

# Fast, scalable prediction of deleterious noncoding variants from functional and population genomic data

Yi-Fei Huang<sup>1</sup>, Brad Gulko<sup>1,2</sup> & Adam Siepel<sup>1</sup>

**Many genetic variants that influence phenotypes of interest are located outside of protein-coding genes, yet existing methods for identifying such variants have poor predictive power. Here we introduce a new computational method, called LINSIGHT, that substantially improves the prediction of noncoding nucleotide sites at which mutations are likely to have deleterious fitness consequences, and which, therefore, are likely to be phenotypically important. LINSIGHT combines a generalized linear model for functional genomic data with a probabilistic model of molecular evolution. The method is fast and highly scalable, enabling it to exploit the ‘big data’ available in modern genomics. We show that LINSIGHT outperforms the best available methods in identifying human noncoding variants associated with inherited diseases. In addition, we apply LINSIGHT to an atlas of human enhancers and show that the fitness consequences at enhancers depend on cell type, tissue specificity, and constraints at associated promoters.**

In the human genome, most nucleotides that are associated with diseases or other phenotypes, or that show signatures of natural selection, fall outside of protein-coding genes<sup>1–3</sup>. Many of these nucleotides appear to fall in *cis*-regulatory elements, including promoters, enhancers, and insulators. Similar observations hold across most animals and plants<sup>4–7</sup>. Recent efforts to characterize noncoding sequences using high-throughput biochemical assays have produced a wealth of data, identified many regulatory elements, and clarified general aspects of gene regulation<sup>8–12</sup>. Nevertheless, a substantial gap remains between the outcomes of these experiments and a detailed understanding of noncoding function, for several reasons. First, these assays generally measure genomic and epigenomic features that are roughly correlated with, but not directly indicative of, regulatory function. Second, they generally have relatively low resolution along the genome, identifying regions that are hundreds of nucleotides long rather than pinpointing single nucleotides. Third, these measures are highly condition specific, and data have only been generated for a small subset of cell types and conditions.

As a consequence, there is a pressing need for computational methods that more precisely predict regulatory function by jointly considering the results of numerous such assays together with

complementary data, such as annotations of protein-coding genes and measures of evolutionary conservation across species. The development of statistical and machine-learning methods that attempt to address this integrative prediction challenge has emerged as an active, fast-moving area of research. Recently published methods in this area can be roughly divided into three categories: (i) machine-learning classifiers that attempt to separate known disease variants from putatively benign variants using a variety of genomic features (for example, GWAVA<sup>13</sup> and FATHMM-MKL<sup>14</sup>); (ii) sequence- and motif-based predictors for the impact of noncoding variants on cell-type-specific molecular phenotypes, such as chromatin accessibility or histone modifications (for example, DeepBind<sup>15</sup>, DeepSEA<sup>16</sup> and Basset<sup>17</sup>); and (iii) evolutionary methods that consider data on genetic variation together with functional genomic data with the aim of predicting the effects of noncoding variants on fitness (for example, CADD<sup>18</sup>, DANN<sup>19</sup>, FunSeq2 (ref. 20), and fitCons<sup>3</sup>). A limitation of the methods in the first category is that they depend strongly on the available training data, which may be limited and may not be representative of the broader class of regulatory sequences of interest. Methods in the second category have the limitation that the importance of the molecular phenotypes at the organismal level is often unclear. Evolutionary methods, by contrast, obtain their signal not primarily from previously assigned class labels but instead from the signatures of natural selection over many generations. They are, therefore, both less data limited and more focused on the phenotypes that truly influence fitness than the other methods. This approach is likely to be particularly powerful for detecting regulatory variants that tend to be under strong purifying selection, such as rare variants associated with severe diseases. Evolution-based methods also naturally integrate over cell types, an important strength when the relevant tissue or cell types for a condition of interest are unknown.

Among the available evolution-based methods, fitCons (ref. 3) is unique in being able to explicitly characterize the influence of natural selection at each genomic site of interest using a full probabilistic evolutionary model and patterns of genetic variation within and between species. FitCons makes a distinction between functional genomic data and comparative genomic data by first defining several hundred clusters of genomic positions with distinct functional genomic ‘fingerprints’ and then estimating the fraction of nucleotides under natural selection within each cluster using polymorphism and divergence

<sup>1</sup>Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, USA. <sup>2</sup>Graduate Field of Computer Science, Cornell University, Ithaca, New York, USA. Correspondence should be addressed to A.S. (asiepel@cshl.edu).

Received 15 August 2016; accepted 13 February 2017; published online 13 March 2017; doi:10.1038/ng.3810

data. These estimates are obtained using the INSIGHT evolutionary model<sup>21,22</sup> and are interpreted as the probabilities that mutations in each cluster of genomic sites will have fitness consequences (fitCons scores). In this manner, fitCons aggregates information about natural selection from large numbers of sites with similar functional profiles based on evolutionary first principles. A major limitation of the method, however, is that it scales poorly with the available functional genomic data. In particular, the number of clusters considered by the method increases exponentially with the number of functional genomic annotations, which keeps it from taking advantage of the growing body of functional genomic data. A related problem is that the restriction to small numbers of genomic features leads to a relatively coarse-grained, blocky pattern of scores along the genome, which does not allow for fine distinctions among nearby nucleotide sites.

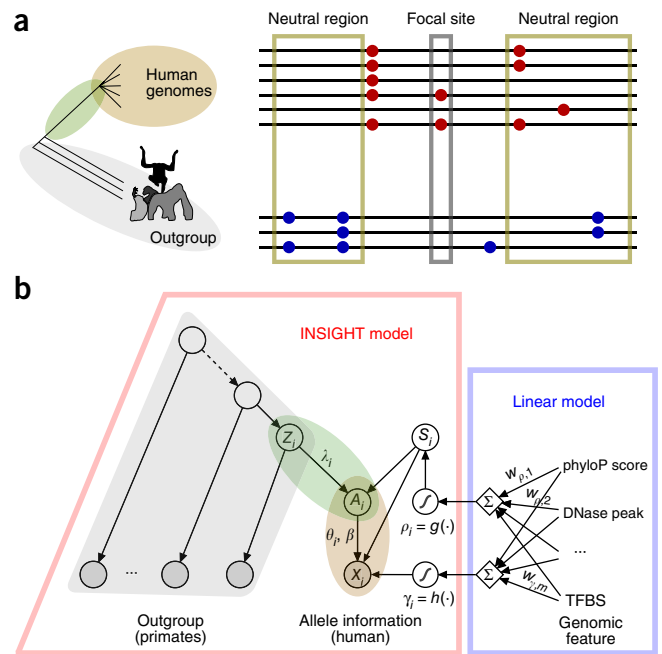
Here we describe a new method, linear INSIGHT (LINSIGHT), that is based on the existing INSIGHT–fitCons framework but that has vastly improved speed, scalability, genomic resolution, and prediction power. The main idea behind LINSIGHT is to bypass the clustering step of fitCons and, instead, couple the probabilistic INSIGHT model directly to a generalized linear model for genomic features. This strategy results in a more streamlined model that scales linearly, rather than exponentially, with the available data and can make direct use of the input data, with no need for discretization. By integrating a large number of genomic features, LINSIGHT provides a precise, high-resolution description of the fitness consequences of noncoding mutations in the human genome. We demonstrate that LINSIGHT outperforms state-of-the-art prediction methods in the task of prioritizing noncoding disease variants from the Human Gene Mutation database (HGMD)<sup>23</sup> and the National Center for Biotechnology Information (NCBI) ClinVar database<sup>24</sup>. Furthermore, we use LINSIGHT to show that the evolutionary constraints on human enhancers depend on their associated tissue types, degree of tissue specificity, and associated promoters, which has important implications for understanding the evolution of *cis*-regulatory elements and for improving variant prioritization methods. Our LINSIGHT scores are available as a track on the Cold Spring Harbor Laboratory mirror of the UCSC Genome Browser (hg19 assembly; <http://genome-mirror.cshl.edu/>). The LINSIGHT software is available through the GitHub repository (<https://github.com/CshlSiepelLab/LINSIGHT>).

