on the parameters in the mixture density network and, therefore, can be precomputed and cached for efficient evaluation of the likelihood function. Assuming independence across sites, the loss function for a mini-batch of sites follows

$$\text{loss} = - \frac{\sum_i \log(\mathcal{L}^i)}{k} \quad (13)$$

in which k is the number of sites in a mini-batch and \mathcal{L}^i is defined in Equation 11 and 12.

References

- Arbiza, L., Gronau, I., Aksoy, B. A., Hubisz, M. J., Gulko, B., Keinan, A., and Siepel, A. (2013). Genome-wide inference of natural selection on human transcription factor binding sites. *Nature Genetics* **45**:723–729.
- Bishop, C. M. (1994). Mixture density networks. Technical report, Aston University.
- Boyko, A. R., Williamson, S. H., Indap, A. R., Degenhardt, J. D., Hernandez, R. D., Lohmueller, K. E., Adams, M. D., Schmidt, S., Sninsky, J. J., Sunyaev, S. R., et al. (2008). Assessing the evolutionary impact of amino acid mutations in the human genome. *PLoS Genet* **4**:e1000083.
- Crank, J. and Nicolson, P. (1947). A practical method for numerical evaluation of solutions of partial differential equations of the heat-conduction type. *Mathematical Proceedings of the Cambridge Philosophical Society* **43**:50–67.
- Evans, S. N., Shvets, Y., and Slatkin, M. (2007). Non-equilibrium theory of the allele frequency spectrum. *Theoretical Population Biology* **71**:109 – 119.
- Findlay, G. M., Daza, R. M., Martin, B., Zhang, M. D., Leith, A. P., Gasperini, M., Janizek, J. D., Huang, X., Starita, L. M., and Shendure, J. (2018). Accurate classification of BRCA1 variants with saturation genome editing. *Nature* **562**:217–222.
- Grantham, R. (1974). Amino acid difference formula to help explain protein evolution. *Science* **185**:862–864.
- Gronau, I., Arbiza, L., Mohammed, J., and Siepel, A. (2013). Inference of natural selection from interspersed genomic elements based on polymorphism and divergence. *Molecular Biology and Evolution* **30**:1159–1171.
- Hartl, D. L., Moriyama, E. N., and Sawyer, S. A. (1994). Selection intensity for codon bias. *Genetics* **138**:227–234.
URL: <http://www.genetics.org/content/138/1/227>
- Kimura, M. (1969). The number of heterozygous nucleotide sites maintained in a finite population due to steady flux of mutations. *Genetics* **61**:893–903.
- Liu, X., Jian, X., and Eric, B. (2013). dbNSFP v2.0: A database of human nonsynonymous snvs and their functional predictions and annotations. *Human Mutation* **34**:E2393–E2402.
- Nair, V. and Hinton, G. E. (2010). Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML'10*, pages 807–814. Omnipress, USA.
- Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R., and Siepel, A. (2010). Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Research* **20**:110–121.
- Press, W., Teukolsky, S., Vetterling, W., and Flannery, B. (1992). *Numerical Recipes in C*. Cambridge University Press, 2nd edition.
- Racimo, F. and Schraiber, J. G. (2014). Approximation to the distribution of fitness effects across functional categories in human segregating polymorphisms. *PLoS Genet* **10**:e1004697.

- Rasmussen, M. D., Hubisz, M. J., Gronau, I., and Siepel, A. (2014). Genome-wide inference of ancestral recombination graphs. *PLOS Genetics* **10**:e1004342.
- Roadmap Epigenomics Consortium, Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., et al. (2015). Integrative analysis of 111 reference human epigenomes. *Nature* **518**:317–330.
- Sawyer, S. A. and Hartl, D. L. (1992). Population genetics of polymorphism and divergence. *Genetics* **132**:1161–1176.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**:1929–1958.
- Van Nostrand, E. L., Pratt, G. A., Shishkin, A. A., Gelboin-Burkhart, C., Fang, M. Y., Sundararaman, B., Blue, S. M., Nguyen, T. B., Surka, C., Elkins, K., et al. (2016). Robust transcriptome-wide discovery of rna-binding protein binding sites with enhanced clip (eclip). *Nature Methods* **13**:508.
- Williamson, S. H., Hernandez, R., Fledel-Alon, A., Zhu, L., Nielsen, R., and Bustamante, C. D. (2005). Simultaneous inference of selection and population growth from patterns of variation in the human genome. *Proceedings of the National Academy of Sciences* **102**:7882–7887.
- Wong, W. C., Kim, D., Carter, H., Diekhans, M., Ryan, M. C., and Karchin, R. (2011). CHASM and SNVBox: toolkit for detecting biologically important single nucleotide mutations in cancer. *Bioinformatics* **27**:2147–2148.
- Xiong, H. Y., Alipanahi, B., Lee, L. J., Bretschneider, H., Merico, D., Yuen, R. K., Hua, Y., Gueroussov, S., Najafabadi, H. S., Hughes, T. R., et al. (2015). The human splicing code reveals new insights into the genetic determinants of disease. *Science* **347**:1254806.

Supplemental Table S1: Genomic features for LASSIE (hg19 assembly)

