

Interpreting *cis*-regulatory mechanisms from genomic deep neural networks using surrogate models

Evan E Seitz¹, David M McCandlish¹, Justin B Kinney^{1,*}, and Peter K Koo^{1,*}

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

*koo@cshl.edu and jkinney@cshl.edu

ABSTRACT

Deep neural networks (DNNs) have greatly advanced the ability to predict genome function from sequence. Interpreting genomic DNNs in terms of biological mechanisms, however, remains difficult. Here we introduce SQUID, a genomic DNN interpretability framework based on surrogate modeling. SQUID approximates genomic DNNs in user-specified regions of sequence space using surrogate models, i.e., simpler models that are mechanistically interpretable. Importantly, SQUID removes the confounding effects that nonlinearities and heteroscedastic noise in functional genomics data can have on model interpretation. Benchmarking analysis on multiple genomic DNNs shows that SQUID, when compared to established interpretability methods, identifies motifs that are more consistent across genomic loci and yields improved single-nucleotide variant-effect predictions. SQUID also supports surrogate models that quantify epistatic interactions within and between *cis*-regulatory elements. SQUID thus advances the ability to mechanistically interpret genomic DNNs.

Deep neural networks (DNNs) are increasingly being used to analyze and biologically interpret functional genomics data. DNNs have demonstrated remarkable success at predicting diverse genomic activities from primary genome sequences, including mRNA expression levels,¹ transcription start site activity,^{2,3} mRNA splicing patterns,⁴ protein-DNA binding,⁵ chromatin accessibility,⁶ and chromatin conformation.⁷ It is widely believed that the success of genomic DNNs reflects their ability to accurately model complex biological mechanisms. However, interpreting genomic DNNs in terms of biological mechanisms remains difficult.

A variety of *post hoc* attribution methods have been developed for aiding the interpretation of genomic DNNs.^{8,9} The most commonly used attribution methods produce *attribution maps*, which quantify the position-specific effects that variant nucleotides in a sequence of interest have on DNN predictions. Attribution maps are often mechanistically interpreted by identifying motifs, i.e. recurrent sequence patterns characteristic of specific biological mechanisms, present within these maps. To provide consistent biological explanations for sequence activity, attribution methods must produce consistent motifs across input sequences that function through shared biological mechanisms.

Established attribution methods for genomic DNNs have major limitations. Different attribution methods use different strategies to quantify position-specific nucleotide effects, and can therefore yield different mechanistic explanations.^{10–12} For example, Saliency Maps¹³ quantify nucleotide effects using the gradient of DNN predictions at the sequence of interest, whereas DeepLIFT¹⁴ propagates activation differences between the sequence of interest and a reference sequence. Moreover, the most widely-used attribution methods in genomics – including Saliency Maps,¹³ DeepLIFT,¹⁴ *in silico* mutagenesis (ISM),¹⁵ SmoothGrad,¹⁶ Integrated Gradients,¹⁷ and DeepSHAP¹⁸ – assume that nucleotide effects on DNN predictions are locally additive. As a result, these attribution methods do not account for the genetic interactions (i.e., specific epistasis^{19–22}), global nonlinearities (i.e., global epistasis^{23–27}), and heteroscedastic noise²⁶ that are often present in functional genomics data.

Here we introduce SQUID (Surrogate Quantitative Interpretability for Deepnets), an interpretability framework for genomic DNNs that overcomes these limitations. SQUID uses surrogate models—simple models with interpretable parameters—to approximate the DNN function within localized regions of sequence space. SQUID applies MAVE-NN,²⁶ a quantitative modeling framework developed for analyzing multiplex assays of variant effects (MAVEs), to *in silico* MAVE datasets generated using the DNN as an oracle. In this way, SQUID models DNN predictions in a user-specified region of sequence space, accounts for the nonlinearities and heteroscedastic noise present in DNN predictions, and (optionally) quantifies specific epistatic interactions. Benchmarking SQUID against existing attribution methods, we find that SQUID more consistently quantifies the binding motifs of transcription factors (TFs), reduces noise in attribution maps, and improves variant-effect predictions. We also find that the domain-specific surrogate models used by SQUID are critical for this improved performance. Finally, we show how SQUID can provide insights into epistatic interactions in *cis*-regulatory elements, and can be used to

study such interactions both locally and globally in sequence space. SQUID thus provides a new and useful way to interpret genomic DNNs.

Results

SQUID: An interpretable surrogate modeling framework for genomic DNNs

SQUID approximates DNNs in user-specified regions of sequence space using surrogate models that have mechanistically interpretable parameters. The SQUID framework comprises three steps (Fig. 1a): (1) generate an *in silico* MAVE dataset comprised of variant sequences and the corresponding DNN predictions; (2) fit a surrogate model to these *in silico* data using MAVE-NN²⁶; and (3) visualize and interpret the surrogate model's parameters. This workflow requires the specification of two key analysis parameters: the region of sequence space over which the DNN is to be approximated, and the mathematical form of the surrogate model. By choosing different values for these analysis parameters, users are able to test different mechanistic hypotheses. SQUID assists users by facilitating the use of analysis parameters that have been found in practice to work well in the design and analysis of MAVE experiments.

First, SQUID generates an *in silico* MAVE dataset. This is done by generating a library of variant sequences, then using the DNN to assign a functional score to each sequence in the library. In this paper we consider two types of libraries: local and global. A local library is generated by partially mutagenizing a specific sequence of interest (e.g., a genomic *cis*-regulatory sequence). A global library is generated by inserting partially-mutagenized versions of a genetic element of interest into random sequences. In what follows, we use a partial mutagenesis rate of 10% per nucleotide (which is common in MAVE experiments, e.g. ref.²⁸) and libraries comprising 100,000 variant sequences (unless otherwise noted).

Next, SQUID fits a surrogate model to the *in silico* MAVE dataset. SQUID uses nonlinear surrogate models developed specifically for modeling MAVE data.²⁶ These models, called latent phenotype models, have three components (Fig. 1b):

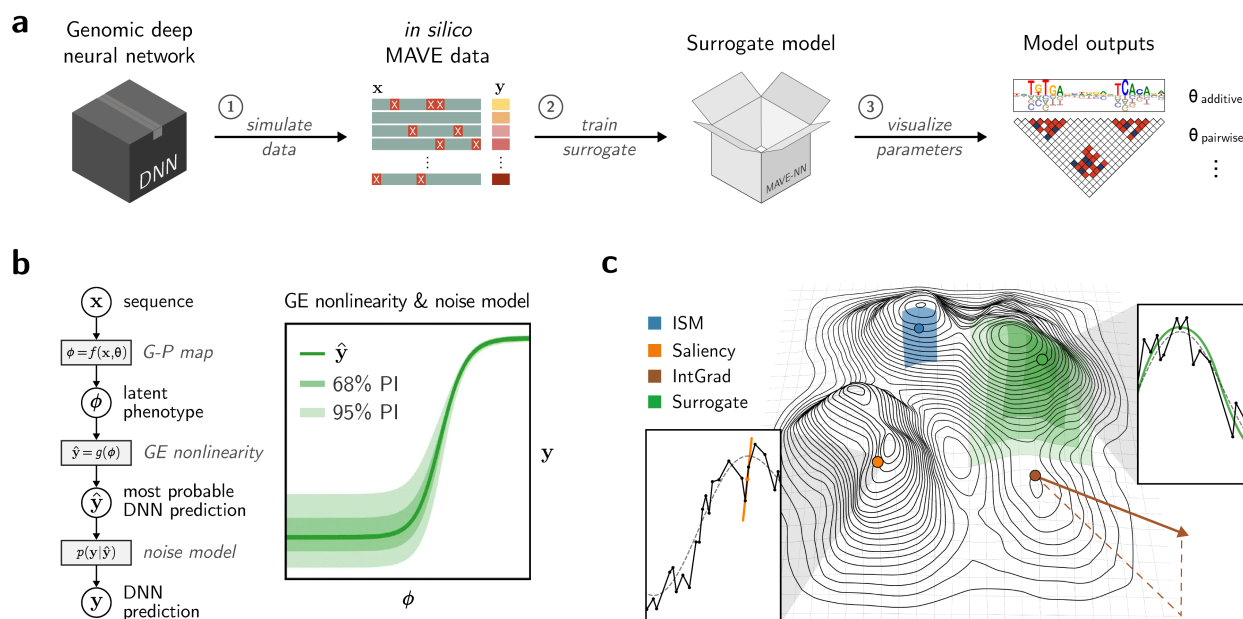


Figure 1. Overview of SQUID. **a**, Schematic of the SQUID modeling framework. Analysis using SQUID comprises three main steps: (1) generate an *in silico* MAVE dataset; (2) train a surrogate model on the MAVE dataset; and (3) visualize parameters of the surrogate model to uncover biological mechanisms. **b**, Structure of the latent phenotype surrogate models supported by SQUID. G-P, genotype-phenotype; GE, global epistasis; PI, prediction interval. **c**, Schematic diagram of a DNN function on a 2-dimensional projection of sequence space. Each point in the plane corresponds to a unique sequence, and elevations represent DNN predictions. Green region schematizes the ability of surrogate models to approximate the DNN function over an extended region of sequence space. Insets show an example DNN function (in 1D profile; black line) centered about a sequence of interest with the ground-truth function (dashed line) overlaid. Left inset illustrates the sensitivity of Saliency Maps to non-smooth local function properties. Right inset illustrates the ability of surrogate models to better approximate ground truth. ISM, *in silico* mutagenesis; IntGrad, Integrated Gradients.

a genotype-phenotype (G-P) map, a global epistasis (GE) nonlinearity, and a noise model. The G-P map projects the input sequence down to a one-dimensional latent phenotype. SQUID, via MAVE-NN, supports the use of additive G-P maps, pairwise-interaction G-P maps, and user-defined G-P maps. The GE nonlinearity, which is modeled using a linear combination of sigmoids, maps the latent phenotype to a most-probable DNN prediction. The noise model then describes how actual DNN predictions are expected to scatter around the predicted most-probable value. SQUID supports a variety of noise models, including heteroscedastic noise models based on the skewed-t distribution of Jones and Faddy²⁹ (which we use in this paper). Surrogate model parameters are inferred from *in silico* MAVE data by maximizing variational information²⁶ (or equivalently, log likelihood) using standard stochastic gradient optimization.