**RESULTS**

**LINSIGHT combines INSIGHT with a scalable linear model**

The original INSIGHT and fitCons methods<sup>3,21,22</sup> infer the selective pressure on noncoding sites, and hence the likely fitness consequences of noncoding mutations, by contrasting patterns of genetic variation at each focal site with the patterns at nearby genomic regions that are likely to be free from the influence of selection (‘neutrally evolving sites’). To address the problem that genetic variation within species and between closely related species (such as humans and chimpanzees) are sparse across the genome, fitCons pools information across the thousands of genomic sites that are assigned to each discrete cluster.

The key idea behind LINSIGHT is, instead, to accomplish this pooling of information across sites indirectly by using a generalized linear model (Fig. 1 and Table 1; see Supplementary Note and Supplementary Tables 1 and 2 for complete details). In particular, the parameters of the INSIGHT model that describe natural selection ( $\rho$  and  $\gamma$ ) are determined as linear-sigmoid functions of the genomic features that are local to each site (the third selection parameter from INSIGHT,  $\eta$ , is omitted because positive selection has a negligible



**Figure 1** Conceptual overview of LINSIGHT. (a) Similar to the fitCons method<sup>3</sup>, LINSIGHT estimates the probabilities that mutations at each genomic site will have fitness consequences, based on patterns of genetic polymorphism within a species (here, humans) and patterns of divergence from closely related outgroup species (chimpanzee, orangutan, and rhesus macaque). Patterns of genetic variation at the focal site and at other sites like it are contrasted with those in neutrally evolving regions nearby. Red circles indicate human SNPs, and blue circles indicate nucleotide substitutions between species. (b) LINSIGHT combines the probabilistic graphical model from INSIGHT<sup>21,22</sup> with a generalized linear model. The selection parameters from INSIGHT,  $\rho$  and  $\gamma$ , are defined in a site-wise manner by linear combinations of local genomic features, followed by sigmoid transformations. The figure summarizes the behavior at a particular focal site  $i$ . The shaded regions in gray, green, and tan indicate corresponding portions of the phylogeny and sequence data (a) and the INSIGHT model (b). See Table 1 for definitions of all parameters and variables.

effect in this setting; see Supplementary Note). Thus, the probability of fitness consequences for mutations at each site  $i$ , denoted  $\rho_i$ , is assumed to depend on the genomic features at that site—such as its RNA expression level (RNA-seq read depth), chromatin accessibility (DNase-I hypersensitive sites), and histone modifications or bound transcription factors (ChIP-seq peaks)—as well as on features based on annotations (for example, distance to the nearest transcription start site (TSS) or a match to a known transcription factor binding sites (TFBS) motif) and comparative genomics (for example, phyloP<sup>25</sup> or phastCons<sup>4</sup> scores). We refer to  $\rho_i$  as the LINSIGHT score at site  $i$ . This scoring strategy has several major advantages—it requires no clustering and no discretization, and it scales linearly with the available genomic features, allowing hundreds of features to be considered. In contrast to fitCons, the scalability of LINSIGHT enables data to be pooled across cell types, and it allows the scores to reach single-nucleotide resolution along the genome. Nevertheless, LINSIGHT continues to benefit from the advantages of the probabilistic INSIGHT model of molecular evolution.

All parameters of the LINSIGHT model are estimated simultaneously from genome-wide data by maximum likelihood using an online stochastic gradient-descent algorithm (Online Methods). The gradients for the feature weights are efficiently computed by

**Table 1 Summary of key model parameters and variables**

Parameters inherited from INSIGHT <sup>a</sup>	
$\rho_i$	Probability that site $i$ is under selection. Interpreted as the LINSIGHT score for site $i$
$\gamma_i$	Expected relative rate of low-frequency-derived alleles at site $i$ given that it is under selection
$\lambda_i$	Neutral substitution rate at site $i$
$\theta_i$	Neutral polymorphism rate at site $i$
$\beta = (\beta_1, \beta_2, \beta_3)$	Fractions of neutral polymorphisms with low-, intermediate-, and high-frequency-derived alleles
Variables inherited from INSIGHT <sup>a</sup>	
$X_i = (X_i^{\text{maj}}, X_i^{\text{min}}, Y_i)$	Observed polymorphism data at site $i$ , including major allele, minor allele, and minor-allele-frequency class
$Z_i$	Human–chimpanzee ancestral allele at site $i$
$A_i$	Human ancestral allele at site $i$
$S_i$	Indicator for whether or not site $i$ is under selection
Components of LINSIGHT's generalized linear model <sup>a</sup>	
$D_i = (d_{i,1}, \dots, d_{i,m})$	Genomic feature vector at site $i$
$W_\rho = (w_{\rho,1}, \dots, w_{\rho,m})$	Weight vector for $\rho$ (free parameters)
$W_\gamma = (w_{\gamma,1}, \dots, w_{\gamma,m})$	Weight vector for $\gamma$ (free parameters)
$g()$	Sigmoid function for $\rho$ (Gompertz)
$h()$	Sigmoid function for $\gamma$ (logistic)

<sup>a</sup>See **Supplementary Note** and **Supplementary Table 1** for full details.

the back-propagation method that is widely used in neural network training<sup>26</sup>. Indeed, the model can be considered a type of neural network, although one without hidden layers. Its main disadvantage relative to fitCons—the assumption of an additive, linear relationship between features and selection parameters—could be addressed by adding hidden layers to the neural network, although we have found its performance to be excellent without this extension. Notably, the amount of data available for training is large in comparison to the number of free parameters, and we have not yet found regularization to be necessary, but it could easily be added if necessary.

**LINSIGHT scores across the human genome are generally consistent with, but often improve on, previous measures of evolutionary conservation**

We applied LINSIGHT to a large public data set—consisting of complete genomic sequences for multiple human individuals and nonhuman primates, comparative genomic data for mammals and vertebrates, and a wide variety of functional genomic data—and we generated LINSIGHT scores for all of the positions across the human reference genome. We considered a total of 48 genomic features, which belonged to three general classes: conservation scores, predicted binding sites, and regional annotations (**Table 2** and **Supplementary Table 3**).