Feature group	Feature name	Type	Reference	Note
Variant category	Stop-gain	Binary	Liu et al. (2013)	Indicate if a mutation results in a gain of stop codon
	Stop-loss	Binary	Liu et al. (2013)	Indicate if a mutation results in a loss of stop codon
	Missense	Binary	Liu et al. (2013)	Indicate if a mutation results in a substitution of amino acid
Sequence conservation	SIFT prediction	Binary	Liu et al. (2013)	Binary prediction of deleteriousness from SIFT
	LRT prediction	Binary	Liu et al. (2013)	Binary prediction of deleteriousness from LRT
	MA prediction	Binary	Liu et al. (2013)	Binary prediction of deleteriousness from Mutation Assessor
	PROVEAN prediction	Binary	Liu et al. (2013)	Binary prediction of deleteriousness from PROVEAN
	SLR score	Binary	Liu et al. (2013)	Raw SLR score
	SIFT score	Numeric	Liu et al. (2013)	Raw SIFT score
	LRT omega	Numeric	Liu et al. (2013)	Raw LRT score
	MA score	Numeric	Liu et al. (2013)	Raw Mutation Assessor score
	PROVEAN score	Numeric	Liu et al. (2013)	Raw PROVEAN score
	Grantham score	Numeric	Grantham (1974)	Raw Grantham score
	HMM entropy	Numeric	Wong et al. (2011)	HMM entropy score from SNVBox
	HMM relative entropy	Numeric	Wong et al. (2011)	HMM relative entropy score from SNVBox
	dscore	Numeric	Liu et al. (2013)	Dscore from PolyPhen-2
	Primate phyloP score	Numeric	Pollard et al. (2010)	Primate phyloP conservation score
	Mammalian phyloP score	Numeric	Pollard et al. (2010)	Mammalian phyloP conservation score
Vertebrate phyloP score	Numeric	Pollard et al. (2010)	Vertebrate phyloP conservation score	
Structural information	PredRSAB	Numeric	Wong et al. (2011)	Probability of the residue being buried
	PredRSAI	Numeric	Wong et al. (2011)	Probability of the residue being intermediately exposed
	PredRSAE	Numeric	Wong et al. (2011)	Probability of the residue being exposed
	PredBFAEF	Numeric	Wong et al. (2011)	Probability that the residue's backbone is flexible
	PredBFAFM	Numeric	Wong et al. (2011)	Probability that the residue's backbone is intermediately flexible
	PredBFASF	Numeric	Wong et al. (2011)	Probability that the residue's backbone is stiff
	PredStabilityH	Numeric	Wong et al. (2011)	Probability that the residue strongly stabilizes folding
	PredStabilityM	Numeric	Wong et al. (2011)	Probability that the residue stabilizes folding
	PredStabilityL	Numeric	Wong et al. (2011)	Probability that the residue destabilizes folding
	PredSSE	Numeric	Wong et al. (2011)	Probability that the secondary structure of the residue is strand
	PredSSH	Numeric	Wong et al. (2011)	Probability that the secondary structure of the residue is helix
	PredSSC	Numeric	Wong et al. (2011)	Probability that the secondary structure of the residue is loop
	Regulatory information	SPIDEX	Numeric	Xiong et al. (2015)
Maximum RNA-seq signal		Numeric	Roadmap Epigenomics Consortium et al. (2015)	Maximum RNA-seq signal from the Roadmap Epigenomics Project

Supplemental Table S2: Model fitting of the mixture density network. All coding sites on chromosome 1 were used as the held-out test data.

Number of hidden layers	Average loss in the held-out test data
No hidden layer (linear)	0.0321644
One hidden layer (Nonlinear)	0.0322484

Supplemental Table S3: Top 10 most enriched Gene Ontology (molecular function) terms among the 1,118 genes under enhanced selection.

Category	Fold enrichment	<i>p</i> -value	FDR
GABA receptor activity (GO:0016917)	5.93	3.63E-04	3.53E-03
voltage-gated potassium channel activity (GO:0005249)	5.07	5.78E-08	1.53E-06
glutamate receptor activity (GO:0008066)	4.98	4.67E-05	6.18E-04
mRNA binding (GO:0003729)	3.85	2.86E-09	1.32E-07
voltage-gated ion channel activity (GO:0005244)	3.73	1.10E-06	2.04E-05
translation initiation factor activity (GO:0003743)	3.46	3.24E-03	2.40E-02
adenylate cyclase activity (GO:0004016)	2.91	4.94E-04	4.57E-03
translation regulator activity (GO:0045182)	2.88	8.03E-04	7.08E-03
chromatin binding (GO:0003682)	2.87	1.86E-06	3.13E-05
ligand-gated ion channel activity (GO:0015276)	2.85	2.01E-05	2.86E-04

Supplemental Table S4: Top 10 most enriched Reactome pathways among the 1,118 genes under enhanced selection.

Category	Fold enrichment	<i>p</i> -value	FDR
Cohesin Loading onto Chromatin (R-HSA-2470946)	10.37	7.87E-05	1.04E-03
CREB phosphorylation through the activation of Adenylate Cyclase (R-HSA-442720)	10.37	3.12E-03	2.02E-02
GABA A receptor activation (R-HSA-977441)	9.22	4.11E-05	5.91E-04
PTK6 Regulates RHO GTPases, RAS GTPase and MAP kinases (R-HSA-8849471)	9.08	1.36E-04	1.58E-03
HuR (ELAVL1) binds and stabilizes mRNA (R-HSA-450520)	8.89	4.55E-04	4.46E-03
Adenylate cyclase activating pathway (R-HSA-170660)	8.30	5.22E-03	2.99E-02
Unblocking of NMDA receptor, glutamate binding and activation (R-HSA-438066)	8.15	3.24E-06	6.18E-05
CREB phosphorylation through the activation of CaMKII (R-HSA-442729)	7.98	1.05E-05	1.75E-04
Interleukin-21 signaling (R-HSA-9020958)	7.78	7.37E-04	6.67E-03
PKA-mediated phosphorylation of CREB (R-HSA-111931)	7.54	1.09E-04	1.36E-03

Supplemental Table S5: Enriched Gene Ontology (molecular function) terms among the 773 genes under relaxed selection.

Category	Fold enrichment	<i>p</i> -value	FDR
metallopeptidase activity (GO:0008237)	3.57	2.08E-03	4.81E-02
ATPase activity, coupled to transmembrane movement of substances (GO:0042626)	3.12	1.80E-03	6.67E-02
oxidoreductase activity (GO:0016491)	2.82	1.04E-10	1.93E-08
catalytic activity (GO:0003824)	1.45	1.81E-09	1.67E-07

Supplemental Table S6: Top 10 most enriched Reactome pathways among the 773 genes under relaxed selection.

