Contrastive losses as generalized models of global epistasis

David H. Brookes * Dyno Therapeutics

Jakub Otwinowski Dyno Therapeutics **Sam Sinai**Dyno Therapeutics

Abstract

Fitness functions map large combinatorial spaces of biological sequences to properties of interest. Inferring these multimodal functions from experimental data is a central task in modern protein engineering. Global epistasis models are an effective and physically-grounded class of models for estimating fitness functions from observed data. These models assume that a sparse latent function is transformed by a monotonic nonlinearity to emit measurable fitness. Here we demonstrate that minimizing supervised contrastive loss functions, such as the Bradley-Terry loss, is a simple and flexible technique for extracting the sparse latent function implied by global epistasis. We argue by way of a fitness-epistasis uncertainty principle that the nonlinearities in global epistasis models can produce observed fitness functions that do not admit sparse representations, and thus may be inefficient to learn from observations when using a Mean Squared Error (MSE) loss (a common practice). We show that contrastive losses are able to accurately estimate a ranking function from limited data even in regimes where MSE is ineffective and validate the practical utility of this insight by demonstrating that contrastive loss functions result in consistently improved performance on empirical benchmark tasks.

1 Introduction

A fitness function maps biological sequences to relevant scalar properties of the sequences, such as binding affinity to a target molecule, or fluorescent brightness. Biological sequences span combinatorial spaces and fitness functions are typically multi-peaked, due to interactions between positions in a sequence. Learning fitness functions from limited experimental data (often a minute fraction of the possible space) can be a difficult task but allows one to predict properties of sequences. These predictions can help identify promising new sequences for experimentation [42] or to guide the search for optimal sequences [6, 8].

Even in the case of where experimental measurements are available for every possible sequence in a sequence space, inferring a model of the fitness function can be valuable for understanding the factors that impact sequences' fitness [16] or how evolution might progress over the fitness landscape [41].

Numerous methods have been developed to estimate fitness functions from experimental data, including classical machine learning techniques [43], deep learning approaches [18], and semi-supervised methods [21]. Additionally, there are many methods that incorporate biological assumptions into the modeling process, such as parameterized biophysical models [29], non-parametric techniques [44, 45], and methods for spectral regularization of neural networks [1]. These latter approaches largely focus on accurately modeling the influence of "epistasis" on fitness functions, which refers to statistical or physical interactions between genetic elements, typically either amino-acids in a protein sequence or genes in a genome.

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

^{*}david.brookes@dynotx.com

"Local" epistasis refers to interactions between a limited number of specific positions in a sequence, and is often modeled using interaction terms in a linear model of a fitness function [30]. "Global" epistasis, on the other hand, refers to the presence of nonlinear relationships that affect the fitness of sequences in a nonspecific manner. A model of global epistasis typically assumes a simple latent fitness function is transformed by a monotonically increasing nonlinearity to produce observed fitness data [36, 31, 39, 34, 37]. Typically, these models assume a particular parametric form of the latent fitness function and nonlinearity, and fit the parameters of both simultaneously. It is most common to assume that the underlying fitness function includes only additive (non-interacting) effects [31], though pairwise interaction effects have been added in some models [39].

Despite their relative simplicity, global epistasis models have been found to be effective at modeling experimentally observed fitness functions [37, 33, 34]. Furthermore, global epistasis is not just a useful modeling choice, but a physical phenomenon that can result from features of a system's dynamics [22] or the environmental conditions in which a fitness function is measured [31]. Therefore, even if one does not use a standard global epistasis model, it is still important to consider the effects of global epistasis when modeling fitness functions.

Due to the monotonicity of the nonlinearity in global epistasis models, the latent fitness function in these models can be interpreted as a parsimonious ranking function for sequences. Herein we show that fitting a model to observed fitness data by minimizing a supervised contrastive, or ranking, loss is a simple and effective method for extracting such a ranking function. In particular, we focus on the Bradley-Terry loss [5], which has been widely used for learning-to-rank tasks [9], and more recently for ordering the latent space of a generative model for protein sequences [12]. Minimizing this loss provides a technique for modeling global epistasis that requires no assumptions on the form of the latent fitness functions or nonlinearity beyond monotonicity, and can easily be applied to any set of observed fitness data.