The distributions of LINSIGHT scores in annotated regions of the noncoding genome are generally consistent with previous observations based on conservation scores<sup>1,4,25</sup>. For example, splice sites were very highly constrained (median LINSIGHT score of 0.956, indicating a 95.6% probability of fitness consequences due to mutations at these nucleotide sites), whereas annotated TFBSs showed reduced, but still substantial, constraint (median LINSIGHT score of 0.240 for TFBSs shared across species, median LINSIGHT score of 0.106 for all TFBSs from the Ensembl Regulatory Build<sup>27</sup>) (**Fig. 2a**). Other promoter regions (median LINSIGHT score of 0.073) and untranslated regions (UTRs; median LINSIGHT scores of 0.128 and 0.076 for 5' and 3' UTRs, respectively) were somewhat less constrained, and unannotated intronic and intergenic regions showed the least amount of constraint (median LINSIGHT scores of 0.044–0.048). As observed previously, 5' UTRs showed somewhat more constraint than 3' UTRs, although both types of UTRs contained subsets of sites that were subject to strong selection (LINSIGHT score > 0.8)<sup>4,25</sup>. The estimate for the more conserved TFBSs (LINSIGHT score = 0.240) was similar to, but slightly lower than, previous estimates that were obtained directly

from experimentally defined TFBSs (~30–40% of sites under selection<sup>22,28</sup>), despite the fact that it was obtained indirectly in this case via the generalized linear model. The genome-wide average of the LINSIGHT scores was ~0.07, suggesting that ~7% of noncoding sites are under evolutionary constraint, which is consistent with numerous previous studies<sup>3,4,29–31</sup>.

Across all noncoding positions in the genome, the LINSIGHT scores were fairly well correlated with those from other recently published methods, particularly within conserved elements, which are enriched for regulatory function (**Supplementary Fig. 1** and **Supplementary Note**).

On the task of identifying likely regulatory elements, the methods that make use of functional genomic data generally perform better than pure conservation methods, and LINSIGHT was among the best at this task (**Supplementary Note**). For example, LINSIGHT had good power to identify transcription factor binding sites from the ORegAnno database<sup>32</sup> (AUC = 0.926), outperformed only by the DeepSEA functional significance score (AUC = 0.965) and FunSeq2 (AUC = 0.950) (**Supplementary Fig. 2**). Thus, despite that it relies on an evolutionary objective function, LINSIGHT maintains good performance in the prediction of regulatory elements.

Consistent with these general trends, LINSIGHT highlighted many of the regions that were identified by conservation methods such as phastCons<sup>4</sup>, phyloP<sup>25</sup>, and GERP++<sup>33</sup>, but it also identified some regions that had relatively low conservation scores yet are likely to have important biological functions. An example is HGMD variant CR065653 in a putative enhancer, which is associated with upregulation of the telomerase reverse transcriptase (*TERT*) gene and which had an increased LINSIGHT score but was not identified by phastCons, phyloP, or GERP++ as being under constraint (**Fig. 2b**). This example also demonstrates that the genomic resolution of the LINSIGHT scores is dramatically better than that of fitCons and approaches the nucleotide resolution of phyloP and GERP++. In addition to enhancers, LINSIGHT identified functional variants in promoter regions (**Supplementary Fig. 3a**) and associated with splicing (**Supplementary Fig. 3b**). Thus, it is useful as a general predictor of functional noncoding sites under evolutionary constraint.

**LINSIGHT accurately identifies disease-associated variants in noncoding regions**

We tested the ability of LINSIGHT to identify noncoding nucleotide positions that are associated with inherited human diseases, using the

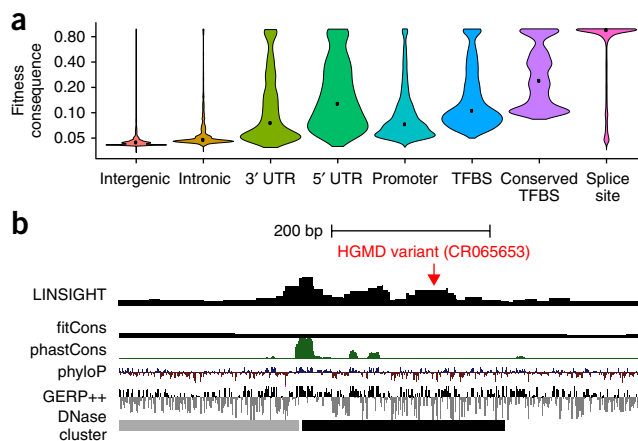
**Table 2 Summary of genomic features used for LINSIGHT scores**

Class	Genomic feature <sup>a</sup>	Spatial resolution
Conservation	phyloP score	High
	phastCons element	High
	SiPhy element	High
	CEGA element	High
Binding site	Conserved TFBS	High
	rVISTA TFBS	High
	SwissRegulon TFBS	High
	Predicted TFBS within ChIP-seq peak	High
	Conserved miRNA binding site	High
Regional annotation	Splicing site predicted by SPIDEX	High
	ChIP-seq peak of transcription factor	Low
Regional annotation	DNase-I hypersensitive site	Low
	UCSC FAIRE peak	Low
	RNA-seq signal	Low
	Histone modification peak	Low
	FANTOM5 enhancer	Low
	Predicted distal regulatory module	Low
	Distance to nearest TSS	Low

<sup>a</sup>Each 'genomic feature' listed here may actually correspond to multiple features in the model. For example, four features are derived from phyloP scores: two from the mammalian phyloP scores and two from the vertebrate phyloP scores. See **Supplementary Table 3** for complete details.

HGMD<sup>23</sup> and ClinVar<sup>24</sup> databases to define positive examples, and common polymorphisms (minor allele frequency (MAF) > 1%), which are unlikely to be functionally important, to define negative examples. For comparison, we evaluated the CADD<sup>18</sup>, Eigen<sup>34</sup>, DeepSEA<sup>16</sup>, FunSeq2 (ref. 20), GWAVA<sup>13</sup>, and phyloP<sup>25</sup> methods on the same task. For each scoring method, we computed false-positive versus true-positive rates for the complete range of score thresholds, displayed the results as 'receiver operating characteristic' (ROC) curves, and measured prediction power by the area-under-the-curve (AUC) statistic. Because the results of these tests can be highly sensitive to the criteria for selecting negative examples, we considered three schemes of increasing stringency<sup>13</sup>: a random sample of negative examples (unmatched), negative examples matched by distance to the nearest TSS (matched TSS), and negative examples matched by specific genomic region (matched region; Online Methods). In all cases, equal numbers of positive and negative examples were considered.

Overall, LINSIGHT outperformed all of the other methods in all comparisons (**Fig. 3**). Its absolute prediction power varied across matching schemes in a predictable manner, being highest in the unmatched comparison (for example, AUC = 0.897 for HGMD) and decreasing in the matched TSS (AUC = 0.759) and matched region (AUC = 0.661) comparisons. The same effect also occurred for most of the other methods, but the methods that make more use of regional information (such as FunSeq2) suffered more as the matching stringency increased. These observations highlight the difficulty of distinguishing functional sites from nearby nonfunctional sites, which is considerably harder than separating regions enriched in functional sites from the genomic background. Nevertheless, LINSIGHT has some power for this challenging task. In almost all cases, the AUCs were considerably higher for ClinVar than for HGMD, apparently because ClinVar is heavily enriched for variants in splice sites, which are relatively easy to identify (**Supplementary Fig. 4**). An exception to this rule was GWAVA, which performed exceptionally well on HGMD (cross-validation AUCs of 0.71–0.97)<sup>13</sup> and much more poorly on ClinVar (AUCs of 0.734–0.884); however, GWAVA was trained using HGMD<sup>13</sup>, and its performance on that data set appears to reflect overfitting (it is not shown in the HGMD ROC plots for this reason).