The parameters of the G-P map are of primary interest for downstream mechanistic interpretations. The parameters of additive G-P maps quantify single-nucleotide effects, and are readily visualized using the same methods (sequence logos³⁰ and heatmaps) normally used for standard attribution maps. The parameters of pairwise-interaction G-P maps quantify epistatic genetic interactions as well as single-nucleotide effects, and can be visualized using block-triangular heat maps (as in ref.²⁶).

SQUID improves the quantification of transcription factor binding motifs

A common goal when interpreting attribution maps for a DNA sequence of interest is to identify functional binding sites for TFs. However, the TF binding motifs observed in attribution maps often vary substantially from sequence to sequence. This variation results in part from the underlying biology, but it can be exacerbated by the specific ways in which attribution methods quantify the behavior of the DNN function in localized regions of sequence space (Fig. 1c). The consistency of binding motifs observed in attribution maps across different genomic sequences therefore provides a way to quantify and compare the performance of different attribution methods.

To benchmark the consistency of binding motifs identified by different attribution methods, we identified sequences in the human genome that contain putative TF binding sites. Specifically, for each TF, we identified genomic instances of the consensus TF binding site sequence. We then aligned these genomic sequences about their putative binding sites and computed, for each genomic sequence, an attribution map that spans the core putative binding site as well as a specified amount of flanking DNA on either side. Finally, we calculated the Euclidean distance between the vector of attribution scores for individual sequences and the mean vector of attribution scores. We refer to this distance as the “attribution error” (Fig. 2a; see Methods for details).

We first applied this benchmarking pipeline to the human TF AP-1 using ResidualBind-32,³¹ a genomic DNN that predicts chromatin accessibility in human cell lines. We compared SQUID to two commonly used attribution methods that had previously been used in ref.³¹ to analyze ResidualBind-32: *in silico* mutagenesis (ISM) and Saliency Maps. We found that, across different genomic sequences that contain the consensus AP-1 binding site TGAGTCA, SQUID recovers AP-1 binding motifs that have markedly-lower attribution errors than the binding motifs recovered by ISM or Saliency Maps (Fig. 2b, left). This result is supported by the examination of attribution maps of individual sequences (Fig. 2b, right). Compared to the attribution maps provided by other methods, the attribution maps provided by SQUID exhibit core regions with greater similarity to the ensemble-averaged motif, and flanking regions with reduced (likely non-biological) scores. This finding was robust to the choice of SQUID analysis parameters (Supplementary Fig. 1). We conclude that SQUID quantifies the AP-1 binding motif from ResidualBind-32 more consistently than does ISM or Saliency Maps.

We next investigated whether the surrogate models for AP-1 identified by SQUID benefited from incorporating a GE nonlinearity. Plotting the effects that mutations in a representative genomic sequence have on DNN predictions, we found that virtually all combinations of 2 or more mutations to the core 7-nt AP-1 site reduced DNN predictions to near-background levels (Fig. 2c). Moreover, the GE nonlinearity learned by SQUID as part of the surrogate model accurately recapitulated this saturation effect. By contrast, surrogate modeling using ridge regression (which is implemented in SQUID as a baseline method and does not incorporate a GE nonlinearity or heteroscedastic noise model; see Methods for details) failed to capture this saturation effect (Figs. 2d and 2e). This finding demonstrates that surrogate modeling of genomic DNNs can benefit from using our domain-specific surrogate models, as opposed to the linear surrogate models that are standard in other fields (e.g., computer vision³²).

We then expanded our analysis to other TFs and to other genomic DNNs (DeepSTARR³³ and BpNet⁵; see Supplementary Table 1). In each benchmark analysis, we compared SQUID to the attribution methods (ISM, Saliency Maps, DeepLIFT, or DeepSHAP) used in the original study describing the DNN being modeled. We found that the attribution maps provided by SQUID consistently yielded binding motifs with markedly-lower attribution error than the binding motifs provided by the other attribution methods (Fig. 3a). These results were robust to the amount of flanking DNA used when computing attribution errors (Fig. 3b). We also found that strong GE nonlinearities were pervasive in the surrogate models inferred by SQUID for the genomic sequences tested (Supplementary Fig. 2). These results suggest that SQUID, quite generally, quantifies TF binding motifs more consistently than do competing attribution methods. These findings also suggest that modeling GE nonlinearities and heteroscedastic noise is important for the accurate surrogate modeling of genomic DNNs.

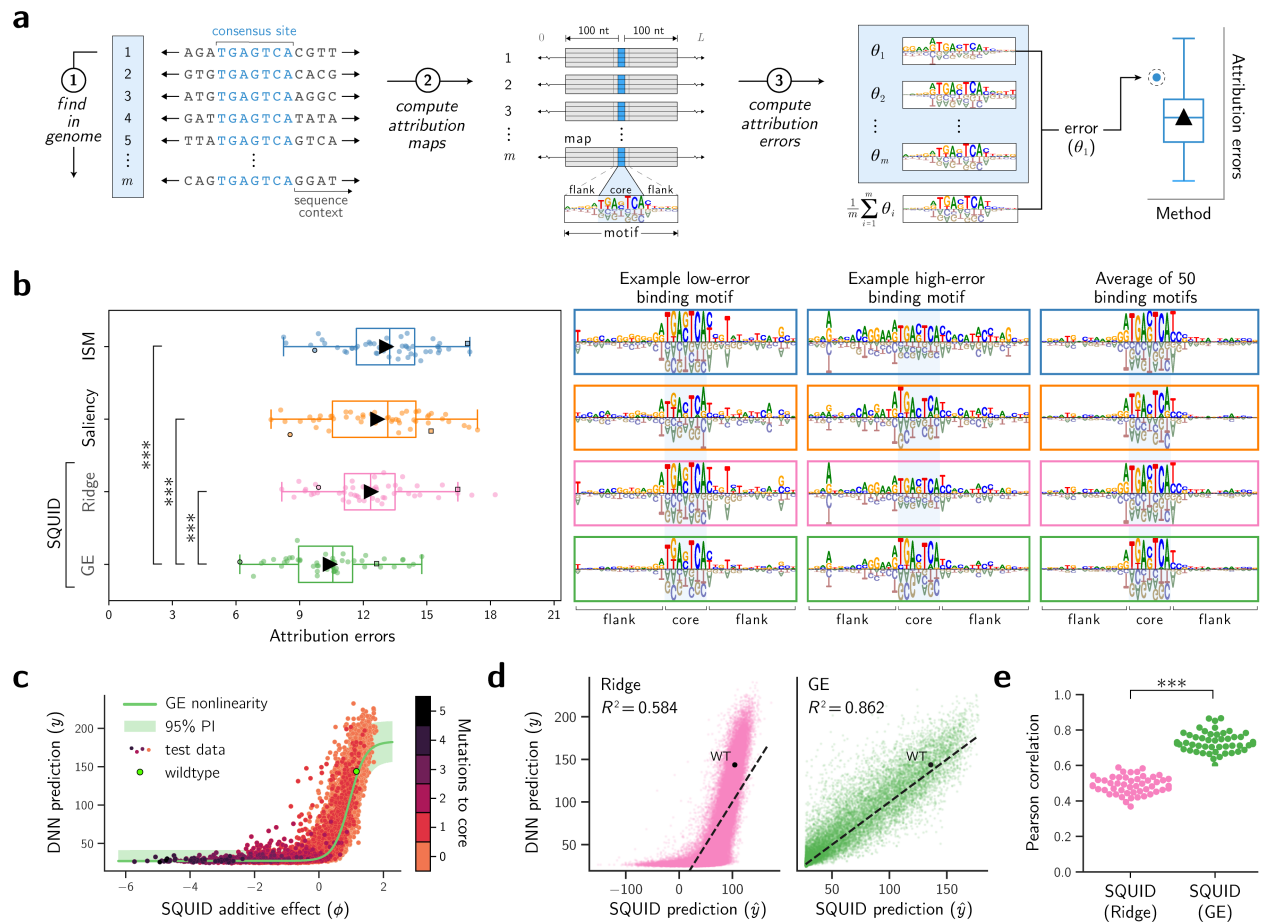


Figure 2. Benchmark analysis of attribution methods. **a**, Our benchmark analysis pipeline consisted of 3 steps: (1) genomic sequences that contained consensus binding sites for a TF of interest were identified; (2) an attribution map spanning the core identified site and 100 nt of flanking sequence on each side was computed for each identified sequence; and (3) a corresponding attribution error was computed (see Methods for details). **b**, Attribution maps and attribution errors for 7-nt consensus AP-1 binding sites (core sequence TGAGTGA; flank size 15 nt), computed using ResidualBind-32 as the DNN. *Left*, Box plots show attribution errors for 50 genomic sequences. Triangles represent means; lines represent median, upper quartile, and lower quartile; whiskers represent 1.5x the inter-quartile range. *Right*, Binding motifs observed for two example genomic sequences, together with the ensemble-averaged binding motif. **c**, Mutational effects predicted by the DNN versus additive effects predicted by SQUID. Dots represent effects of different sets of mutations to a representative genomic sequence. The GE nonlinearity and 95% PI inferred by SQUID are shown for comparison. **d**, DNN predictions versus the predictions of two surrogate models inferred by SQUID for a representative genomic sequence. One surrogate model (GE, right) has a GE nonlinearity, the other (Ridge, left) does not. Dots represent test sequences from the *in silico* MAVE dataset. Diagonal lines represent equality between DNN and surrogate model predictions. WT, wild-type sequence; R^2 , squared Pearson correlation coefficient. **e**, R^2 values computed as in panel **d** for 50 different sequences of interest. p -values in panels **b** and **e** were computed using a one-sided Mann-Whitney U test; ***, $p < 0.001$. TF, transcription factor; DNN, deep neural network; GE, global epistasis; PI, prediction interval; ISM, *in silico* mutagenesis.