Category	Fold enrichment	<i>p</i> -value	FDR
Melanin biosynthesis (R-HSA-5662702)	15.87	1.42E-04	2.29E-02
Eicosanoids (R-HSA-211979)	12.34	1.79E-05	5.35E-03
Fructose metabolism (R-HSA-5652084)	11.33	4.03E-04	4.45E-02
Laminin interactions (R-HSA-3000157)	6.04	2.07E-04	2.41E-02
Cytochrome P450 - arranged by substrate type (R-HSA-211897)	5.41	1.17E-06	8.20E-04
Phase I - Functionalization of compounds (R-HSA-211945)	4.31	8.54E-07	8.95E-04
Collagen formation (R-HSA-1474290)	3.84	3.82E-05	8.01E-03
Collagen biosynthesis and modifying enzymes (R-HSA-1650814)	3.71	5.19E-04	4.94E-02
Diseases of metabolism (R-HSA-5668914)	3.27	4.77E-04	4.77E-02
Biological oxidations (R-HSA-211859)	2.82	6.63E-06	3.48E-03

Supplemental Table S7: Top 10 most enriched Gene Ontology (molecular function) terms among the 1,118 genes under enhanced selection. To control gene length as a possible confounding factor, foreground genes are matched with background genes by the number of potential mutations.

Category	Fold enrichment	<i>p</i> -value	FDR
Histone-lysine N-methyltransferase activity (GO:0018024)	19.76	3.46E-03	3.04E-02
Histone methyltransferase activity (GO:0042054)	19.76	3.46E-03	2.99E-02
Transcriptional repressor activity, RNA polymerase II transcription regulatory region sequence-specific DNA binding (GO:0001227)	11.12	3.94E-05	5.45E-04
Voltage-gated potassium channel activity (GO:0005249)	9.88	1.56E-09	5.04E-08
Glutamate receptor activity (GO:0008066)	7.14	9.05E-06	1.41E-04
Glutamate binding (GO:0016595)	7.14	9.05E-06	1.37E-04
Signal sequence binding (GO:0005048)	7.06	1.09E-04	1.39E-03
Voltage-gated cation channel activity (GO:0022843)	6.59	7.25E-04	7.80E-03
Protein serine/threonine/tyrosine kinase activity (GO:0004712)	6.35	3.88E-04	4.37E-03
MAP kinase kinase activity (GO:0004708)	6.35	3.88E-04	4.27E-03

Supplemental Table S8: Top 10 most enriched Reactome pathways among the 1,118 genes under enhanced selection. To control gene length as a possible confounding factor, foreground genes are matched with background genes by the number of potential mutations.

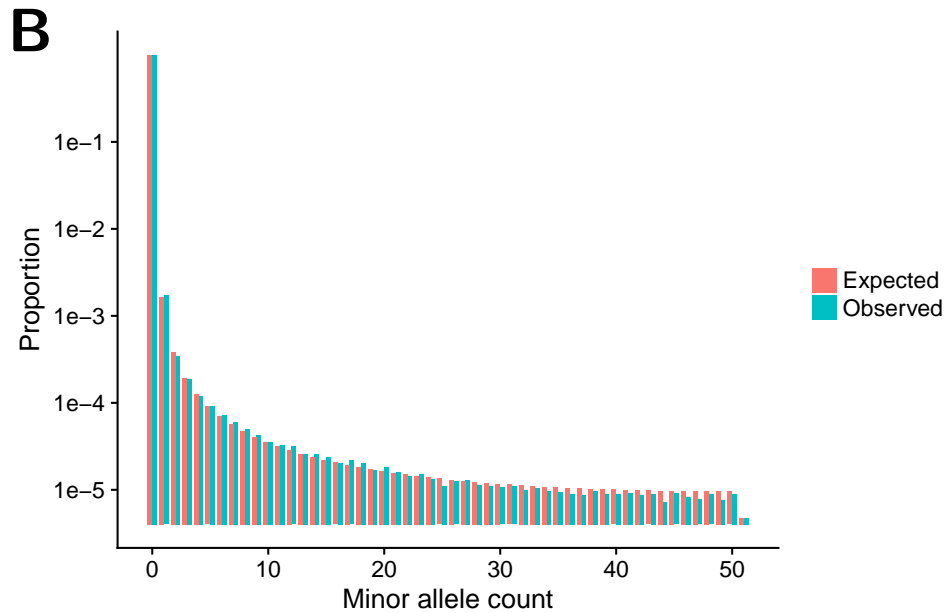
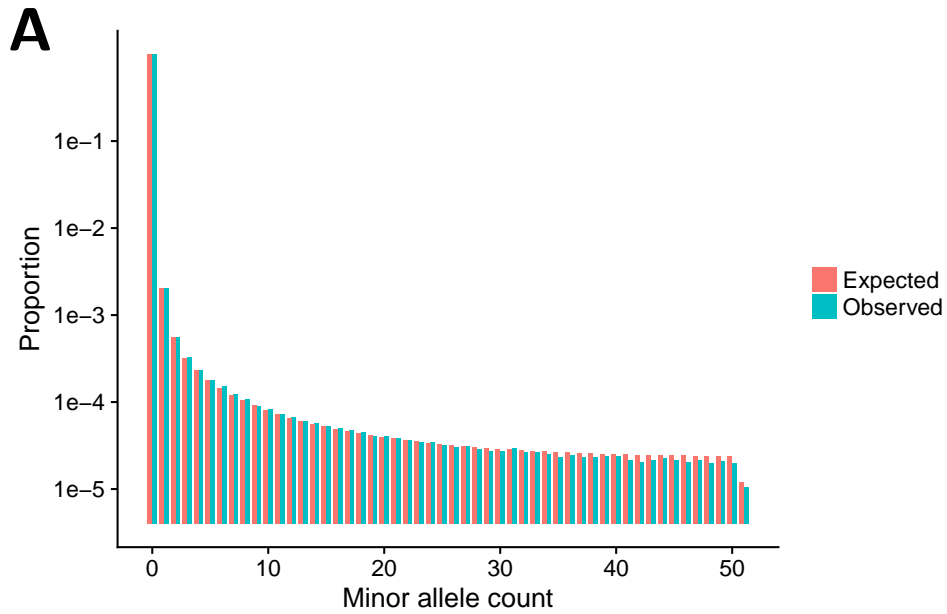
Category	Fold enrichment	<i>p</i> -value	FDR
PKA-mediated phosphorylation of CREB (R-HSA-111931)	39.53	4.84E-06	6.66E-05
GABA A receptor activation (R-HSA-977441)	39.53	4.84E-06	6.61E-05
PKA activation (R-HSA-163615)	34.59	2.57E-05	3.14E-04
PTK6 Regulates RHO GTPases, RAS GTPase and MAP kinases (R-HSA-8849471)	34.59	2.57E-05	3.12E-04
HuR (ELAVL1) binds and stabilizes mRNA (R-HSA-450520)	29.64	1.35E-04	1.29E-03
Glucagon signaling in metabolic regulation (R-HSA-163359)	29.64	1.35E-04	1.28E-03
Signaling by Activin (R-HSA-1502540)	29.64	1.35E-04	1.28E-03
PKA activation in glucagon signalling (R-HSA-164378)	29.64	1.35E-04	1.27E-03
Activated PKN1 stimulates transcription of AR (androgen receptor) regulated genes KLK2 and KLK3 (R-HSA-5625886)	29.64	1.35E-04	1.27E-03
ERCC6 (CSB) and EHMT2 (G9a) positively regulate rRNA expression (R-HSA-427389)	29.64	3.21E-08	9.03E-07