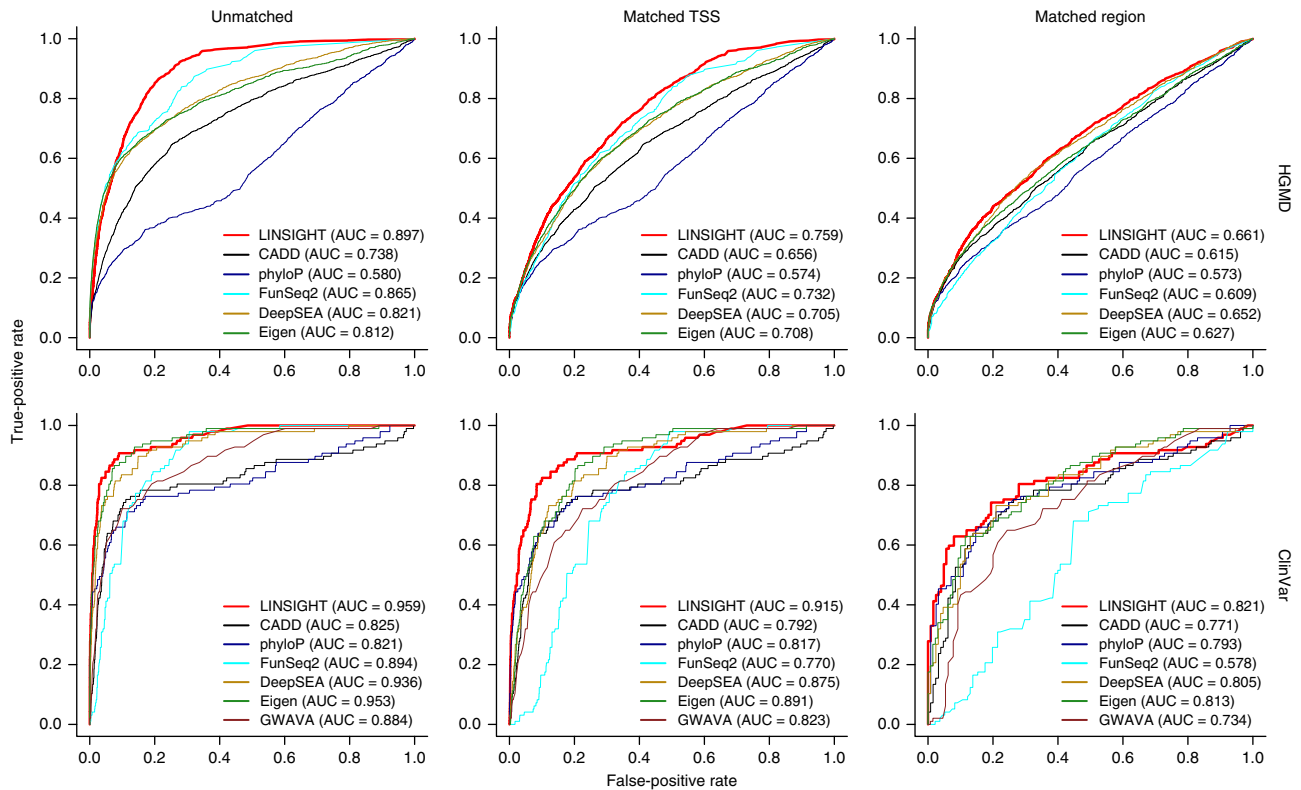


**Figure 2** Summary of LINSIGHT scores across the noncoding human genome (3.001 billion nucleotide sites). **(a)** Distributions of LINSIGHT scores for various genomic regions. Intergenic regions, intronic regions, UTRs, and 1-kb promoters were defined using GENCODE annotations (version 19); TFBSs were predicted from ChIP-seq peaks (Ensembl Regulatory Build); conserved TFBSs were obtained from the UCSC Genome Browser. Within each violin plot, the width represents density, and the black dot represents the median LINSIGHT score. Note the logarithmic vertical scale. **(b)** UCSC Genome Browser display showing LINSIGHT scores alongside those from fitCons, phastCons, phyloP, and GERP++. LINSIGHT integrates functional genomic data together with conservation scores and other features to provide a high-powered, high-resolution measure of potential function. In this example, it is the only method to highlight a variant from HGMD (CR065653) that is associated with upregulation of the *TERT* gene. See **Supplementary Figure 3** for additional examples.

This dependency on the training set for GWAVA demonstrates one of the pitfalls of pure classification strategies and highlights a strength of the evolution-based strategy, which does not require a training set. Nevertheless, phyloP performed quite poorly on the HGMD data set (**Fig. 3**), showing that scores based exclusively on evolution are of limited usefulness in this task.

The performance advantage of LINSIGHT was maintained when performance was measured using precision-recall curves in place of standard ROC curves (**Supplementary Fig. 5**) and when rare variants were used in place of common variants as negative examples (**Supplementary Figs. 6 and 7**). These performance advantages were statistically significant in most cases, with a few exceptions that mostly stemmed from the small size of the ClinVar data set (**Supplementary Tables 4 and 5**). In addition, a more detailed comparison with CADD showed that training CADD's logistic regression model using LINSIGHT's features resulted in improved performance but not enough to make it competitive with LINSIGHT (**Supplementary Table 6**). Thus, the excellent performance of LINSIGHT in these tests seems to derive both from its use of a broad collection of informative features along the genome and its probabilistic model of evolution.

To gain insight into which genomic features were most informative, we systematically omitted groups of related features and reassessed the prediction performance of LINSIGHT (**Supplementary Note**). Briefly, we found that regional features, such as ChIP-seq peaks and DNase-I-hypersensitive sites (**Table 2**), were broadly useful in distinguishing genomic regions enriched for functional variants from the genomic background, but conservation scores were most important in separating functional sites from nearby nonfunctional sites (**Supplementary Fig. 8**). Predicted binding sites were most informative in promoter regions.



**Figure 3** Prediction power of various computational methods for distinguishing disease-associated noncoding variants from variants not likely to have phenotypic effects. True-positive and false-positive rates are the proportions of disease-associated and neutral variants, respectively, having scores that exceeded each threshold, as the threshold is varied. Power was quantified using the area-under-the-curve (AUC) statistic. Results are shown for positive examples from the HGMD<sup>23</sup> (1,495 variants) (top) and ClinVar<sup>24</sup> (101 variants not in HGMD) (bottom) databases. Only autosomal variants were included, and duplicated variants were removed. Common SNPs (MAF > 1%) were used as negative examples and were either randomly selected (unmatched) (left), matched to positive examples by distance to nearest TSS (matched TSS) (middle), or matched to positive examples within 1 kb along the genome (matched region) (right). The numbers of positive and negative examples were balanced by subsampling, which was performed 100 times to obtain the average true-positive and false-positive rates. Analysis with LINSIGHT was compared with that of CADD<sup>18</sup>, phyloP<sup>25</sup>, FunSeq2 (ref. 20), DeepSEA<sup>16</sup>, Eigen<sup>34</sup>, and GWAVA<sup>13</sup>. FitCons was not included because it performs poorly on this task due to its low genomic resolution and cell-type specificity. GWAVA results are not shown for the HGMD data set because GWAVA was trained on this data set.