The ability of SQUID to better quantify TF binding motifs is exemplified in Figure 3c. Shown are attribution maps for the mouse TF Nanog, computed using BPNNet, a consensus binding site of AGCCATCAA, and 50 nt of flanking DNA. When attribution maps are averaged across genomic loci, ISM, DeepLIFT, and SQUID all produce attribution maps that reveal both the Nanog binding motif and a more subtle preference for periodically-spaced AT-rich sequences in the flanking DNA (which likely reflects the interaction of nucleosomes with the DNA double helix). However, the attribution maps for specific genomic sequences often exhibit many other noticeable features. Some of these features may be functional in principle, but many appear to be spurious and likely reflect non-biological noise in the attribution map. The attribution map provided by SQUID appears to

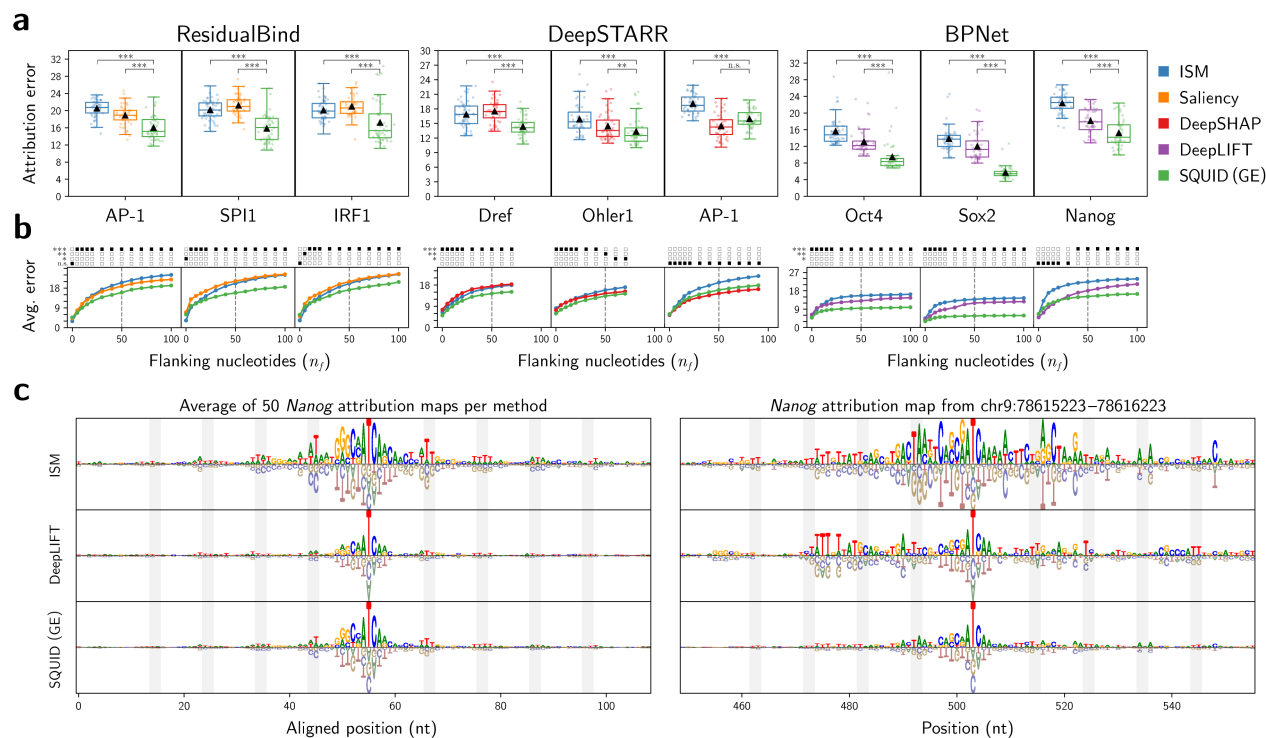


Figure 3. Benchmark analysis across TFs, DNNs, and flank sizes. **a**, Attribution error analysis for various TFs and genomic DNNs. Results are visualized as in Figure 2b. Each test used 50 sequences from either the human genome (for ResidualBind-32), the mouse genome (for BPNet), or the fly genome (for DeepSTARR), together with consensus TF binding site listed in Supplementary Table 1, and flanked by 50 nt of DNA. **b**, Mean attribution scores computed as in panel **a**, but using variable lengths of flanking DNA. The boxes above each plot indicate the largest (i.e., least significant) p-value, computed as in Figure 2b. n.s., $p \geq 0.05$; *, $0.01 \leq p < 0.05$; **, $0.001 \leq p < 0.01$; ***, $p < 0.001$. Dashed line, flank size used in panel **a**. **c**, Binding motifs computed for *Nanog* via attribution analysis of BPNet. *Left*, Attribution maps averaged across 50 mouse loci containing (and centered on) the consensus *Nanog* binding site AGCCATCAA. *Right*, Attribution map observed for a single such locus. Gray bars indicate 10.5 nt periodicity on either side of the consensus *Nanog* binding site. TF, transcription factor. DNN, deep neural network. ISM, *in silico* mutagenesis. GE, global epistasis.

exhibit less noise in the sequences that directly flank core motifs than the attribution maps given by ISM or DeepLIFT. Similar observations hold in analyses of other TFs (Supplementary Figure 3). These findings raise the possibility that the increased consistency of binding motifs identified by SQUID is due, at least in part, to the reduction of non-biological noise in attribution maps.

SQUID reduces noise in attribution maps

Non-biological noise in attribution maps can arise from the inherent roughness in the DNN function, a phenomenon that is often associated with benign overfitting.^{34,35} Benign overfitting is of particular concern for Saliency Maps, since it can adversely affect the accuracy of DNN gradients.^{36–38} Benign overfitting is similarly expected to impact other attribution methods, since these methods essentially quantify how small changes in sequence space affect DNN predictions.^{10,37} We reasoned that, because SQUID integrates information over an extended region of sequence space, the attribution maps provided by SQUID might be less noisy than the maps provided by other attribution methods.

We therefore asked whether SQUID can reduce the attribution map noise caused by benign overfitting. To answer this question, we trained a DNN to classify ChIP-seq peaks for the human TF GABPA (see Methods) and saved DNN parameters both before benign overfitting (i.e., using early stopping) and after benign overfitting (Fig. 4a). Benign overfitting is apparent from the plateau in validation-set performance that occurs when near-perfect classification performance is achieved on the training set.³⁹ We then selected 100 random sequences from the test set, and for each sequence and each of four attribution methods (SQUID, DeepSHAP, SmoothGrad, and Saliency Maps), quantified the differences between the attribution maps obtained using the pre-overfitting DNN parameters versus the post-overfitting DNN parameters. Figure 4b shows that SQUID

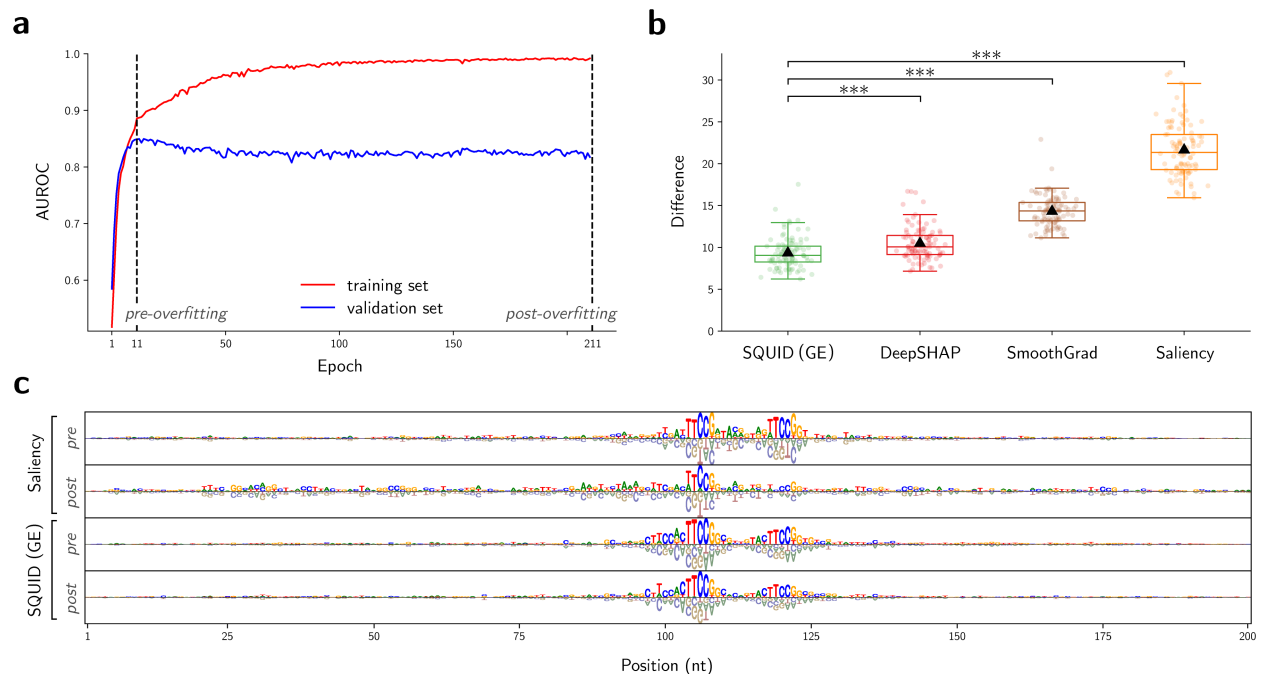


Figure 4. Attribution method performance during benign overfitting. **a**, DNN performance as a function of training epoch. DNN was a 3-layer convolutional neural network trained to classify 200-nt ChIP-seq peaks for the human TF GABPA. Tests used DNN parameters from epoch 11 (pre-overfitting) and epoch 211 (post-overfitting). AUROC, area under the receiver-operator characteristic curve. **b**, Differences between attribution maps for the DNN with pre-overfitting parameters versus post-overfitting parameters, as quantified by the Euclidean distance between attribution map vectors. Results are shown for 100 genomic sequences in the ChIP-seq peak test set. Data is visualized and statistical tests were performed as in Figure 2b. **c**, Attribution maps obtained for a representative test sequence using pre-overfitting and post-overfitting DNN parameters. DNN, deep neural network; TF, transcription factor; GE, global epistasis.

provided attribution maps that changed substantially less over the course of benign overfitting. Figure 4c illustrates this behavior for a representative sequence. These results support the hypothesis that the attribution maps provided by SQUID are more robust to the adverse effects of benign overfitting than are other attribution methods. These results also suggest that SQUID, more generally, yields attribution maps that have lower noise than maps computed by other attribution methods.