Supplemental Table S9: Enriched Gene Ontology (molecular function) terms among the 773 genes under relaxed selection. To control gene length as a possible confounding factor, foreground genes are matched with background genes by the number of potential mutations.

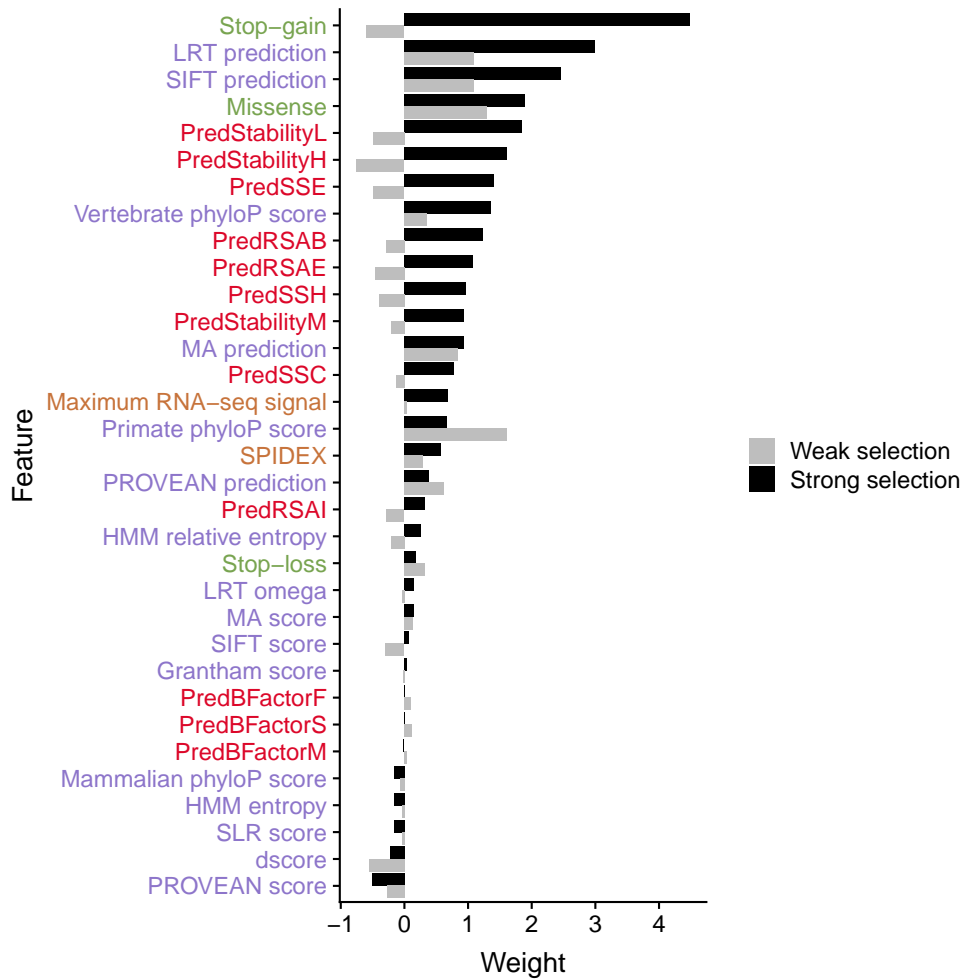
Category	Fold enrichment	<i>p</i> -value	FDR
Oxidoreductase activity (GO:0016491)	3.05	5.94E-10	2.86E-07

Supplemental Table S10: Top 10 most enriched Reactome pathways among the 773 genes under relaxed selection. To control gene length as a possible confounding factor, foreground genes are matched with background genes by the number of potential mutations.