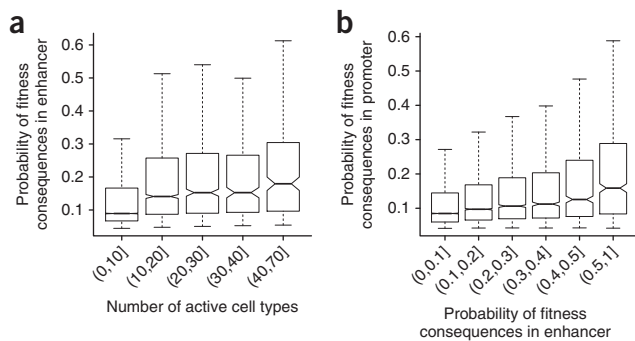
**The evolutionary constraints on enhancers are context dependent**

LINSIGHT is also potentially useful for studying the influence of natural selection on noncoding sequences. As compared with other measures of selection, LINSIGHT has the advantages of considering both functional genomic and population genomic data, of detecting the influence of selection on relatively recent time scales (that is, since the human–chimpanzee divergence), and of providing a model-based, easily interpretable measure of fitness consequences. With these advantages in mind, we used LINSIGHT to gauge the degree of evolutionary constraint on enhancers in the human genome, considering in particular the relationships of constraint with the number and type of active cell types, and with constraint at the target promoter of each enhancer. We analyzed nearly 30,000 enhancers (median length 293 bp) from a recently published atlas of active enhancers in dozens of human cell types and tissues, which were identified on the basis of their transcriptional signatures<sup>35</sup>. This approach of annotating enhancers based on enhancer-associated RNAs (eRNAs) has been shown to identify elements that have active roles in gene regulation in a cell-type-specific fashion with high genomic resolution<sup>35–37</sup>.

First, we examined the relationship between the LINSIGHT scores and the number of cell types in which each enhancer is active. We found that the LINSIGHT scores were significantly positively correlated with the number of active cell types (Spearman’s  $\rho = 0.284$ ,

$P < 10^{-15}$ ; **Fig. 4a**), indicating that a broader spectrum of activity across cell types is associated with stronger purifying selection. To ensure that this observation reflected real differences in selective pressure and not simply correlations with the epigenomic features considered by LINSIGHT, we retrained LINSIGHT using only conservation scores and predicted binding sites and obtained essentially identical results (**Supplementary Fig. 9a**). Furthermore, a partial correlation test indicated that the LINSIGHT scores were still strongly correlated with the number of cell types when controlling for eRNA expression levels averaged across all FANTOM5 libraries (partial Spearman’s  $\rho = 0.24$ ;  $P < 10^{-15}$ ). These findings paralleled similar findings for protein-coding genes<sup>38–40</sup> and TFBSs<sup>22</sup> and likely reflect a general correlation between pleiotropy and constraint (see Discussion).

Second, we examined the relationship between the LINSIGHT score and the tissue type in which each enhancer is active, focusing on enhancers active in a single tissue type. We found that tissue-specific enhancers associated with sensory perception (olfactory region and parotid gland), the immune system (lymph node), digestion (stomach), and male reproduction (penis and testis) had the lowest LINSIGHT scores, whereas tissue-specific enhancers associated with tissues such as smooth muscle, the skin, and the urinary tract and bladder had the highest LINSIGHT scores (**Supplementary Fig. 10**). These findings were also broadly consistent with findings for



**Figure 4** Evolutionary constraints on enhancers. **(a)** Probability of fitness consequences for mutations in enhancers (measured by average LINSIGHT score) is positively correlated with the number of cell types in which each enhancer is active (Spearman's rank correlation coefficient  $\rho = 0.284$ ; two-tailed  $P < 10^{-15}$  by  $t$ -test). Results are shown for 29,303 enhancers in 69 cell types. Labels on the x axis of the form '(a, b]' represent ranges from a (exclusive) to b (inclusive) of active cell types. **(b)** Probability of fitness consequences for mutations in enhancers is positively correlated with probability of fitness consequences for mutations in associated promoters (Spearman's rank correlation coefficient  $\rho = 0.150$ ; two-tailed  $P < 10^{-15}$  by  $t$ -test). In each box plot, the bottom and top of the box, and the horizontal bar inside it, represent the first quartile, second quartile, and median, respectively. The whiskers represent 1.5-fold interquartile ranges. Results are shown for 25,067 enhancer–promoter pairs.

protein-coding genes, which have indicated that genes involved in the sensory, immune, dietary, and male reproductive systems are associated with relaxation of constraint and/or positive selection<sup>40,41</sup>. Notably, enhancers that are active in tissues associated with female reproduction (for example, the uterus, female gonad, and vagina) seemed to be under substantially more constraint than those active in tissues associated with male reproduction. Finally, we compared the LINSIGHT scores at enhancer–promoter pairs that were predicted from co-expression across tissues<sup>35</sup>. The LINSIGHT scores for these paired enhancers and promoters were weakly but significantly correlated (**Fig. 4b** and **Supplementary Fig. 9b**), indicating that the same types of evolutionary pressures tend to act at both members of each pair. Taken together, these results indicate that the evolutionary constraints on enhancers are dependent on several factors, including their degree of tissue specificity, the particular tissues in which they are active, and the evolutionary constraints associated with their target promoters.

## DISCUSSION

As sequencing costs fall and appreciation for regulatory variation grows, whole-genome sequencing is rapidly supplanting exome sequencing as the primary technique for identifying and characterizing genetic variants that have phenotypic consequences. Hence, there is an increasing need for computational methods that can effectively prioritize noncoding variants based on their likelihood of phenotypic importance. Here we address this problem with a new computational method, called LINSIGHT, that combines the evolutionary model of our previously developed INSIGHT method with a generalized linear model for functional genomic data and genome annotations, resulting in substantially improved scalability, resolution, and power. We have generated LINSIGHT scores across the human genome, making use of a large collection of publicly available population, comparative, and functional genomic data; we found the scores to be consistent with previously available scores in many respects, but to improve on them in other respects. In particular, on the task of identifying human disease-associated variants from the HGMD and ClinVar databases, LINSIGHT offered the best performance

of several methods that we tested, across a range of types of variants and test designs. Notably, LINSIGHT required no training set of known regulatory or disease variants, and, therefore, it is expected to have better generalization properties than 'supervised' machine-learning classifiers.

In conceptual terms, the new LINSIGHT method is closely related to our previous fitCons method<sup>3</sup>, with the primary difference being that LINSIGHT pools data across sites implicitly through the use of its generalized linear model, whereas fitCons pools data by explicitly clustering sites according to discretized functional genomic signatures. In effect, LINSIGHT trades the restrictions of a linearity assumption for the benefits of computational speed, a reduced parameterization, and scalability to very large numbers of genomic features. Notably, the new model design also has a number of important side benefits. First, it avoids the need for discretization of the genomic features. In addition, as the number of features grows larger, the genomic resolution of the scores naturally becomes much finer, approaching the nucleotide-level resolution of conservation scores. Finally, the generalized linear model can readily be extended to a 'deep' neural network through the addition of hidden layers. Although it remains to be seen how much this extension will help in practice, in principle it can capture the types of nonlinearity and interactions between features that have been observed in this setting (for examples, see refs. 3,42).

Our approach to characterizing noncoding variants is based on the premise that natural selection in the past, at individual nucleotide sites, provides useful information about phenotypic importance in the present. This assumption clearly will not hold in all cases. For example, variants that increase the risk for post-reproductive diseases or that influence phenotypes dependent on the modern human environment will not necessarily show signs of historical purifying selection. In addition, traits that are dependent on highly epistatic loci or on the aggregate contributions of large numbers of loci may have difficult-to-detect marginal contributions to fitness at individual nucleotides. Nevertheless, our results indicate that the evolution-based approach is useful for many phenotypes of interest. Furthermore, in comparison to the available high-throughput experimental methods, evolution-based methods have the crucial advantage of measuring the importance of genetic variants in real organisms in their natural environments over many generations.