Further, we use an entropic uncertainty principle to show that global epistasis can result in observed fitness functions that cannot be represented using a sparse set of epistatic interactions. In particular, this uncertainty principle shows that a fitness function that is sufficiently concentrated in the fitness domain—meaning that a small number of sequences have fitness values with relatively large magnitudes—can not be concentrated in the Graph Fourier bases that represent fitness functions in terms of local epistatic interactions [38, 40, 7]. We show that global epistasis nonlinearities tend to concentrate observed fitness functions in the fitness domain, thus preventing a sparse representation in the epistatic domain. This insight has the implication that observed fitness functions that have been affected by global epistasis may be difficult to estimate with undersampled training data and a Mean Squared Error (MSE) loss. We hypothesize that estimating the latent ranking fitness function using a contrastive loss can be done more data-efficiently than estimating the observed fitness function using MSE, and conduct simulations that support this hypothesis. Additionally, we demonstrate the practical importance of these insights by showing that models trained with the Bradley-Terry loss outperform those trained with MSE loss on nearly all empirical benchmark tasks [14].

2 Background

2.1 Fitness functions and the Graph Fourier transform

A fitness function $f: \mathcal{S} \to \mathbb{R}$ maps a space of sequences \mathcal{S} to a scalar property of interest. In the case where \mathcal{S} contains all combinations of elements from an alphabet of size q at L sequence positions, then the fitness function can be represented exactly in terms of increasing orders of local epistatic interactions. For binary sequences (q=2), this representation takes the form:

$$f(x) = \beta_0 + \sum_{i=1}^{L} \beta_i x_i + \sum_{ij} \beta_{ij} x_i x_j + \sum_{ijk} \beta_{ijk} x_i x_j x_k + ...,$$

where $x_i \in \{-1, 1\}$ represent elements in the sequence and each term in the expansion represents a (local) epistatic interaction with weight $\beta_{\{i\}}$, with the expansion continuing up to L^{th} order terms. Analogous representations can be constructed for sequences with any size alphabet q using Graph Fourier bases[38, 40, 7]. These representations can be compactly expressed as:

$$f = \Phi \beta, \tag{1}$$

where f is a length q^L vector containing the fitness values of every sequence in \mathcal{S} , Φ is a $q^L \times q^L$ orthogonal matrix representing the Graph Fourier basis, and β is a length q^L vector containing the weights corresponding to all possible epistatic interactions. We refer to f and β as representing the fitness function in the fitness domain and the epistatic domain, respectively. Note that we may apply the inverse transformation of Eq. 1 to any complete observed fitness function, g to calculate the epistatic representation of the observed data, $g_g = \Phi^T g$. Similarly, if \hat{f} contains the predictions of a fitness model for every sequence in a sequence space, then $\hat{\beta} = \Phi^T \hat{f}$ is the epistatic representation of the model.

A fitness function is considered sparse, or concentrated, in the epistatic domain when β contains a relatively small number of elements with large magnitudes, and many elements equal to zero or with small magnitudes. In what follows, we may refer to a fitness function that is sparse in the epistatic domain as simply being a "sparse fitness function". A number of experimentally-determined fitness functions have been observed to be sparse in the epistatic domain [32, 17, 7]. Crucially, the sparsity of a fitness function in the epistatic domain determines how many measurements are required to estimate the fitness function using Compressed Sensing techniques that minimize a MSE loss function [7]. Herein we consider the effect that global epistasis has on a sparse fitness function. In particular, we argue that global epistasis results in observed fitness functions that are dense in the epistatic domain, and thus require a large amount of data to accurately estimate by minimizing a MSE loss function. However, in these cases, there may be a sparse ranking function that can be efficiently extracted by minimizing a contrastive loss function.