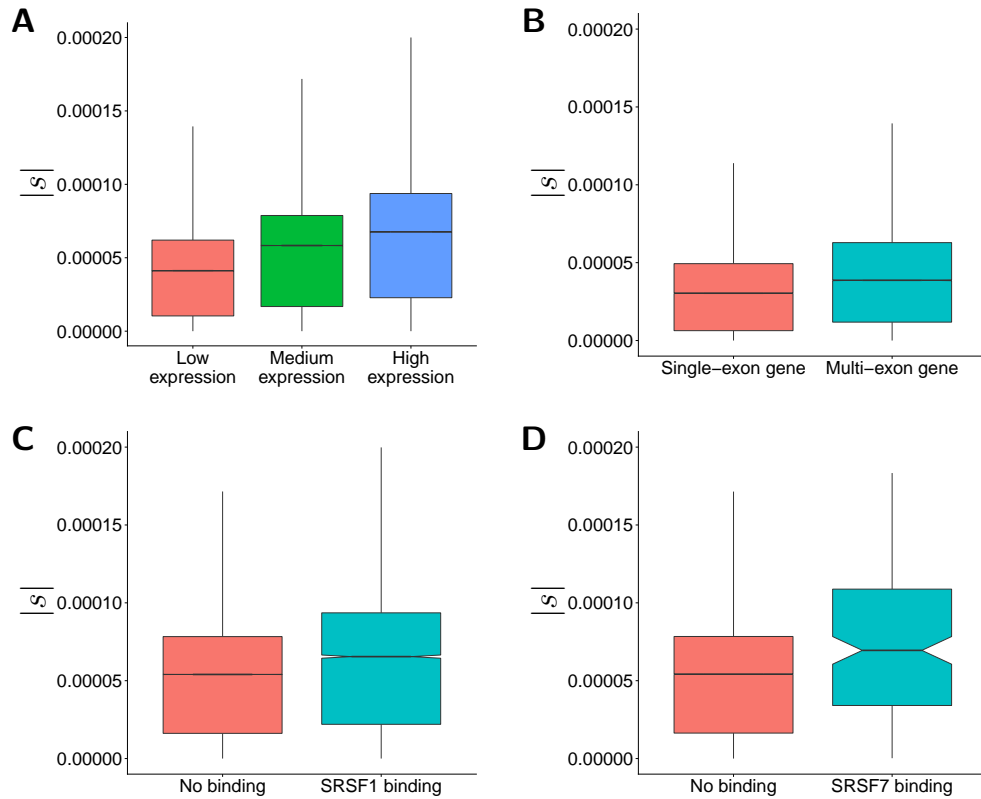
Category	Fold enrichment	<i>p</i> -value	FDR
Eicosanoids (R-HSA-211979)	34.85	2.45E-05	7.47E-03
Cytochrome P450 - arranged by substrate type (R-HSA-211897)	6.79	2.45E-06	8.99E-04
Arachidonic acid metabolism (R-HSA-2142753)	6.40	3.66E-04	5.16E-02
Histidine, lysine, phenylalanine, tyrosine, proline and tryptophan catabolism (R-HSA-6788656)	6.40	3.66E-04	4.79E-02
Phase I - Functionalization of compounds (R-HSA-211945)	4.98	2.08E-06	9.52E-04
Biological oxidations (R-HSA-211859)	4.34	9.74E-08	8.92E-05
Fatty acid metabolism (R-HSA-8978868)	4.07	2.48E-05	6.50E-03
Metabolism of amino acids and derivatives (R-HSA-71291)	3.25	1.01E-06	6.14E-04
Metabolism of carbohydrates (R-HSA-71387)	2.41	2.42E-04	4.02E-02
Metabolism (R-HSA-1430728)	2.21	5.84E-17	1.07E-13



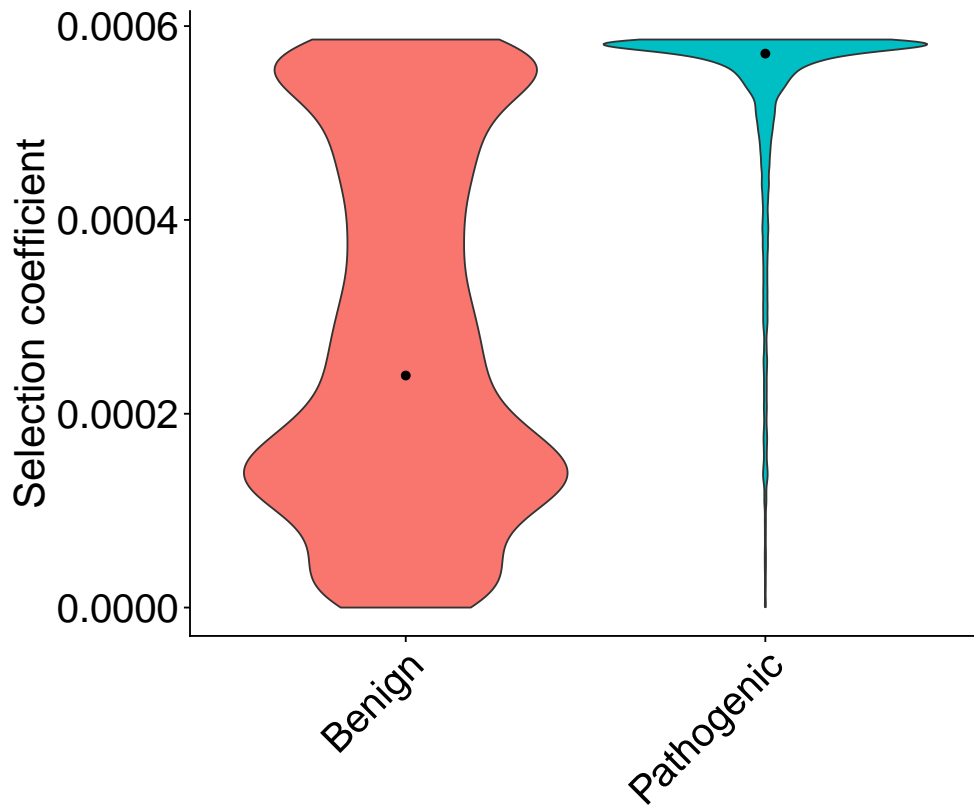
Supplemental Figure S1: Comparison of folded site-frequency spectra between the Poisson Random Field model and the observed data. **(A)** The expected site-frequency spectrum from the demographic model provides an excellent fit to the observed site-frequency spectrum in putative neutral regions. **(B)** The expected site-frequency spectrum from the three-component selection model provides an excellent fit to the observed site-frequency spectrum in coding regions.