We used LINSIGHT to examine the influence of negative selection on enhancers and considered the relationships between constraint on enhancers and numbers of active cell types, tissue of activity, and constraint at associated promoters. LINSIGHT is potentially useful for addressing these questions because it should be much more robust to evolutionary turnover than conventional conservation-based methods, and some classes of enhancers are known to turn over more quickly than others<sup>42</sup>. We found that, in general, the trends in constraint at enhancers paralleled those previously reported for protein-coding genes. For example, constraint increased with breadth of activity across cell types and decreased in tissues that were associated with rapid evolution, such as olfactory regions, the immune system, and male reproduction. Constraint also seemed to be correlated at enhancer–promoter pairs. These observations regarding the specific ways in which evolutionary constraints on enhancers depend on genomic context may be useful in improving the prediction power for the fitness consequences of noncoding mutations.

As has been suggested for protein-coding genes<sup>38</sup>, it seems plausible that the positive correlation between the strength of constraint and the number of active cell types can be explained by pleiotropy: enhancers that are active in more cell types are more likely to participate in multiple regulatory networks, perhaps with distinct roles

involving the binding of different factors and/or the use of different binding sites within each enhancer. As a result, they may be subject to greater constraint. Nevertheless, many open questions remain about the influences of constraint on enhancers, and it will be important to examine these questions further in light of rapidly improving enhancer annotations, data describing enhancer–promoter interactions<sup>43–45</sup>, and observations of complex evolutionary behavior at enhancers<sup>46</sup>.

## METHODS

Methods, including statements of data availability and any associated accession codes and references, are available in the [online version of the paper](#).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

We thank I. Gronau for comments on the manuscript and members of the Siepel laboratory for helpful discussions. This research was supported by the US National Institutes of Health (NIH) grants GM102192 (A.S.) and HG008901 (A.S.). The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH.

## AUTHOR CONTRIBUTIONS

Y.-F.H. and A.S. conceived and designed the study; Y.-F.H. designed and implemented the LINSIGHT method; Y.-F.H. and B.G. analyzed the data; A.S. supervised the research; Y.-F.H. and A.S. wrote the manuscript with review and feedback from B.G.

## COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

1. Mouse Genome Sequencing Consortium. Initial sequencing and comparative analysis of the mouse genome. *Nature* **420**, 520–562 (2002).
2. Hindorf, L.A. *et al.* Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc. Natl. Acad. Sci. USA* **106**, 9362–9367 (2009).
3. Gulko, B., Hubisz, M.J., Gronau, I. & Siepel, A. A method for calculating probabilities of fitness consequences for point mutations across the human genome. *Nat. Genet.* **47**, 276–283 (2015).
4. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* **15**, 1034–1050 (2005).
5. Wang, Y. *et al.* Sequencing and comparative analysis of a conserved syntenic segment in the Solanaceae. *Genetics* **180**, 391–408 (2008).
6. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* **478**, 476–482 (2011).
7. Haudry, A. *et al.* An atlas of over 90,000 conserved noncoding sequences provides insight into crucifer regulatory regions. *Nat. Genet.* **45**, 891–898 (2013).
8. ENCODE Project Consortium. Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
9. Gerstein, M.B. *et al.* Integrative analysis of the *Caenorhabditis elegans* genome by the modENCODE project. *Science* **330**, 1775–1787 (2010).
10. Roy, S. *et al.* Identification of functional elements and regulatory circuits by *Drosophila* modENCODE. *Science* **330**, 1787–1797 (2010).
11. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
12. Roadmap Epigenomics Consortium. Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
13. Ritchie, G.R.S., Dunham, I., Zeggini, E. & Flicek, P. Functional annotation of noncoding sequence variants. *Nat. Methods* **11**, 294–296 (2014).
14. Shihab, H.A. *et al.* An integrative approach to predicting the functional effects of noncoding and coding sequence variation. *Bioinformatics* **31**, 1536–1543 (2015).
15. Alipanahi, B., DeLong, A., Weirauch, M.T. & Frey, B.J. Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat. Biotechnol.* **33**, 831–838 (2015).
16. Zhou, J. & Troyanskaya, O.G. Predicting effects of noncoding variants with deep-learning-based sequence model. *Nat. Methods* **12**, 931–934 (2015).
17. Kelley, D.R., Snoek, J. & Rinn, J.L. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* **26**, 990–999 (2016).
18. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
19. Quang, D., Chen, Y. & Xie, X. DANN: a deep-learning approach for annotating the pathogenicity of genetic variants. *Bioinformatics* **31**, 761–763 (2015).
20. Fu, Y. *et al.* FunSeq2: a framework for prioritizing noncoding regulatory variants in cancer. *Genome Biol.* **15**, 480 (2014).
21. Gronau, I., Arbiza, L., Mohammed, J. & Siepel, A. Inference of natural selection from interspersed genomic elements based on polymorphism and divergence. *Mol. Biol. Evol.* **30**, 1159–1171 (2013).
22. Arbiza, L. *et al.* Genome-wide inference of natural selection on human transcription factor binding sites. *Nat. Genet.* **45**, 723–729 (2013).
23. Stenson, P.D. *et al.* The human gene mutation database: building a comprehensive mutation repository for clinical and molecular genetics, diagnostic testing, and personalized genomic medicine. *Hum. Genet.* **133**, 1–9 (2013).
24. Landrum, M.J. *et al.* ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980–D985 (2014).
25. Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R. & Siepel, A. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110–121 (2010).
26. Rumelhart, D., Hinton, G. & Williams, R. Learning representations by back-propagating errors. *Nature* **323**, 533–536 (1986).
27. Zerbino, D.R., Wilder, S.P., Johnson, N., Juettemann, T. & Flicek, P.R. The Ensembl Regulatory Build. *Genome Biol.* **16**, 56 (2015).
28. Gaffney, D.J., Blehman, R. & Majewski, J. Selective constraints in experimentally defined primate regulatory regions. *PLoS Genet.* **4**, e1000157 (2008).
29. Chiaromonte, F. *et al.* The share of human genomic DNA under selection estimated from human–mouse genomic alignments. *Cold Spring Harb. Symp. Quant. Biol.* **68**, 245–254 (2003).
30. Meader, S.J., Ponting, C.P. & Lunter, G. Massive turnover of functional sequence in human and other mammalian genomes. *Genome Res.* **20**, 1335–1343 (2010).
31. Rands, C.M., Meader, S., Ponting, C.P. & Lunter, G. 8.2% of the human genome is constrained: variation in rates of turnover across functional element classes in the human lineage. *PLoS Genet.* **10**, e1004525 (2014).
32. Lesurf, R. *et al.* ORegAnno 3.0: a community-driven resource for curated regulatory annotation. *Nucleic Acids Res.* **44**, D126–D132 (2016).
33. Davydov, E.V. *et al.* Identifying a high fraction of the human genome to be under selective constraint using GERP. *PLoS Comput. Biol.* **6**, e1001025 (2010).
34. Ionita-Laza, I., McCallum, K., Xu, B. & Buxbaum, J.D. A spectral approach integrating functional genomic annotations for coding and noncoding variants. *Nat. Genet.* **48**, 214–220 (2016).
35. Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455–461 (2014).
36. Core, L.J. *et al.* Analysis of nascent RNA identifies a unified architecture of transcription initiation regions at mammalian promoters and enhancers. *Nat. Genet.* **46**, 1311–1320 (2014).
37. Andersson, R., Sandelin, A. & Danko, C.G. A unified architecture of transcriptional regulatory elements. *Trends Genet.* **31**, 426–433 (2015).
38. Duret, L. & Mouchiroud, D. Determinants of substitution rates in mammalian genes: expression pattern affects selection intensity but not mutation rate. *Mol. Biol. Evol.* **17**, 68–74 (2000).
39. Drummond, D.A., Bloom, J.D., Adami, C., Wilke, C.O. & Arnold, F.H. Why highly expressed proteins evolve slowly. *Proc. Natl. Acad. Sci. USA* **102**, 14338–14343 (2005).
40. Kosiol, C. *et al.* Patterns of positive selection in six mammalian genomes. *PLoS Genet.* **4**, e1000144 (2008).
41. Voight, B.F., Kudravalli, S., Wen, X. & Pritchard, J.K. A map of recent positive selection in the human genome. *PLoS Biol.* **4**, e72 (2006).
42. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554–566 (2015).
43. Rao, S.S. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665–1680 (2014).
44. Guo, Y. *et al.* CRISPR inversion of CTCF sites alters genome topology and enhancer–promoter function. *Cell* **162**, 900–910 (2015).
45. Tang, Z. *et al.* CTCF-mediated human 3D genome architecture reveals chromatin topology for transcription. *Cell* **163**, 1611–1627 (2015).
46. Wunderlich, Z. *et al.* Kruppel expression levels are maintained through compensatory evolution of shadow enhancers. *Cell Rep.* **12**, 1740–1747 (2015).