2.2 Global epistasis models

A model of global epistasis assumes that noiseless fitness measurements are generated according to the model:

$$y = g\left(f(\boldsymbol{x})\right),\tag{2}$$

where f is a latent fitness function, g is a monotonically increasing nonlinear function. In most cases, f is assumed to include only first or second order epistatic terms and the nonlinearity is explicitly parameterized using, for example, spline functions [31] or sums of hyperbolic tangents [39]. The restriction that f includes only low-order terms is somewhat arbitrary, as higher-order local epistatic effects have been observed in fitness data (see, e.g., Wu et al. [41]). In general we may consider f to be any fitness function that is sparse in the epistatic domain, and global epistasis then refers to the transformation of a sparse fitness function by a monotonically-increasing nonlinearity.

Global epistasis models of the form of Eq. 2 have proved effective at capturing the variation observed in empirical fitness data [26, 37, 31, 33, 34, 39], suggesting that global epistasis is a common feature of natural fitness functions. Further, it has been shown that global epistasis results from first-principles physical considerations that are common in many biological systems. In particular, Husain and Murugan [22] show that global epistasis arises when the physical dynamics of a system is dominated by slow, collective modes of motion, which is often the case for protein dynamics. Aside from intrinsic/endogenous sources, the process of measuring fitness can also introduce nonlinear effects that are dependent on the experiment and not on the underlying fitness function. For example, fitness data is often left-censored, as many sequence have fitness that falls below the detection threshold of an assay. Finally, global diminishing-returns epistatic patterns have been observed widely in both single and multi-gene settings where the interactions are among genes rather than within a gene [26, 34, 4].

Together, these results indicate that global epistasis is an effect that can be expected in empirically-observed fitness functions. In what follows, we argue that global epistasis is detrimental to effective modeling of fitness functions using standard techniques. In particular, global epistasis manifests itself by producing observed data that is dense in the epistatic domain. In other words, when observed fitness data is produced through Eq. 2 then the epistatic representation of this fitness function (calculated through application of Eq. 1), is not sparse. Further we argue that this effect of global epistasis makes it to difficult to model such observed data by minimizing standard MSE loss functions with a fixed amount of data. Further, we argue that fitting fitness models aimed at extracting the latent fitness function from observed data is a more efficient use of observed data that results in improved predictive performance (in the ranking sense).

While the models of global epistasis described thus far could be used for this purpose, they have the drawback that they assume a constrained form of both g and f, which enforces inductive biases that may affect predictive performance. Here we propose a flexible alternative to modeling global epistasis that makes no assumptions on the form of f or g. In particular, we interpret the latent fitness function f as a parsimonious ranking function for sequences, and the problem of modeling global epistasis as recovering this ranking function. A natural method to achieve this goal is to fit a model of f to the observed data by minimizing a contrastive, or ranking, loss function. These loss functions are designed to learn a ranking function and, as we will show, are able to recover a sparse fitness function that has been transformed by global epistasis to produce observed data. An advantage of this approach to modeling global epistasis is that the nonlinearity g is modeled non-parametrically, and is free to take any form, while the latent fitness function can be modeled by any parametric model, for example, convolutional neural networks (CNNs) or fine-tuned language models, which have been found to perform well in fitness prediction tasks [14]. An accurate ranking model also enables effective optimization, as implied by the results in Chan et al. [12].

2.3 Contrastive losses

Contrastive losses broadly refer to loss functions that compare multiple outputs of a model and encourage those outputs to be ordered according to some criteria. In our case, we desire a loss function that encourages model outputs to be ranked according to observed fitness values. An example of such a loss function is the Bradley-Terry (BT) loss [5, 9], which has the form:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i,j:y_i > y_j} \log \left[1 + e^{-(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - f_{\boldsymbol{\theta}}(\boldsymbol{x}_j))} \right], \tag{3}$$

where f_{θ} is a model with parameters θ , x_i are model inputs and y_i are the corresponding labels of those inputs. This loss compares every pair of data points and encourages the model output $f_{\theta}(x_i)$ to be greater than $f_{\theta}(x_j)$ whenever $y_i > y_j$; in other words, it encourages the model outputs to be ranked according to their labels. A number of distinct but similar loss functions have been proposed in the learning-to-rank literature [13] and also for metric learning [19]. An example is the Margin ranking loss [20], which replaces the logistic function in the sum of Eq. 3 with a hinge function. In our experiments, we focus on the BT loss of Eq. 3 as we found it typically results in superior predictive performance