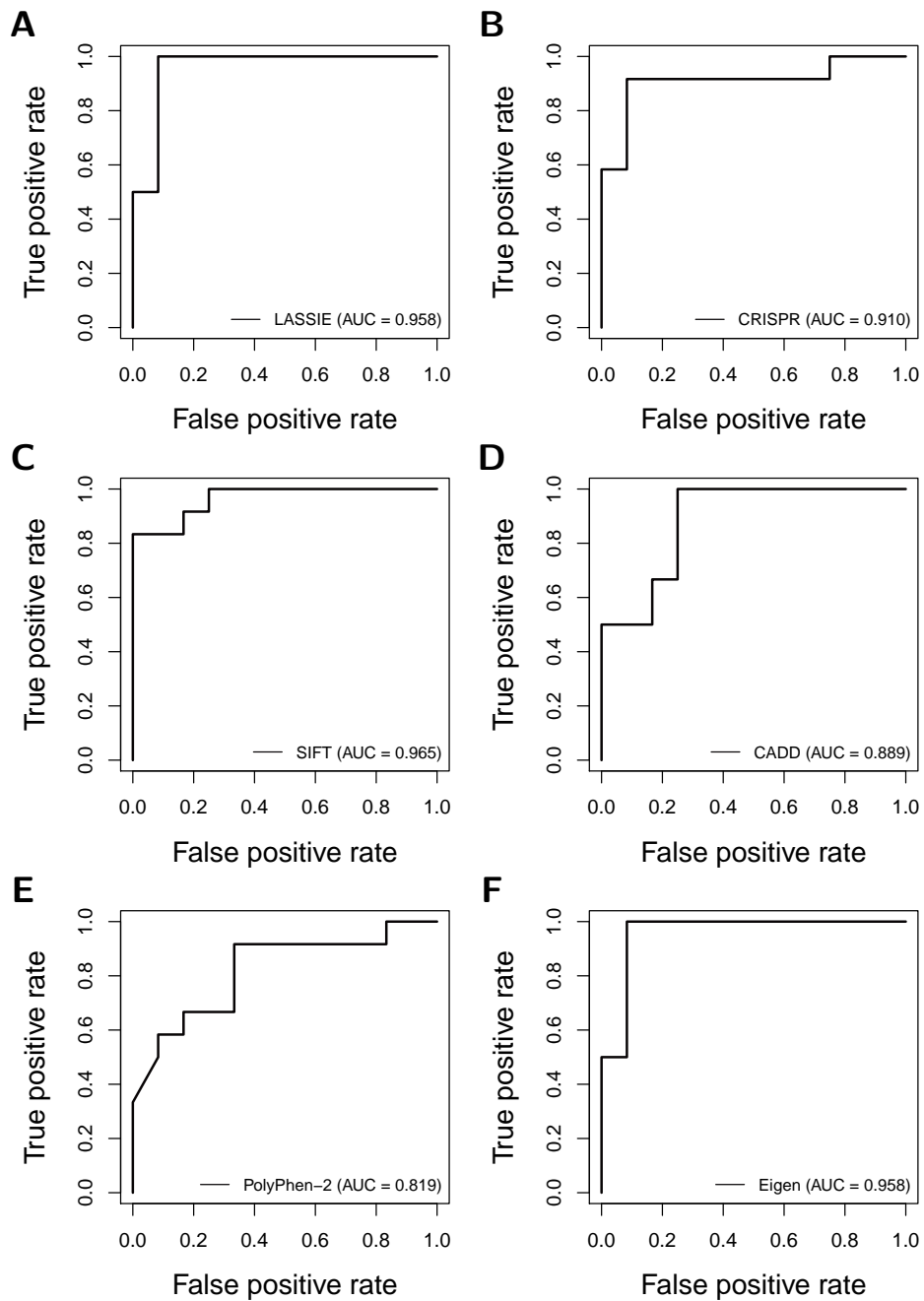
Supplemental Figure S2: Feature weights estimated by LASSIE. Blue and red bars depict the weights associated with strong and weak selection, respectively. A positive weight suggests that the corresponding feature is positively correlated with weak or strong selection. The colors of feature names correspond to four feature groups: variant category (green), sequence conservation (purple), structural information (red), and regulatory information (orange).



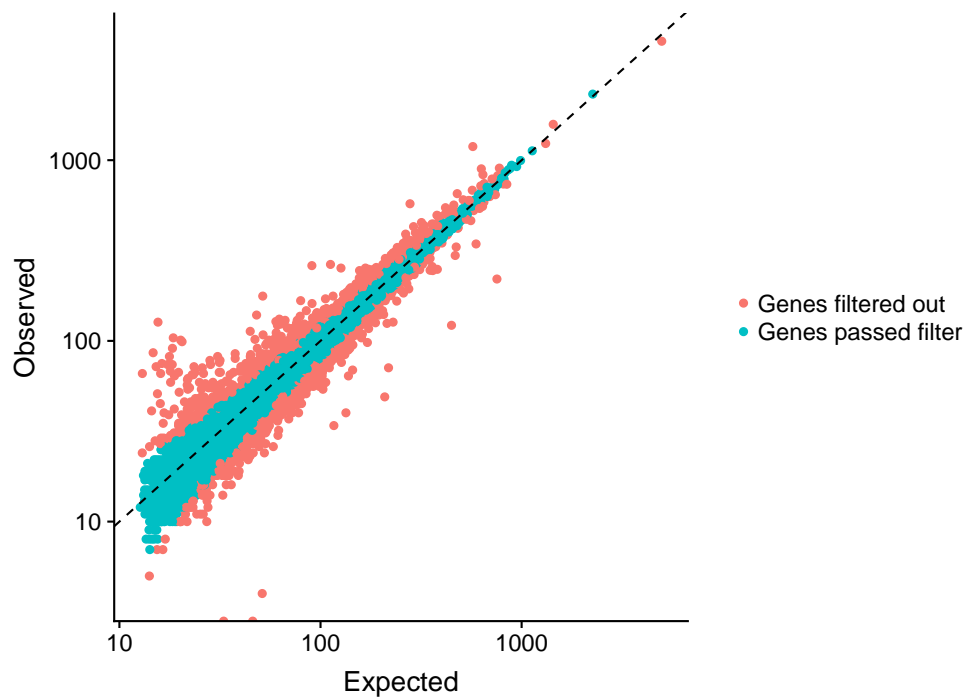
Supplemental Figure S3: Distribution of selection coefficients across synonymous mutations. **(A)** Negative selection on synonymous mutations is positively correlated with gene expression level (Spearman's rank correlation coefficient $\rho = 0.286$; two-tailed $p < 10^{-15}$ by t -test). **(B)** Negative selection on synonymous mutations is stronger in multi-exon genes than single-exon genes (two-tailed $p < 10^{-15}$ by Wilcoxon rank-sum test). **(C)** Negative selection on synonymous mutations is stronger in SRSF1 binding sites than non-binding sites (two-tailed $p < 10^{-15}$ by Wilcoxon rank-sum test). **(D)** Negative selection on synonymous mutations is strong in SRSF7 binding sites than non-binding sites (two-tailed $p = 3.036 \times 10^{-10}$ by Wilcoxon rank-sum test). The SRSF1 and SRSF7 binding sites were obtained from a previous study (Van Nostrand et al., 2016). In each box plot, the bottom, the top, and the internal horizontal bar of each box depict the first quartile, the third quartile, and the median, respectively. The whiskers represent the 1.5-fold interquartile ranges.