SQUID better identifies putative weak TF binding sites.

Weak TF binding sites play critical roles in eukaryotic gene regulation.^{40–42} Functional signals from weak binding sites, however, can be difficult to distinguish from noise in attribution maps. Having shown that SQUID reduces noise in attribution maps relative to other attribution methods, we hypothesized that SQUID would also better identify weak yet functional TF binding sites. To test this hypothesis, we quantified how well attribution maps generated for putative weak TF binding sites matched the TF binding motifs identified in Figure 3a. For each TF of interest, we randomly selected 150 putative binding sites in the genome having 0, 1, or 2 mutations relative to the consensus binding site (see Supplementary Table 1). We also recorded the genomic sequence containing each selected site padded by 50 nt of flanking DNA on either side. We then computed the score assigned to each selected site using a TF-specific position weight matrix (PWM),⁴³ as in Figure 5a. Different attribution methods were then used to compute attribution maps for each genomic sequence, after which each sequence was then assigned an attribution error quantified by the Euclidean distance between the attribution map and the corresponding ensemble-averaged attribution map from Figure 3a. Figure 5b shows that, as expected, the resulting errors for all attribution methods increased as PWM score decreased. However, the attribution errors obtained by SQUID were in general consistently lower than the attribution errors obtained by competing methods. This finding is confirmed by visually examining attribution maps for selected sites (Supplementary Fig. 4). We conclude that SQUID is better than competing methods at identifying signatures of TF binding at weak yet functional TF binding sites.

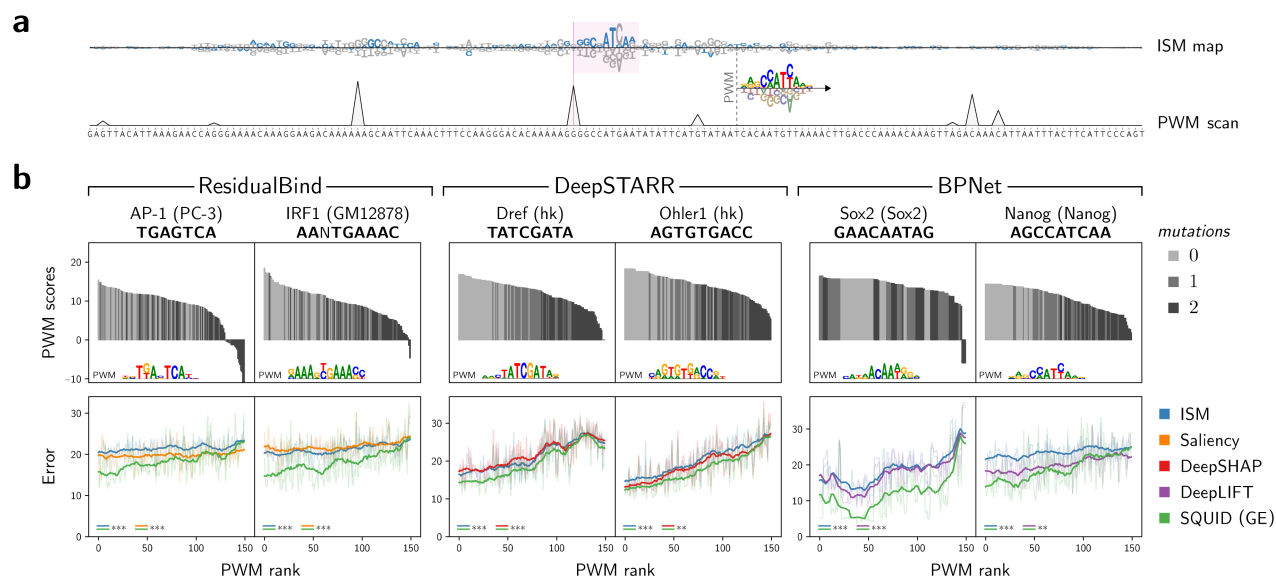


Figure 5. Benchmark analysis of attribution errors for putative weak TF binding sites. **a**, *Upper*, Attribution map for a representative genomic sequence containing multiple putative weak binding sites for the mouse TF Nanog. Attribution map is for the BPNNet DNN and was computed using ISM. Blue, wild-type nucleotides; gray, non-wild-type nucleotides. *Middle*, PWM for Nanog. *Lower*, PWM scores across the genomic sequence. Only positive PWM scores are shown. **b**, For each TF and DNN, plots show attribution errors for 150 putative TF binding sites plotted against putative binding site strength. Bold lines indicate signals smoothed with a sliding window of 20 nt. Stars indicate p-values computed using the one-sided Mann-Whitney U test: **, $0.001 \leq p < 0.01$; ***, $p < 0.001$. PWM scores for each of the 150 putative sites are shown above, along with a logo representation of the PWM used. Each site is represented by a gray bar shaded according to the number of mutations (0, 1, or 2) in the core of the putative site.

SQUID improves the zero-shot prediction of single-nucleotide variant effects

A major goal of genomic DNNs and their attribution methods is to predict which genetic variants are pathogenic. Having observed that SQUID reduces noise in attribution maps, mitigates the adverse effects of benign overfitting, and better identifies TF binding motifs, we hypothesized that SQUID would provide improved variant effect predictions. To test this hypothesis, we used SQUID and other attribution methods to predict the effects that single-nucleotide variants (SNVs) have on the activity of *cis*-regulatory elements for 15 different disease-associated loci in the human genome, for which the effects of SNVs had been measured using massively parallel reporter assays (MPRAs)^{44,45}. In technical terms, this is a zero-shot prediction task, since none of the methods tested were trained on MPRA data. At each locus, the Pearson correlation coefficient was computed between the attribution scores and measured SNV effects. We performed this benchmark analysis using three genomic DNNs previously reported for predicting regions of open chromatin: Enformer,² Basenji-32,⁴⁶ and ResidualBind-32.³¹ We found that SQUID yielded substantially higher correlations between predicted and measured SNV effects across the 15 assayed loci (Table 1 and Supplementary Fig. 5). The results suggest that attribution maps provided by SQUID are generally able to predict variant effects better than the attribution maps provided by competing methods. The improved performance of SQUID relative to ISM, which directly uses DNN predictions to quantify SNV effect sizes, suggests that surrogate modeling may provide a general way of improving the variant effect predictions of genomic DNNs themselves.

SQUID illuminates epistatic interactions

Understanding the role of epistatic interactions within gene regulatory sequences is a major goal in the study of *cis*-regulatory codes. An important advantage of surrogate modeling over other DNN interpretability approaches is that surrogate models with different mathematical forms can be used to answer different questions about DNN behavior. To test whether SQUID could be used to study epistatic interactions in *cis*-regulatory sequences, we implemented a surrogate model that describes all possible pairwise interactions between nucleotides within a sequence (in addition to the additive contributions from individual nucleotides). We then used this pairwise-interaction model to quantify the effects of pairs of putative AP-1 binding sites learned by ResidualBind-32.³¹ We identified 50 genomic sequences having pairs of putative AP-1 binding sites spaced 4 to 20 nt apart, used SQUID to infer the parameters of the pairwise-interaction model about each of the 50 sequences, and then averaged the

DNN		Average Pearson correlation				Statistical significance		
Architecture	Task	Saliency	ISM	SQUID (Ridge)	SQUID (GE)	GE vs. Saliency	GE vs. ISM	GE vs. Ridge
Enformer	DNASE	0.2977	0.4498	0.4486	0.4801	***	**	**
Basenji-32	ATAC-seq	0.2727	0.3575	0.3700	0.4036	***	***	**
ResBind-32	ATAC-seq	0.2846	0.3388	0.3567	0.3912	**	***	*

Table 1. Attribution method performance on a zero-shot variant-effect prediction task. Performance of four attribution methods, applied to three DNNs, on the CAGI5 variant-effect prediction challenge.⁴⁴ Numbers indicate the average correlation observed across the 15 genomic loci assayed by MPRA in ref.⁴⁴. Bold text indicates the best-performing method for each DNN. P-values report results from a paired Wilcoxon signed-rank test on the 15 locus-specific correlation values: *, $0.01 \leq p < 0.05$; **, $0.001 \leq p < 0.01$; ***, $p < 0.001$. MPRA, massively parallel reporter assay; DNN, deep neural network; ISM, *in silico* mutagenesis; GE, global epistasis.

values of the pairwise-interaction parameters corresponding to both inter-site and intra-site interactions across the 50 sequences (Fig. 6a).