The BT loss was recently used by Chan et al. [12] to order the latent space of a generative model for protein sequences such that certain regions of the latent space corresponding to sequences with higher observed fitness values. In this case, the BT loss is used in conjunction with standard generative modeling losses. In contrast, here we analyze the use of the BT loss alone in order to learn a ranking function for sequences given corresponding observed fitness values.

A key feature of the contrastive loss in Eq. 3 is that it only uses information about the ranking of observed labels, rather than the numerical values of the labels. Thus, the loss is unchanged when the observed values are transformed by a monotonic nonlinearity. We will show that this feature allows this loss to recover a sparse latent fitness function from observed data that has been affected by global epistasis, and enables more data-efficient learning of fitness functions compared to a MSE loss.

3 Results

Our results are aimed at demonstrating three properties of contrastive losses. First, we show that given complete, noiseless fitness data (i.e. noiseless fitness values associated with every sequence in the sequence space) that has been affected by global epistasis, minimizing the BT loss enables a model to nearly exactly recover the sparse latent fitness function f. Next, we consider the case of incomplete data, where the aim is to predict the relative fitness of unobserved sequences. In this regime, we find through simulation that minimizing the BT loss enables models to achieve better predictive performance then minimizing the MSE loss when the observed data has been affected with global epistasis. We argue by way of a fitness-epistasis uncertainty principle that this is due to the fact that nonlinearities tend to produce fitness functions that do not admit a sparse representation in the epistatic domain, and thus require more data to learn with MSE loss. Finally, we demonstrate the practical significance of these insights by showing that minimizing the BT loss results in improved performance over MSE loss in nearly all tasks in empirical benchmarks [14].

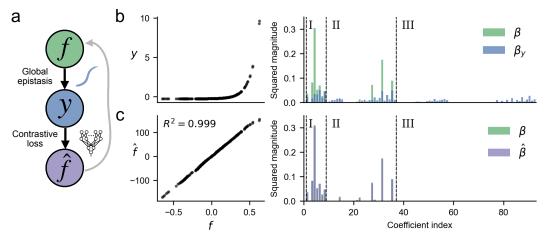


Figure 1: Recovery of latent fitness function from complete fitness data by minimizing Bradley-Terry loss. (a) Schematic of simulation. (b) Comparison between latent (f) and observed (y) fitness functions in fitness (left) and epistatic (right) domains. The latent fitness function is sampled from the NK model with L=8 and K=2 and the global epistasis function is $g(f)=\exp(10\cdot f)$. Each point in the scatter plot represents the fitness of a sequence, while each bar in the bar plot (right) represents the squared magnitude of an epistatic interaction (normalized such that all squared magnitudes sum to 1), with roman numerals indicating the order of interaction. Only epistatic interactions up to order 3 are shown. The right plot demonstrates that global epistasis produces a dense representation in the epistatic domain compared to the representation of the latent fitness in the epistatic domain. (c) Comparison between latent and estimated (\hat{f}) fitness functions in fitness and epistatic domains.

3.1 Recovery from complete data

We first consider the case of "complete" data, where fitness measurements are available for every sequence in the sequence space. The aim of our task in this case is to recover a sparse latent fitness function when the observed measurements have been transformed by an arbitrary monotonic nonlinearity. In particular, we sample a sparse fitness function f from the NK model [24], a popular model of fitness functions that has been shown to recapitulate the sparsity properties of some empirical fitness functions [7]. The NK model has three parameters: L, the length of the sequences, q, the size of the alphabet for sequence elements, and K, the maximum order of (local) epistatic interactions in the fitness function. Roughly, the model randomly assigns K-1 interacting positions to each position in the sequence, resulting in a sparse set of interactions in the epistatic domain. The weights of each of the assigned interactions are then drawn from a independent unit normal distributions.