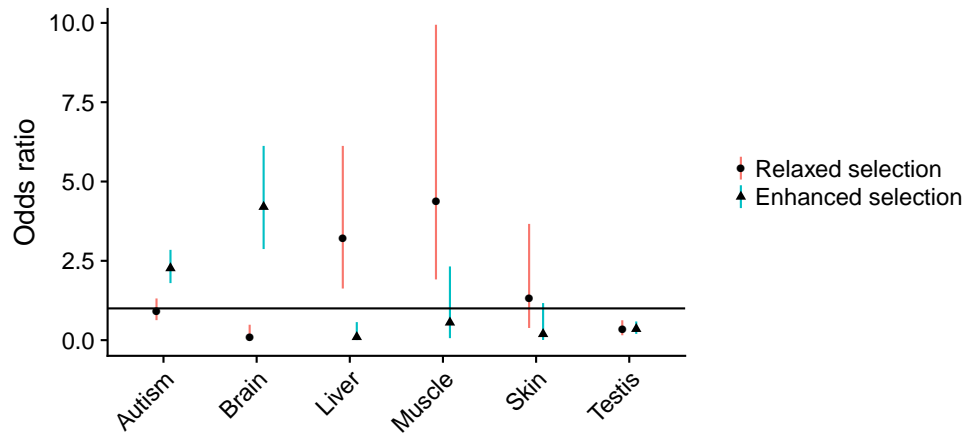
Supplemental Figure S4: Comparison of the distributions of selection coefficients between pathogenic and benign variants from the ClinVar database.



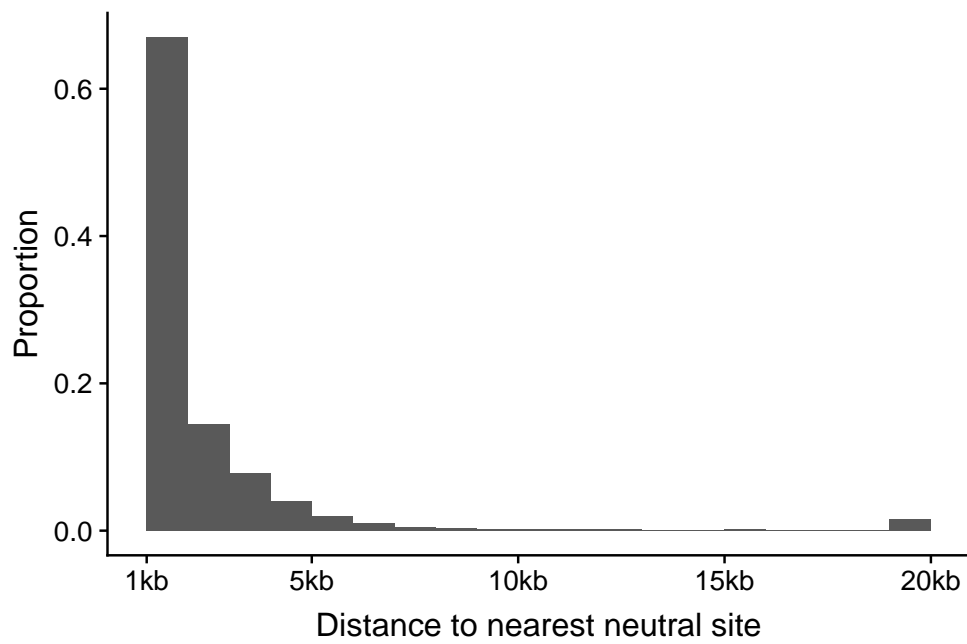
Supplemental Figure S5: Prediction power of different computational and experimental methods for separating pathogenic variants from benign variants in the *BRCA1* gene. The CRISPR scores were obtained from Findlay et al. (2018).



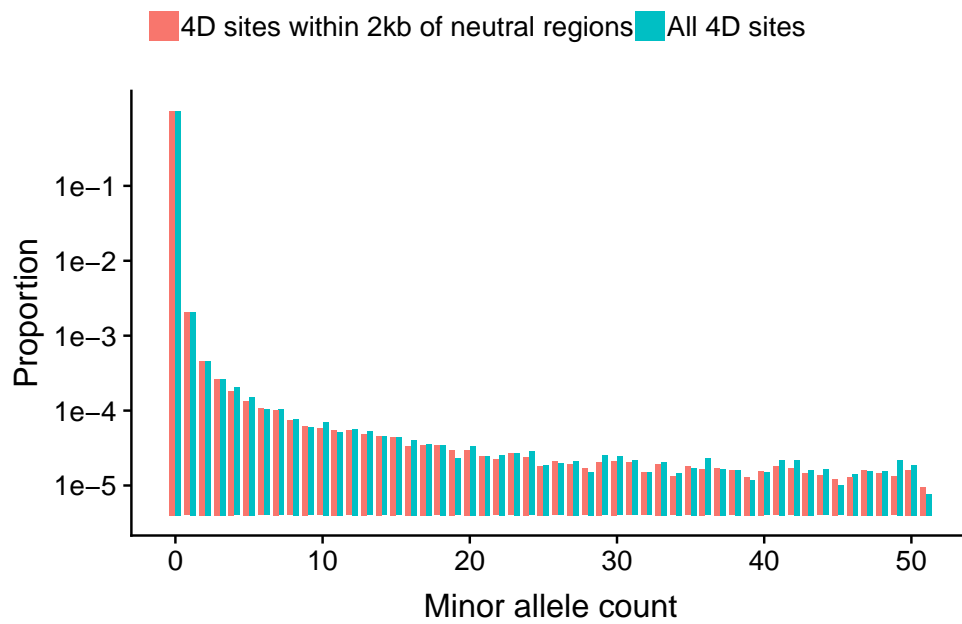
Supplemental Figure S6: Comparison of expected and observed numbers of synonymous mutations across all protein coding genes in the ExAC data set. Each dot represents a single protein coding gene. A gene was filtered out if it is enriched or depleted with synonymous mutations in the ExAC data set (Poisson-Binomial test; FDR rate ≤ 0.2). Genes with less than 200 potential synonymous mutations were also filtered out and were not shown in the plot. The dashed diagonal line represents the case of the expected number of mutations being equal to the observed number of mutations.