The results are shown in Figures 6b and 6c. We found that pairwise-interaction models consistently performed better on test data than additive surrogate models inferred in a similar manner, and that both types of models benefited from having a GE nonlinearity (Fig. 6b). The pairwise-interaction models thus provide more accurate approximations of the DNN, suggesting that the pairwise-interaction parameters in these models are likely to be meaningful. Examining the resulting context-averaged pairwise-interaction parameters (Fig. 6c), we observed strong positive intra-site interactions and strong negative inter-site interactions for critical mutations, i.e., mutations away from the preferred nucleotide at any of the 6 most selective positions of the AP-1 motif. Intuitively, the positive interactions between a pair of critical mutations within the same AP-1 site arise because a single critical mutation is sufficient to abrogate AP-1 binding. Thus, a second critical mutation within the same site has no effect, rather than the additional negative effect that would be predicted by an additive model. By contrast, the negative inter-site interactions between critical mutations indicate a degree of redundancy between nearby AP-1 binding sites, since mutations that abrogate both AP-1 sites result in lower activity than is expected from adding the effects of abrogating each AP-1 site individually.⁴⁷

While these averages are informative about overall trends, analyses of specific sequences can provide additional insight. In the analysis above, we observed one sequence that contained 3 putative AP-1 binding sites, with different combinations of these sites displaying either positive or negative epistatic interactions (Fig. 6d). This example shows that context-specific interactions can be complex, and that other factors like binding site orientation and spacing may play an important role.⁴⁸ We note that including the GE nonlinearity in the pairwise-interaction model was essential for identifying these complex interactions. When a linear pairwise-interaction model was used, we instead observed positive interactions of similar magnitude between every pair of sites (Fig. 6e). The reason is that, in the linear pairwise model, the pairwise-interaction parameters are co-opted to describe a global nonlinearity instead of the nucleotide-specific interactions they are intended to model. Specifically, for this genomic sequence, the DNN exhibits a rapidly saturating effect of removing AP-1 binding sites on the functional score (Supplementary Fig. 6), such that disrupting a second binding site has a much smaller effect on DNN predictions than would be predicted based on the effect of disrupting just one of the three sites. We thus see that modeling GE nonlinearities is essential when interpreting genomic DNNs in terms of epistatic interactions between genetic elements.

SQUID supports global DNN interpretations

In the above analyses, SQUID used *in silico* MAVE libraries generated by partially mutagenizing a specific sequence of interest in the genome. As a result, the surrogate models inferred by SQUID provided DNN interpretations that are limited to localized regions of sequence space. In previous work, however, we proposed a DNN interpretation method called global importance analysis (GIA),⁴⁹ in which the DNN is evaluated on completely random sequences containing embedded genetic elements, such as putative TF binding sites. We therefore investigated whether SQUID could provide useful global DNN interpretations using a modified version of the sequence libraries used in GIA.

First we asked whether SQUID could provide global interpretations for individual TFs. For each TF of interest, we generated an *in silico* MAVE library by embedding partially mutagenized versions of the consensus TF binding site within random DNA

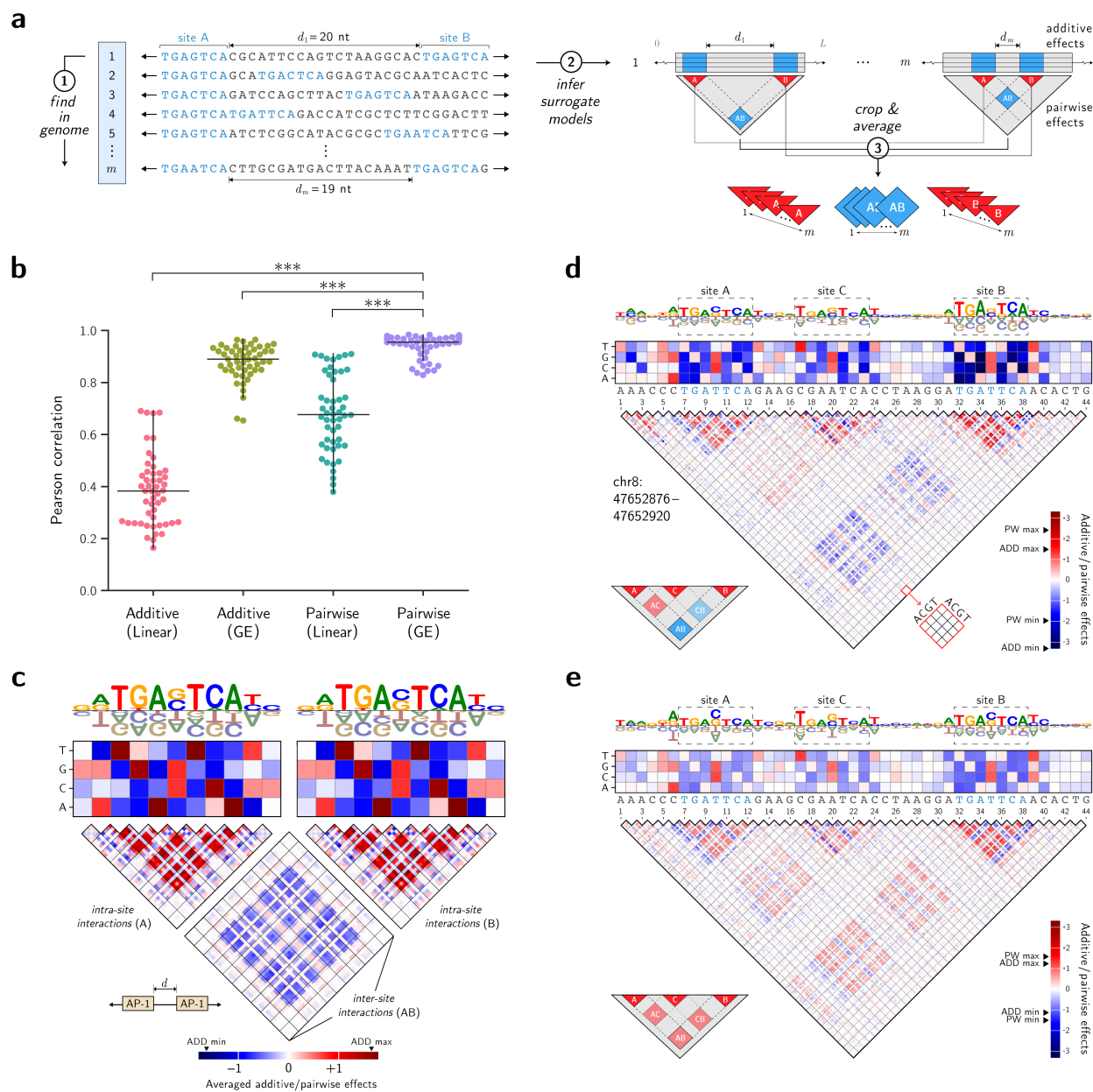


Figure 6. SQUID captures epistatic interactions. **a**, Pipeline for analyzing context-averaged epistatic interactions. The pipeline consists of 3 steps: (1) pairs of consensus binding sites (A and B) are identified in genomic sequences; (2) pairwise-interaction models are inferred for each identified sequence; and (3) the values of surrogate model parameters describing intra-site interactions (A, B) and inter-site interactions (AB) are averaged across sequence contexts. **b**, Performance of surrogate models for 50 genomic sequences having two consensus AP-1 binding sites each. Results are shown for both additive and pairwise-interaction models, each with (GE) and without (linear) a GE nonlinearity. Correlation values were computed between surrogate model predictions and DNN predictions on *in silico* test data. Overlaid lines represent the median, upper, and lower quartiles. P-values were computed using the one-sided Mann-Whitney U test. ***, $p < 0.001$. **c**, Surrogate model parameters quantifying intra-site and inter-site interactions, averaged across the 50 genomic sequence contexts. Pairwise-interaction models having GE nonlinearities were used in this analysis. **d,e**, Parameters of pairwise-interaction models, either with (panel **d**) or without (panel **e**) a GE nonlinearity, determined for a genomic locus having three putative AP-1 binding sites. GE, global epistasis; DNN, deep neural network.

We next investigated whether SQUID could provide global insights into the epistatic interactions between pairs of TF binding sites. We created an *in silico* MAVE library in which partially-mutagenized version of the Nanog and Sox2 consensus binding sites were inserted into random DNA sequences a fixed distance apart (Fig. 7c). We then used SQUID to infer a pairwise-interaction model based on this *in silico* MAVE library. Fig. 7d shows that, similar to the findings in Fig. 6, critical mutations within Nanog and Sox2 binding sites exhibited positive intra-site epistatic interactions and negative inter-site epistatic interactions. We then repeated this analysis for Nanog and Sox2 binding sites separated 0 to 32 nt and observed that the strength of inter-site epistatic interactions varied in a sinusoidal manner consistent with the periodicity of the DNA double helix (Fig. 7e). We conclude that global surrogate modeling with SQUID can provide useful characterizations of epistatic interactions between TFs in a manner that is independent of genomic sequence context.

Discussion

Here we introduced SQUID, a framework for interpreting genomic DNNs. SQUID uses surrogate models to approximate DNN functions in user-defined regions of sequence space. The parameters of these surrogate models can then be mechanistically interpreted: additive surrogate models can be interpreted as attribution maps, and pairwise-interaction surrogate models can be interpreted as quantifying epistatic interactions. Applying SQUID to a variety of genomic DNNs, we observed that the attribution maps obtained by SQUID more robustly identify TF binding motifs and provide better variant-effect predictions than the attribution maps obtained using other DNN interpretability methods. We also observed that SQUID is able to quantify epistatic interactions that are otherwise obscured by global nonlinearities in the DNN.

SQUID works by using the DNN of interest as a forward simulator of MAVE experiments. SQUID then uses the MAVE-NN modeling framework²⁶ to infer surrogate models from the resulting *in silico* data. This surrogate modeling approach has multiple important advantages over standard DNN interpretability methods.