## ONLINE METHODS

**Genomic features.** The genomic features used by LINSIGHT can be divided into three categories: conservation scores, predicted binding sites, and regional annotations (Table 2 and Supplementary Table 3). Conservation scores included phyloP scores<sup>25</sup>, phastCons elements<sup>4</sup>, SiPhy omega elements<sup>6,47</sup>, and CEGA elements<sup>48</sup>. Except for SiPhy, each score type was represented by multiple data tracks—for example, phastCons tracks for vertebrate, mammalian, and primate alignments (Supplementary Table 3). Predicted binding sites included transcription factor binding sites (TFBS) and RNA binding sites. Predicted TFBSs were obtained from the conserved TFBS track in the UCSC Genome Browser<sup>49</sup>, the rVISTA database<sup>50</sup>, SwissRegulon<sup>51</sup>, FunSeq2 (ref. 20), and the Ensembl Regulatory Build<sup>27</sup>. RNA binding sites included splice sites predicted by SPIDEX<sup>52</sup> (<http://www.deepgenomics.com/spidex/>) and miRNA target sites predicted by TarBase<sup>53</sup>. The regional annotations were based on a variety of sources, including ChIP-seq and RNA-seq data from the ENCODE<sup>11</sup> and Roadmap Epigenomics<sup>12</sup> projects, enhancers from FANTOM5 (ref. 35), predicted distal regulatory modules from FunSeq2 (ref. 20), and the distances to the nearest TSSs based on GENCODE gene models<sup>54</sup>. All features and the resulting LINSIGHT scores were expressed in genomic coordinates for the hg19 assembly of the human genome.

**Polymorphism and divergence data.** The polymorphism and divergence data used by the INSIGHT component of the LINSIGHT model were borrowed from previous analyses<sup>3,21,22</sup>. Briefly, we obtained human single-nucleotide polymorphisms (SNPs) from high-coverage genome sequences for 54 unrelated individuals from the ‘69 Genomes’ data set from Complete Genomics (<http://www.completegenomics.com/public-data/69-Genomes/>), after eliminating nucleotide sites with more than two alleles. Outgroup alleles were defined by the aligned chimpanzee, orangutan, and rhesus macaque reference genomes from UCSC. Several filters were applied to these data to reduce technical errors from alignment, sequencing, and genotype inference; for example, we removed simple repeats, recent transposable elements, recent segmental duplications, and regions not in syntenic alignment across primates<sup>22</sup>. Putatively neutral regions were identified by starting with all of the aligned regions and then removing the coding sequences, conserved noncoding sequences, and their proximal flanking regions. These regions were used to estimate neutral divergence and polymorphism rates in the human lineage in a block-wise manner across the genome, to account for local variation in mutation rates<sup>21</sup>. To allow for uncertainty in the human–chimpanzee most recent common ancestor (MRCA), we integrated over a distribution of ancestral alleles inferred after fitting a standard phylogenetic model to the outgroup sequences<sup>21</sup>.

**Generalized linear model.** The selection parameters in the INSIGHT model,  $\rho$  and  $\gamma$ , were defined as linear-sigmoid functions of the local genomic features at each nucleotide site  $i$ . Specifically, if  $D_i$  is a column vector of genomic features at site  $i$ , then

$$\rho_i = g(W_\rho \times D_i) \quad \text{and} \quad \gamma_i = h(W_\gamma \times D_i), \quad (1)$$

where the row vectors  $W_\rho$  and  $W_\gamma$  consist of feature weights (free parameters in the model), and  $g(\cdot)$  and  $h(\cdot)$  are sigmoid functions that map all input values to the range (0,1). For  $h(\cdot)$ , we used the standard logistical function,  $h(x) = 1/(1 + e^{-x})$ . For  $g(\cdot)$ , however, we used the asymmetric Gompertz sigmoid function<sup>55</sup>,  $g(x) = \exp(-3\exp(-x))$ , which ensured that the gradients were not too small when  $\rho_i$  was close to zero and allowed for accelerated convergence during model fitting.

**Fitting the LINSIGHT model to the data.** The weights for all of the genomic features were estimated by approximately maximizing the log likelihood of the INSIGHT model with respect to our genome-wide data set. We began by considering all genomic positions that were not excluded by our data-quality filters. Because our focus was on noncoding regions, we additionally excluded coding regions that were annotated by GENCODE (release 19). Instead of a traditional ‘batch’ learning algorithm, which would require either storing all of the data in computer memory or reading it from disk many times, we used an ‘online’ stochastic gradient-descent algorithm<sup>56</sup>. The algorithm processed the genome sequentially in ‘minibatches’ of 100 successive nucleotides, each time updating

the parameter vector in the direction of the gradient of the log-likelihood function, with learning rates of 0.001 and 0.01 for  $\rho$  and  $\gamma$ , respectively. Gradients were computed analytically by propagating partial derivatives through the linear-sigmoid component of the model using the chain rule (back-propagation). The learning procedure was truncated after 20 passes through the entire data set. The entire process took less than 1 d on a desktop computer. The genome-wide LINSIGHT scores are available from the Cold Spring Harbor Laboratory mirror of the UCSC Genome Browser (hg19 assembly).

**Comparison with other methods.** Our benchmarking scheme for prioritization of disease-associated variants closely followed the one introduced previously<sup>13</sup>. The HGMD and ClinVar noncoding disease variants and three sets of negative controls were obtained directly from this study. The negative controls consisted of: (i) a randomly selected subset of human common variants which was 100-fold larger than the set of HGMD variants (unmatched); (ii) a subset of human common variants that were matched to the disease variants by the exact distance-to-nearest TSS (matched TSS) (although each negative example was not necessarily near the same TSS as the matched disease variant); and (iii) a subset of human common variants that were required to be within 1 kb of the matched disease variants (matched region). The two matched sets accounted for the enrichment of known disease variants near coding genes. We later defined three additional sets of negative controls by the same strategy but using singleton variants from the 1000 Genomes Project phase 3 data<sup>57</sup> instead of common variants, to ensure that our results were not driven by differences in allele frequency between the disease variants and negative controls. In all cases, we subsampled the negative sets to balance the numbers of positive and negative examples. To reduce stochasticity, subsampling was performed 100 times, and average performance statistics were reported.