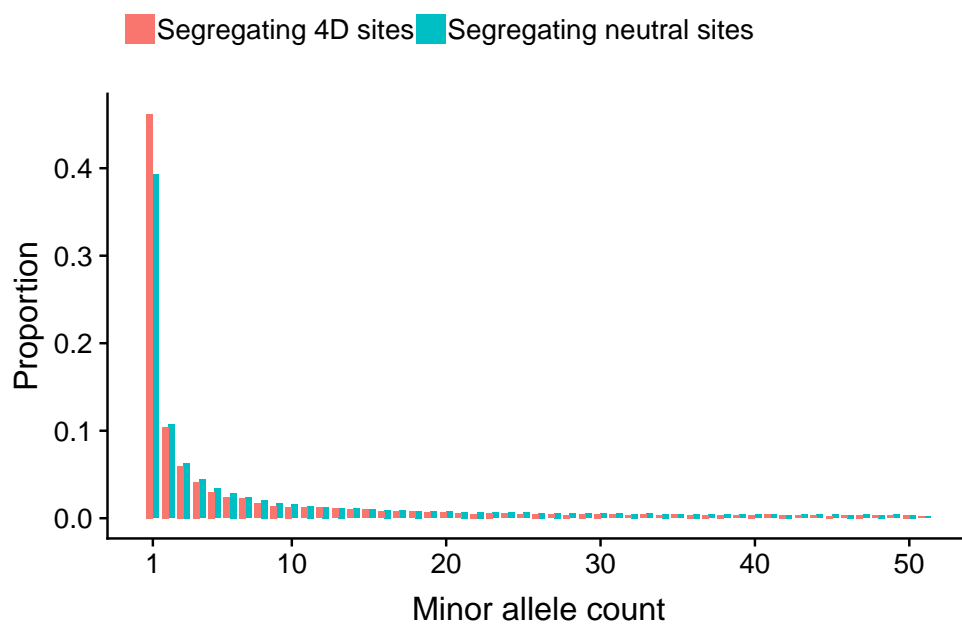
Supplemental Figure S7: Groups of genes enriched for enhanced or relaxed selection. To control gene length as a possible confounding factor, foreground genes are matched with background genes by the number of potential mutations.



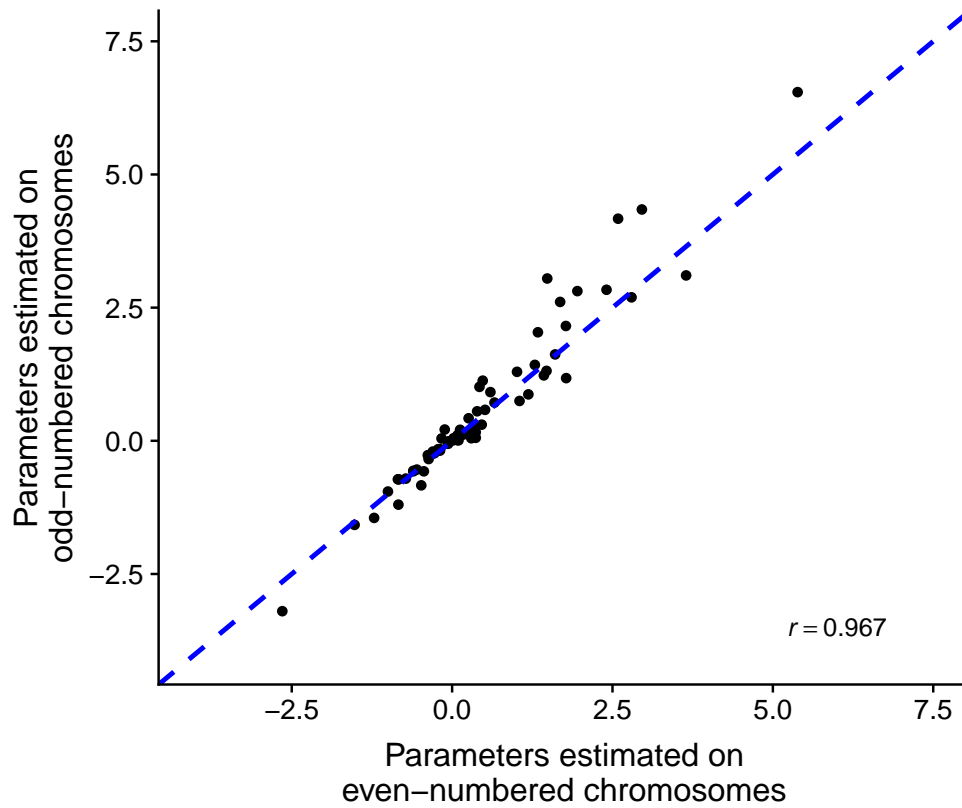
Supplemental Figure S8: Distribution of distances between coding sites and nearest putative neutral site. 1.5% coding sites are more than 20kb away from any putative neutral site and their distances are truncated at 20kb.



Supplemental Figure S9: Comparison of folded site-frequency spectra between all 4D synonymous sites and 4D sites within 2kb of any neutral sites.



Supplemental Figure S10: Comparison of folded site-frequency spectra between segregating 4D sites and neutral sites. In comparison with segregating neutral sites, segregating 4D sites are enriched with singleton variants, suggesting negative selection on synonymous mutations.



Supplemental Figure S11: Correlation between LASSIE parameters separately estimated on even- and odd-numbered chromosomes (Pearson's correlation coefficient $r = 0.967$). The dashed diagonal line represents the case of the estimated parameters being equal between even- and odd-numbered chromosomes.