First, SQUID is model agnostic: it does not require access to the parameters or gradients of the DNN. Rather, SQUID simply uses the DNN of interest as a black-box oracle. This allows SQUID to be applied to arbitrary genomic DNNs, regardless of their computational implementation.

Second, SQUID smooths over fluctuations in DNN predictions that are likely to be random rather than functional. This is because the parameters of the surrogate models that SQUID infers are fit to DNN predictions over an extended region of sequence space. Our results showed that this increased smoothness can cause the surrogate models to have improved accuracy relative to the parent DNN.

Third, SQUID leverages domain-specific knowledge to improve the utility of surrogate models. While the use of surrogate models for interpreting DNNs has been established previously in other fields, these applications typically use linear models.³² By using the latent phenotype models supported by MAVE-NN, the surrogate models inferred by SQUID explicitly account for global nonlinearities and heteroscedastic noise, both of which are ubiquitous in functional genomic data and in genomic DNNs. By explicitly accounting for these influences, SQUID is able to remove their confounding effects on the inferred surrogate model parameters.

Finally, SQUID supports a variety of surrogate models. In this work we demonstrated the use of additive and pairwise-interaction models, but SQUID also supports use of surrogate models having user-specified mathematical form. SQUID thus has the ability to infer surrogate models that reflect higher-order epistatic interactions,^{21,22,25,50} or specific biophysical hypotheses (e.g. thermodynamic models^{51–55} or kinetic models^{56–59}).

The biggest drawback of SQUID relative to other DNN interpretability methods is its higher computational demands. This is due to the need for many forward passes through the DNN during simulation of the *in silico* MAVE dataset, as well as the need to fit the surrogate model parameters to these simulated data. We therefore suggest that SQUID is likely to be more useful for in-depth analyses of specific sequences of interest (e.g., disease-associated loci) rather than large-scale genome-wide analyses. In this context, however, the advantages of SQUID described above—especially SQUID's ability to support surrogate models of arbitrary mathematical form—enables new ways of biologically interpreting genomic DNNs.

Methods

The SQUID framework

An overview of SQUID’s workflow is given in Supplementary Figure 7. Briefly, SQUID takes as input a sequence of interest and a specified surrogate model. An *in silico* MAVE dataset is generated with the `InSilicoMAVE` object, with specifications of the mutagenesis strategy given by a `Mutagenizer` object, and processed model predictions given by a `Predictor` object. The *in silico* MAVE dataset is then fit with surrogate models which are defined as objects in `squid/surrogate_zoo.py`.

- **Mutagenizer.** An *in silico* MAVE dataset is generated by sampling a library of sequences using random partial mutagenesis of a sequence of interest. We modulate the size of the sequence-space region from which this library is drawn using two hyperparameters: the sequence that defines the region of interest, which has length L , and the mutation rate r . The resulting number of mutations in each individual sequence is a Poisson distributed random variable having mean Lr . `squid/mutagenizer.py` contains objects that apply the chosen mutagenesis strategy, with the object `RandomMutagenesis` executing the random partial mutagenesis in this study.
- **Predictor.** SQUID currently requires that DNNs provide scalar-valued outputs. However, some genomic DNNs output high-dimensional predicted profiles, not scalar predictions. For example, `ResidualBind-32` predicts 15 chromatin accessibility profiles of length 64, each profile corresponding to a different cell type and each profile element representing a binned position with a resolution of 32 nt. `squid/predictor.py` contains objects, including `ScalarPredictor` and `ProfilePredictor`, that reduce model predictions to scalar values. For profile-based predictions, SQUID also offers an approach to reduce profiles using principal component analysis (PCA). Specifically, profiles are projected onto their first principal component, with sign chosen so that the wild-type sequence has higher-than average score (see Supplementary Fig. 8).
- **InSilicoMAVE.** The object `InSilicoMAVE`, defined in `squid/mave.py`, is a data structure that takes the mutagenizer and predictor objects and generates the *in silico* MAVE dataset for a sequence of interest.
- **Surrogate Zoo.** The `squid/surrogate_zoo.py` module currently offers a linear model (`SurrogateLinear`) as well as models based on MAVE-NN (`SurrogateMAVENN`). In `SurrogateRidgeCV`, the model is computed using the ridge regression from `sklearn` (`linear_model`, `RidgeCV`). `SurrogateMAVENN` supports models based on the built-in modeling capabilities of MAVE-NN. The mathematical forms of these models are described in detail in ref.²⁶

Surrogate models

The surrogate models supported by MAVE-NN for use in SQUID comprise three parts (Fig. 1b): a G-P map, a GE nonlinearity, and a noise model. Here we summarize the mathematical forms of these model components. In what follows, \mathbf{x} represents a sequence of interest, $x_{l:c}$ is a one-hot encoding of \mathbf{x} (i.e., is equal to 1 if \mathbf{x} has character c at position l and is equal to 0 otherwise), and y is the DNN-predicted scalar activity of \mathbf{x} .

- The **G-P map**, $f(\mathbf{x}, \theta)$, which has parameters θ , maps a sequence \mathbf{x} to a latent phenotype ϕ . There are two types of G-P maps: additive G-P maps and pairwise-interaction G-P maps. Additive G-P maps have a constant parameter θ_0 , additive parameters $\theta_{l:c}$, and the mathematical form

$$f_{\text{additive}}(\mathbf{x}, \theta) = \theta_0 + \sum_{l=1}^L \sum_c \theta_{l:c} x_{l:c} .$$

Pairwise-interaction G-P maps have a constant parameter, additive parameters, pairwise-interaction parameters $\theta_{l:c,l':c'}$, and the mathematical form

$$f_{\text{pairwise}}(\mathbf{x}, \theta) = \theta_0 + \sum_{l=1}^L \sum_c \theta_{l:c} x_{l:c} + \sum_{l=1}^{L-1} \sum_{l'}^L \sum_{c,c'} \theta_{l:c,l':c'} x_{l:c} x_{l':c'} .$$

- The **GE nonlinearity**, $g(\phi)$, maps the latent phenotype ϕ to a most-probable scalar DNN prediction \hat{y} . By default, $g(\phi)$ is defined to be an over-parameterized linear combination of hyperbolic tangents. For models “without” a GE nonlinearity, $g(\phi)$ is defined to be a linear function of ϕ .
- The **noise model**, $p(y|\hat{y})$, describes the expected distribution of DNN predictions y about the most-probable prediction \hat{y} . The noise model can be defined using a Gaussian distribution, a Student’s t-distribution, or the skewed t-distributed of Jones and Faddy.²⁹ The shape parameters of this distribution can either be independent of \hat{y} (for homoscedastic noise) or a polynomial function of \hat{y} (for heteroscedastic noise).

Deep learning models

This study used six DNNs: ResidualBind-32,³¹ Basenji-32,³¹ DeepSTARR,³³ Enformer,² BPNet,⁵ and a baseline CNN that predicts ChIP-seq data for the human TF GABPA. Here we briefly describe each DNN and how that DNN was used in our study to compute a prediction y for each sequence x when generating *in silico* MAVE data.

- **ResidualBind-32** predicts ATAC-seq profiles across 15 human cell lines.³¹ ResidualBind-32 takes as input a DNA sequence of length 2048 nt and outputs 15 profiles (one for each cell line) where each profile comprises 64 bins, with each bin spanning 32 nt. The published ResidualBind-32 parameters were used to compute these profiles. In our attribution error analyses, y was set equal to the sum of predicted binned ATAC-seq signals over all 64 bins for the single output channel corresponding to the cell type most associated with the TF of interest (see Supplementary Table 1). In our variant-effect analysis, y was set equal to the sum of predicted binned ATAC-seq signals over all 64 bins, using a profile averaged across all 15 output channels.
- **Basenji-32** predicts ATAC-seq profiles across 15 human cell lines.⁴⁶ The input and output of Basenji-32 is identical to that of ResidualBind-32, and the published Basenji-32 parameters were used to compute predicted ATAC-seq profiles. Analyses performed using Basenji-32 were performed the same way as for ResidualBind-32.
- **DeepSTARR** predicts *Drosophila* enhancer activity as assayed by UMI-STARR-seq.³³ DeepSTARR takes as input a DNA sequence of length 249 nt and outputs two scalar-valued predictions for enhancer activity for developmental (Dev) and housekeeping (Hk) regulatory programs. The published DeepSTARR parameters were used to predict enhancer activity. In each analysis, y was computed using the regulatory program most associated with the TF of interest (see Supplementary Table 1).
- **Enformer** predicts many different types of functional genomic tracks (e.g., ChIP-seq, DNase-seq, ATAC-seq, and CAGE) across the human and mouse genomes.² Enformer takes as input a DNA sequence of length 393,216 nt and (for humans) outputs 5,313 profiles (one for each track) where each profile comprises 128 bins, each bin spanning 32 nt, representing the central 114,688 nt of the input sequence. The published Enformer parameters were used to compute these profiles. In our variant-effect analysis, we used human predictions, where y was computed as in the original study by cropping all 674 “cell-type agnostic” DNase profiles to a 10-bin (1280 nt) region centered about the variant of interest, then summing across bins in the mean cropped profile.
- **BPNet** predicts nucleotide-resolution ChIP-nexus binding profiles for four TFs (Oct4, Sox2, Klf4, and Nanog) in mouse embryonic stem cells.⁵ BPNet takes as input a DNA sequence of length 1000 nt and outputs a 1000-valued positive (+) and negative (-) strand profile for each of the four TFs (8 profiles in total per input). The parameters of BPNet were retrained as specified in the original release,⁶⁰ and the resulting model was confirmed to recapitulate the published model BPNet-OSKN⁶¹ by a visual inspection of attribution maps. Analyses for different prediction tasks used different profiles (see Supplemental Table 2). In all BPNet analyses except those in Figure 7, y was computed using the profile contribution score defined in the original paper.⁵ In the BPNet analysis for Figure 7, y was instead computed using PCA as described above and in Supplementary Figure 8.
- The **baseline CNN** predicts ChIP-seq peaks for the human TF GABPA in GM12878 cells. This model takes as input a DNA sequence of length 200 nt and outputs a single probability. In our analysis of the effects of benign overfitting, y was computed as the logit of the output probability.