For comparison, we downloaded precomputed CADD (v1.3)<sup>18</sup>, GWAVA (v1.0)<sup>13</sup>, FunSeq2 (v2.1.0)<sup>20</sup>, and Eigen<sup>34</sup> (Oct. 11, 2015) scores from the source websites. In all cases, we used GWAVA scores based on training with variants matched by distance to nearest TSS<sup>13</sup>. In addition, we obtained mammalian phyloP<sup>25</sup> scores based on the 46-way whole-genome alignment for hg19 from the UCSC Genome Browser<sup>49</sup>, and we computed DeepSEA functional significance scores for both disease variants and negative controls by using the online DeepSEA web service<sup>16</sup> (computed on Nov 3, 2016). The DeepSEA functional significance scores integrate individual tissue-specific DeepSEA scores based on polymorphism data; these were used in all of the comparisons because the tissue types associated with disease variants and ORegAnno TFBSs are typically unknown. Note that two of the methods considered, CADD and DeepSEA, provide allele-specific predictions, whereas the other methods assign identical scores to all alternative variants. When evaluating CADD and DeepSEA on the ClinVar data set, we used the score corresponding to the annotated disease-associated allele. When evaluating these methods on the HGMD data set, however, no disease-associated allele was provided, so we used the maximum score for the three alternative alleles.

**Classification of disease-associated variants by genomic location.** For analyses that considered the genomic locations of disease-associated variants, we divided the variants in the HGMD and ClinVar databases into four categories based on their locations relative to the gene models from GENCODE (release 19). These categories were: (i) ‘promoter’ variants, which are located within 1 kb upstream of the 5′-most annotated TSS of any protein-coding gene; (ii) ‘splicing’ variants, which are located within 20 bp of any annotated splice junction; (iii) ‘UTR’ variants, which are located within the annotated 5′ or 3′ UTR of any protein-coding gene; and (iv) all ‘other’ variants. Each variant was assigned to the first category whose criteria it fulfilled in the order: splicing > UTR > promoter > other.

**Quantification of the contributions of genomic feature classes.** We measured the relative contributions of the conservation scores, predicted binding sites, and regional annotations by removing all of the features of each class (Table 2), retraining the LINSIGHT model without those features, and evaluating the AUC of the reduced model. The ‘contribution’ of each class of features was defined as the AUC for the full model minus the AUC for the reduced model, averaged across 100 independent subsamples of the negative



controls described above. Notice that, although this difference in AUCs was generally positive, it could be negative due to stochastic effects. This analysis was performed on a merged set of HGMD and ClinVar variants, separately for promoter, splicing, UTR, and other regions.

**Analysis of evolutionary constraints on enhancers.** To study evolutionary constraints on enhancers, we used the comprehensive atlas of human enhancers based on enhancer RNAs (eRNAs) that was recently provided by the FANTOM5 project<sup>35</sup>. The evolutionary constraint for each enhancer was quantified by taking the average LINSIGHT score across all nucleotide sites in the enhancer. We examined the relationship between this measure of constraint and the number of cell types in which each enhancer was active (according to a detectable eRNA signature). We also defined a subset of enhancers as tissue specific based on apparent activity in only a single tissue type, and we examined the relationship between tissue of activity and degree of constraint. Finally, we obtained putative enhancer–TSS pairs (based on correlated patterns of expression across tissues) from the FANTOM5 website and examined the correlation in constraint at the enhancer and promoter in each pair, defining the promoter as the 1-kb region upstream of the TSS. In cases where an enhancer was associated with multiple TSSs, the TSS with the highest correlation coefficient was selected.

**Statistical analysis.** To examine the relationship between evolutionary constraints on enhancers and tissue specificity, we calculated the Spearman's rank correlation coefficient between the average LINSIGHT score for each enhancer and its number of active cell types. To quantify the statistical significance of the correlation, a two-tailed *P* value was computed using the standard asymptotic *t* approximation implemented in the 'cor.test' function in R ( $P < 10^{-15}$ ;  $n = 29,303$ ). The same method was used to quantify the statistical significance of the correlation between the average LINSIGHT scores at enhancer–promoter pairs ( $P < 10^{-15}$ ;  $n = 25,067$ ). Furthermore, to investigate the relationship between the average LINSIGHT score in an enhancer and the number of active cell types when controlling for average eRNA expression level, the partial Spearman's  $\rho$  and a two-tailed *P* value were computed using the ppcor package<sup>58</sup> ( $P < 10^{-15}$ ;

$n = 29,303$ ). To investigate whether the difference between two AUCs was statistically significant, the DeLong test was used to compute two-tailed *P* values<sup>59</sup>.

**Code availability.** The LINSIGHT code is available at <https://github.com/CshSiepelLab/LINSIGHT>.

**Data availability.** The training data and pre-computed LINSIGHT scores are available at <http://compgen.cshl.edu/~yihuang/LINSIGHT/>.

47. Garber, M. *et al.* Identifying novel constrained elements by exploiting biased substitution patterns. *Bioinformatics* **25**, i54–i62 (2009).
48. Dousse, A., Junier, T. & Zdobnov, E.M. CEGA—a catalog of conserved elements from genomic alignments. *Nucleic Acids Res.* **44**, D96–D100 (2016).
49. Kent, W.J. *et al.* The human genome browser at UCSC. *Genome Res.* **12**, 996–1006 (2002).
50. Dubchak, I. *et al.* Whole-genome rVISTA: a tool to determine enrichment of transcription factor binding sites in gene promoters from transcriptomic data. *Bioinformatics* **29**, 2059–2061 (2013).
51. Pachkov, M., Balwierz, P.J., Arnold, P., Ozonov, E. & van Nimwegen, E. SwissRegulon, a database of genome-wide annotations of regulatory sites: recent updates. *Nucleic Acids Res.* **41**, D214–D220 (2013).
52. Xiong, H.Y. *et al.* The human splicing code reveals new insights into the genetic determinants of disease. *Science* **347**, 1254806 (2014).
53. Vlachos, I.S. *et al.* DIANA-TarBase v7.0: indexing more than half a million experimentally supported miRNA:mRNA interactions. *Nucleic Acids Res.* **43**, D153–D159 (2015).
54. Harrow, J. *et al.* GENCODE: the reference human genome annotation for the ENCODE Project. *Genome Res.* **22**, 1760–1774 (2012).
55. Gompertz, B. On the nature of the function expressive of the law of human mortality, and on a new mode of determining the value of life contingencies. *Philos. Trans. R. Soc. Lond.* **115**, 513–583 (1825).
56. Duchi, J., Hazan, E. & Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.* **12**, 2121–2159 (2011).
57. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
58. Kim, S. ppcor: an R package for a fast calculation to semi-partial correlation coefficients. *Commun. Stat. Appl. Methods* **22**, 665–674 (2015).
59. DeLong, E.R., DeLong, D.M. & Clarke-Pearson, D.L. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* **44**, 837–845 (1988).