We then transform the sampled fitness function f with a monotonic nonlinearity g to produce a set of complete observed data, $y_i = g(f(\boldsymbol{x}_i))$ for all $\boldsymbol{x}_i \in \mathcal{S}$. The goal of the task is then to recover the function f given all (\boldsymbol{x}_i, y_i) pairs. In order to do so, we model f using a two layer fully connected neural network and fit the parameters of this model by performing stochastic gradient descent (SGD) on the BT loss, using the Spearman correlation between model predictions and the y_i values to determine convergence of the optimization. We then compare the resulting model, \hat{f} , to the latent fitness function f in both the fitness and epistatic domains, using the forward and inverse transformation of Eq. 1 to convert between the two domains.

Fig. 1 shows the results of one of these tests. In this case, we used an exponential function to represent the global epistasis nonlinearity. The exponential function exaggerates the effects of global epistasis in the epistatic domain and thus better illustrates the usefulness of contrastive losses, although the nonlinearities in empirical fitness functions tend to have a more sigmoidal shape [31]. Fig. 1b shows that the global epistasis nonlinearity substantially alters the representations of the observed data y in both the fitness and epistatic domains, as compared to the latent fitness function f. Nonetheless, Fig. 1c demonstrates that the model fitness function \hat{f} created by minimizing the BT loss is able to nearly perfectly recover the sparse latent fitness function (where recovery is defined as being equivalent up to an affine transformation). This is a somewhat surprising result, as there are many fitness functions that correctly rank the fitness of sequences, and it is not immediately clear why minimizing the BT loss produces this particular sparse latent fitness function. However, this example makes clear that fitting a model by minimizing the BT loss can be an effective strategy for recovering a sparse latent fitness function from observed data that has been affected by global epistasis. Similar results from

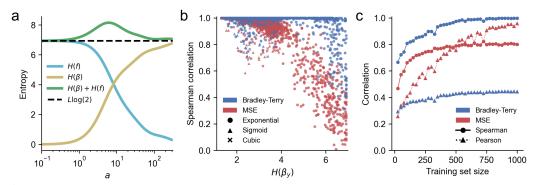


Figure 2: (a) Demonstration of the fitness-epistasis uncertainty principle for a latent fitness function transformed by $g(f) = \exp(a \cdot f)$ for various settings of a. The dashed black line indicates the lower bound on the sum of the entropies of the fitness and epistatic representations of the fitness function (b) Test-set Spearman correlation for models trained with MSE and BT losses on incomplete fitness data transformed by various nonlinearities, compared to the entropy of the fitness function in the epistatic domain. Each point corresponds to a model trained on 256 randomly sampled training points from an L=10, K=2 latent fitness function which was then transformed by a nonlinearity. (c) Convergence of models fit with BT and MSE losses to observed data generated by transforming an L=10, K=2 latent fitness function by $g(f)=\exp(10\cdot f)$. Each point represents the mean test set correlation over 200 training set replicates.

additional examples of this task using different nonlinearities and latent fitness functions are shown in Appendix B.1.

3.2 Fitness-epistasis uncertainty principle

Next, we consider the case where fitness data is incomplete. Our aim is to understand how models trained with the BT loss compare to those trained with MSE loss at predicting the relative fitness of unobserved sequence using different amounts of subsampled training data. We take a signal processing perspective on this problem, and consider how the density of a fitness function in the epistatic domain affects the ability of a model to accurately estimate the fitness function given incomplete data. In particular, we demonstrate that global epistasis tends to increase the density of fitness functions in the epistatic domain, and use an analogy to Compressive Sensing (CS) to hypothesize that more data is required to effectively estimate these fitness functions when using an MSE loss [7]. In order to support this claim, we first examine the effects of global epistasis on the epistatic domain of fitness functions.