The baseline CNN has not been previously published. Briefly, this model was trained to distinguish binary-labeled sequences: ChIP-Seq peaks of GABPA in human GM12878 cells (positive labels) and DNase-seq peaks of GM12878 cells that did not overlap with any GABPA peaks (negative labels). Data was acquired from <https://zenodo.org/record/7011631> (data/GABPA_200.h5); 11,022, 1,574, and 3,150 sequences in the training, validation, and test sets were respectively used. The CNN takes as input a DNA sequence of length 200 nt and outputs a probability. The hidden layers consist of three convolutional blocks each with max pooling (size 4 with stride 4), followed by a fully-connected hidden layer with 128 units and an output layer to a single node with sigmoid activations. The number of filters and the kernel sizes from the first convolutional layer to the third are: (32, 15), (64, 5), and (96, 5). The same padding was used for each convolutional layer. All hidden layers used ReLU activations and dropout,⁶² with rates given in order by 0.1, 0.2, 0.3, and 0.5. The CNN was trained to minimize the binary cross-entropy loss function using Adam⁶³ with default parameters and a batch size of 64. Early stopping was implemented to save model parameters at the epoch that corresponded to the largest area-under the ROC curve (AUROC) on the validation set, yielding the pre-overfitting DNN parameters. From that point, the CNN was trained for an additional 200 epochs, yielding the post-overfitting DNN parameters.

Attribution methods

Our analyses used attribution maps computed using a variety of methods, implemented as follows.

- **In silico mutagenesis (ISM) scores** were computed by evaluating the scalar DNN prediction for every single-nucleotide variant of the sequence of interest.
- **Saliency Maps scores** were computed by evaluating the gradient of the scalar DNN prediction at the sequence of interest with respect to the one-hot encoding of that sequence.
- **DeepSHAP scores** were computed as previously described for DeepSTARR.³³
- **DeepLIFT scores** were computed as previously described for BPNet.⁵
- **SmoothGrad scores** were computed by averaging Saliency Maps over 50 noisy encodings of the sequence of interest. Each noisy encoding was computed by adding Gaussian noise (mean zero, standard deviation 0.25) to each of the $4L$ matrix elements of the one-hot encoding for the sequence of interest.

Standardization of attribution maps

Prior to comparing attribution maps, we standardized these maps to remove non-identifiable and/or non-meaningful degrees of freedom. Using $v_{l:c}$ to denote the attribution map value for character c at position l , this standardization was carried out as follows.

- For **plotting sequence logos** and **computing attribution errors**: Attribution map values were standardized using the transformation $v_{l:c} \rightarrow (v_{l:c} - \bar{v}_l) / \sigma$, where $\bar{v}_l = \frac{1}{4} \sum_{c'} v_{l:c'}$ and $\sigma^2 = \frac{1}{4L} \sum_{c,l} (v_{l:c} - \bar{v}_l)^2$. This transformation is essentially a gradient correction at each position,⁶⁴ followed by a normalization with the square root of the total variation of the attribution scores across the sequence.
- For **variant-effect analysis** (Table 1): Attribution map values were standardized using the transformation $v_{l:c} \rightarrow v_{l:c} - v_{l:s_l^{\text{wt}}}$, where s^{wt} is the wild-type sequence and s_l^{wt} is the character at position l in this sequence.
- For **plotting heatmaps** of additive and pairwise-interaction model parameters (Fig. 6c,d,e and Fig. 7d): Parameters were standardized as in ref.²⁶ using the “empirical gauge”.

Attribution error computations

Attribution errors were computed as follows.

- For **consensus TF binding sites** (Fig. 2 and Fig. 3), we first located all instances of the consensus TF binding site in the genome; the consensus sites used for each TF are listed in Supplemental Table 1. Using a baseline attribution method (Saliency Maps for ResidualBind-32, ISM for DeepSTARR and BPNet), we then ranked and manually pruned these genomic sequences to identify $m = 50$ putative functional and spatially-isolated TF binding sites. For each of the m genomic sequences, we then computed an attribution map spanning the putative TF binding site plus 100 nucleotides on either side (i.e., flanks). We then cropped these m attribution maps to span n_f nucleotides on either side of the consensus TF binding site. For Figure 2 we used $n_f = 15$; for Figure 3 we used $n_f = 0, 5, 10, 20, 30, \dots, 100$. For each of the m cropped attribution maps, the attribution error was defined to be the Euclidean distance between the cropped attribution map and the average of the m cropped attribution maps.
- For **weak TF binding sites** (Fig. 5), we first located all instances of variants of the consensus TF binding site in the genome having up to 2 naturally-occurring mutations. For each group of putative binding sites having 0, 1, or 2 mutations, we then ranked and manually pruned these genomic sequences to identify $m = 50$ putative functional and spatially-isolated TF binding sites as above. For each of the $3m$ genomic sequences, we then computed an attribution map spanning the putative TF binding site plus 100 nucleotides on either side. We then cropped these $3m$ attribution maps to span $n_f = 50$ nucleotides on either side of the consensus TF binding site. For each of the $3m$ cropped attribution maps, the attribution error was defined to be the Euclidean distance between the cropped attribution map and the average of the m cropped attribution maps from the ensemble of binding sites having 0 mutations.

Context-averaged epistatic interactions

To compute the context-averaged epistatic interactions shown in Figure 6c, first, we located all pairs of consensus AP-1 binding sites in the genome separated by no more than 20 nucleotides. We then ranked and pruned these pairs to identify $m = 50$ genomic sequences for further analysis. For each of the m genomic sequences, we used SQUID to infer a pairwise-interaction model spanning the pair of putative sites and 6 nucleotides on either side. Parameter values were standardized by expressing them in the empirical gauge.²⁶ Additive parameters for each site were then cropped and averaged over the m sequences. The intra-site and inter-site pairwise-interaction parameters were similarly cropped and averaged. This process is schematized in Figure 6a.

Variation-effect analysis

The ability of different attribution methods to predict variant effects, quantified in Table 1, was benchmarked as follows.

- The **variant-effect data** used to perform these benchmark studies is MPRA data from ref.⁴⁴. These data comprise measurements for the effects of nearly all SNVs in 15 disease-associated regulatory sequences, with each sequence ranging from 187 nt to 600 nt in length. For each assayed SNV, variant effect was quantified as the difference in measured activity between the variant and wild-type sequences.
- To compute **attribution map predictions of variant effects**, we extracted the genomic sequences centered about each of the 15 assayed regulatory sequences (2,048 nt per sequence for tests of ResidualBind-32 and Basenji-32; 393,216 nt per sequence for tests of Enformer). For each attribution method, the effect of an SNV, from character s_l^{wt} to character c at position l in each regulatory sequence s^{wt} , was quantified as $v_{l:c} - v_{l:s_l^{\text{wt}}}$. When inferring surrogate models using SQUID, we mutagenized the regulatory sequence and 400 nt of flanking sequence on either side.

Data availability

The data used in this paper is available on Zenodo.⁶⁵

Code availability

SQUID is an open-source Python package based on TensorFlow.⁶⁶ SQUID can be installed via pip (<https://pypi.org/project/squid-nn>) or GitHub (<https://github.com/evanseitz/squid-nn>). Documentation for SQUID is provided on ReadTheDocs (<https://squid-nn.readthedocs.io>). The code for performing all analyses in this paper is available on GitHub as well (<https://github.com/evanseitz/squid-manuscript>), and a static snapshot of this code is available on Zenodo.⁶⁵

References

1. Linder, J., Srivastava, D., Yuan, H., Agarwal, V. & Kelley, D. R. Predicting RNA-seq coverage from DNA sequence as a unifying model of gene regulation. *bioRxiv* 2023–08 (2023).
2. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat. Methods* **18**, 1196–1203 (2021).
3. Dudnyk, K., Shi, C. & Zhou, J. Sequence basis of transcription initiation in human genome. *bioRxiv* (2023).
4. Jaganathan, K. *et al.* Predicting splicing from primary sequence with deep learning. *Cell* **176**, 535–548.e24 (2019).
5. Avsec, Ž. *et al.* Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat. Genet.* **53**, 354–366 (2021).
6. Chen, K. M., Wong, A. K., Troyanskaya, O. G. & Zhou, J. A sequence-based global map of regulatory activity for deciphering human genetics. *Nat. Genet.* **54**, 940–949 (2022).
7. Zhou, J. Sequence-based modeling of three-dimensional genome architecture from kilobase to chromosome scale. *Nat. Genet.* **54**, 725–734 (2022).
8. Koo, P. K. & Ploenzke, M. Deep learning for inferring transcription factor binding sites. *Curr. Opin. Syst. Biol.* **19**, 16–23 (2020).
9. Novakovsky, G., Dexter, N., Libbrecht, M. W., Wasserman, W. W. & Mostafavi, S. Obtaining genetics insights from deep learning via explainable artificial intelligence. *Nat. Rev. Genet.* **24**, 125–137 (2022).
10. Han, T., Srinivas, S. & Lakkaraju, H. Which explanation should I choose? A function approximation perspective to characterizing post hoc explanations. *arXiv* (2022).

11. Hooker, S., Erhan, D., Kindermans, P.-J. & Kim, B. A benchmark for interpretability methods in deep neural networks. *Adv. neural information processing systems* **32** (2019).
12. Ancona, M., Ceolini, E., Öztireli, C. & Gross, M. Towards better understanding of gradient-based attribution methods for deep neural networks. *arXiv* (2017).
13. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *Workshop at International Conference on Learning Representations* (2014).
14. Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML'17*, 3145–3153 (JMLR.org, 2017).
15. Zhou, J. & Troyanskaya, O. G. Predicting effects of noncoding variants with deep learning-based sequence model. *Nat. Methods* **12**, 931–934 (2015).
16. Smilkov, D., Thorat, N., Kim, B., Viégas, F. & Wattenberg, M. SmoothGrad: Removing noise by adding noise (2017). [1706.03825](https://arxiv.org/abs/1706.03825).
17. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. *arXiv* (2017).
18. Lundberg, S. M. & Lee, S.-I. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems*, vol. 30 (Curran Associates, Inc., 2017).
19. Starr, T. N. & Thornton, J. W. Epistasis in protein evolution. *Protein science* **25**, 1204–1218 (2016).
20. Weinreich, D. M., Lan, Y., Wylie, C. S. & Heckendorn, R. B. Should evolutionary geneticists worry about higher-order epistasis? *Curr. Opin. Genet. & Dev.* **23**, 700–707 (2013).
21. Aghazadeh, A. *et al.* Epistatic net allows the sparse spectral regularization of deep neural networks for inferring fitness functions. *Nat. communications* **12**, 5225 (2021).
22. Zhou, J. *et al.* Higher-order epistasis and phenotypic prediction. *Proc. Natl. Acad. Sci.* **119** (2022).
23. Domingo, J., Baeza-Centurion, P. & Lehner, B. The causes and consequences of genetic interactions (epistasis). *Annu. Rev. Genomics Hum. Genet.* **20**, 433–460 (2019).
24. Otwinowski, J., McCandlish, D. M. & Plotkin, J. B. Inferring the shape of global epistasis. *Proc. Natl. Acad. Sci.* **115** (2018).
25. Poelwijk, F. J., Krishna, V. & Ranganathan, R. The context-dependence of mutations: A linkage of formalisms. *PLOS Comput. Biol.* **12**, e1004771 (2016).
26. Tareen, A. *et al.* MAVE-NN: learning genotype-phenotype maps from multiplex assays of variant effect. *Genome Biol.* **23** (2022).
27. Tonner, P. D., Pressman, A. & Ross, D. Interpretable modeling of genotype–phenotype landscapes with state-of-the-art predictive power. *Proc. Natl. Acad. Sci.* **119**, e2114021119 (2022).
28. Kinney, J. B., Murugan, A., Callan Jr, C. G. & Cox, E. C. Using deep sequencing to characterize the biophysical mechanism of a transcriptional regulatory sequence. *Proc. Natl. Acad. Sci.* **107**, 9158–9163 (2010).
29. Jones, M. C. & Faddy, M. J. A skew extension of the t-distribution, with applications. *J. Royal Stat. Soc. Ser. B: Stat. Methodol.* **65**, 159–174 (2003).
30. Tareen, A. & Kinney, J. B. Logomaker: Beautiful sequence logos in Python. *Bioinformatics* **36**, 2272–2274 (2019).
31. Toneyan, S., Tang, Z. & Koo, P. Evaluating deep learning for predicting epigenomic profiles. *Nat. Mach. Intell.* **4**, 1088–1100 (2022).
32. Ribeiro, M. T., Singh, S. & Guestrin, C. "Why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (ACM, 2016).
33. de Almeida, B. P., Reiter, F., Pagani, M. & Stark, A. DeepSTARR predicts enhancer activity from DNA sequence and enables the de novo design of synthetic enhancers. *Nat. Genet.* **54**, 613–624 (2022).
34. Bartlett, P. L., Long, P. M., Lugosi, G. & Tsigler, A. Benign overfitting in linear regression. *Proc. Natl. Acad. Sci.* **117**, 30063–30070 (2020).
35. Chatterji, N. S. & Long, P. M. Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. *The J. Mach. Learn. Res.* **22**, 5721–5750 (2021).

36. Wang, Z. *et al.* Smoothed geometry for robust attribution. *Adv. neural information processing systems* **33**, 13623–13634 (2020).
37. Alvarez-Melis, D. & Jaakkola, T. S. Towards robust interpretability with self-explaining neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS'18*, 7786–7795 (Curran Associates Inc., Red Hook, NY, USA, 2018).
38. Majdandzic, A. *et al.* Selecting deep neural networks that yield consistent attribution-based interpretations for genomics. In *Machine Learning in Computational Biology*, 131–149 (PMLR, 2022).
39. Zhang, C., Bengio, S., Hardt, M., Recht, B. & Vinyals, O. Understanding deep learning requires rethinking generalization (2017). [1611.03530](https://arxiv.org/abs/1611.03530).
40. Papagianni, A. *et al.* Capicua controls Toll/IL-1 signaling targets independently of RTK regulation. *Proc. Natl. Acad. Sci.* **115**, 1807–1812 (2018).
41. Crocker, J. *et al.* Low affinity binding site clusters confer hox specificity and regulatory robustness. *Cell* **160**, 191–203 (2015).
42. Farley, E. K. *et al.* Suboptimization of developmental enhancers. *Science* **350**, 325–328 (2015).
43. Castro-Mondragon, J. A. *et al.* JASPAR 2022: The 9th release of the open-access database of transcription factor binding profiles. *Nucleic Acids Res.* **50**, D165–D173 (2021).
44. Kircher, M. *et al.* Saturation mutagenesis of twenty disease-associated regulatory elements at single base-pair resolution. *Nat. Commun.* **10** (2019).
45. Shigaki, D. *et al.* Integration of multiple epigenomic marks improves prediction of variant impact in saturation mutagenesis reporter assay. *Hum. mutation* **40**, 1280–1291 (2019).
46. Kelley, D. R. *et al.* Sequential regulatory activity prediction across chromosomes with convolutional neural networks. *Genome Res.* **28**, 739–750 (2018).
47. Kim, S. & Wsocka, J. Deciphering the multi-scale, quantitative cis-regulatory code. *Mol. Cell* (2023).
48. Georgakopoulos-Soares, I. *et al.* Transcription factor binding site orientation and order are major drivers of gene regulatory activity. *Nat. Commun.* **14**, 2333 (2023).
49. Koo, P. K., Majdandzic, A., Ploenzke, M., Anand, P. & Paul, S. B. Global importance analysis: An interpretability method to quantify importance of genomic features in deep neural networks. *PLoS Comput. Biol.* **17**, e1008925 (2021).
50. Weinreich, D. M., Lan, Y., Jaffe, J. & Heckendorn, R. B. The influence of higher-order epistasis on biological fitness landscape topography. *J. Stat. Phys.* **172**, 208–225 (2018).
51. Ackers, G. K., Johnson, A. D. & Shea, M. A. Quantitative model for gene regulation by lambda phage repressor. *Proc. Natl. Acad. Sci.* **79**, 1129–1133 (1982).
52. Bintu, L. *et al.* Transcriptional regulation by the numbers: Models. *Curr. Opin. Genet. & Dev.* **15**, 116–124 (2005).
53. Segal, E. & Widom, J. From DNA sequence to transcriptional behaviour: A quantitative approach. *Nat. Rev. Genet.* **10**, 443–456 (2009).
54. Sherman, M. S. & Cohen, B. A. Thermodynamic state ensemble models of cis-regulation. *PLoS Comput. Biol.* **8**, e1002407 (2012).
55. Faure, A. J. *et al.* Mapping the energetic and allosteric landscapes of protein binding domains. *Nature* **604**, 175–183 (2022).
56. Tareen, A. & Kinney, J. B. Biophysical models of cis-regulation as interpretable neural networks. (2019).
57. Estrada, J., Wong, F., DePace, A. & Gunawardena, J. Information integration and energy expenditure in gene regulation. *Cell* **166**, 234–244 (2016).
58. Scholes, C., DePace, A. H. & Sánchez, Á. Combinatorial gene regulation through kinetic control of the transcription cycle. *Cell Syst.* **4**, 97–108.e9 (2017).
59. Park, J. *et al.* Dissecting the sharp response of a canonical developmental enhancer reveals multiple sources of cooperativity. *eLife* **8** (2019).
60. Žiga Avsec & Weilert, M. kundajelab/bpnet-manuscript: Publication release (2020). Code available from zenodo.org/records/4294814.

61. Avsec, Z. *et al.* The Kipoi repository accelerates community exchange and reuse of predictive models for genomics. *Nat. Biotechnol.* 1 (2019).
62. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *The J. Mach. Learn. Res.* **15**, 1929–1958 (2014).
63. Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv* (2014).
64. Majdandzic, A., Rajesh, C. & Koo, P. K. Correcting gradient-based interpretations of deep neural networks for genomics. *Genome Biol.* **24**, 1–13 (2023).
65. Seitz, E. evanseitz/squid-manuscript: SQUID manuscript workflow with outputs, DOI: [10.5281/zenodo.10047748](https://doi.org/10.5281/zenodo.10047748) (2023). Code and data available from zenodo.org/records/10047748.
66. Abadi, M. *et al.* TensorFlow: Large-scale machine learning on heterogeneous systems (2015). Software available from tensorflow.org.

Acknowledgements

We thank Amber Tang, Shushan Toneyan, Mahdi Kooshkbaghi, Chandana Rajesh, Jakub Kaczmarzyk, and Carlos Martí for helpful discussions. This work was supported in part by: the Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory; NIH grants R01HG012131 (PKK, ESS, JBK, DMM), R01HG011787 (JBK, ESS, DMM), R01GM149921 (PKK), R35GM133777 (JBK), and R35GM133613 (DMM); and an Alfred P. Sloan Foundation Research Fellowship (DMM). Computations were performed using equipment supported by NIH grant S10OD028632.

Author contributions

ESS, DMM, JBK, and PKK conceived of the study. EES wrote the software and performed the analysis. ESS designed the analysis with help from DMM, JBK, and PKK. JBK and PKK supervised the study. ESS, DMM, JBK, and PKK wrote the manuscript.

Competing interests

The authors declare no competing interests.