Fig. 1b provides an example where a monotonic nonlinearity applied to a sparse fitness increases the density of the the fitness function in the epistatic domain. In particular, we see that many "spurious" local epistatic interactions must appear in order to represent the nonlinearity (e.g. interactions of order 3, when we used an NK model with K=2). This effect can be observed for many different shapes of nonlinearities [36, 3]. We can understand this effect more generally using uncertainty principles, which roughly show that a function cannot be concentrated on a small number of values in two different representations. In particular, we consider the discrete entropic uncertainty principle proved by Dembo et al. [15]. When applied to the transformation in Eq. 1, this uncertainty principle states:

$$H(\mathbf{f}) + H(\boldsymbol{\beta}) \ge -2L\log m,$$
 (4)

where $H(x) = -\sum_i \frac{x_i^2}{||x||^2} \log \frac{x_i^2}{||x||^2}$ is the entropy of the normalized squared magnitudes of a vector and $m = 1/\sqrt{q}$ when q = 2, $m = 1/(q - \sqrt{q})$ when q = 3 and $m = 1 - 1/(q - \sqrt{q})$ otherwise. Low entropy indicates that a vector is concentrated on a small set of elements. Thus, the fitness-epistasis uncertainty principle of Eq. 4 shows that fitness functions cannot be concentrated in both the fitness and epistatic domains. A sparse fitness function (in the epistatic domain) must therefore be "spread out" (i.e. dense) in the fitness domain, and vice-versa.

The importance of this result for understanding global epistasis is that applying a nonlinearity to a dense vector will often have the effect of concentrating the vector on a smaller number of values. This can most easily be seen for convex nonlinearities like the exponential shown in Fig. 1a, but is also true of many other nonlinearities (see Theorem 1 in Appendix C for a sufficient condition for a nonlinearity to decrease entropy in the fitness domain and Appendix D for additional experiments

demonstrating the uncertainty principle). If this concentration in the fitness domain is sufficiently extreme, then the epistatic representation of the fitness function, β , must be dense in order to satisfy Eq. 4. Fig. 2a demonstrates the uncertainty principle by showing how the entropy in the fitness and epistatic domains decrease and increase, respectively, as more extreme nonlinearities are applied to a sparse latent fitness function.

The uncertainty principle quantifies how global epistasis affects a fitness function by preventing a sparse representation in the epistatic domain. From a CS perspective, this has direct implications for modeling the fitness function from incomplete data. In particular, if we were to model the fitness function using CS techniques such as LASSO regression with the Graph Fourier basis as the representation, then it is well known that the number of noiseless data points required to perfectly estimate the function scales as $\mathcal{O}(S \log N)$ where S is the sparsity of the signal in a chosen representation and N is the total size of the signal in the representation [11]. Therefore, when using these techniques, fitness functions affected by global epistasis will require more data to effectively model. Notably, the techniques for which these scaling laws apply minimize a MSE loss functions as part of the estimation procedure. Although these scaling laws only strictly apply to CS modeling techniques, we hypothesize that similar principles may apply when using neural network models and SGD training procedures. In particular, our results suggest that more data is required to fit neural networks to fitness functions with dense epistatic representations using the MSE loss. In the next section we present the results of simulations that support this hypothesis by showing that the entropy of the epistatic representation is negatively correlated with the predictive performance of models trained with an MSE loss on a fixed amount of fitness data. Further, these simulations show that models trained with the BT loss are robust to the dense epistatic representations produced global epistasis, and converge faster to maximum predictive performance as they are provided more fitness data compared to models trained with an MSE loss.

3.3 Simulations with incomplete data

We next present simulation results aimed at showing that global epistasis adversely effects the ability of models to effectively learn fitness functions from incomplete data when trained with MSE loss and that models trained with BT loss are more robust to the effects of global epistasis.

In our first set of simulations, we tested the ability models to estimate a fitness function of L=10binary sequences given one quarter of the fitness measurements (256 measurements out of total of $2^{10} = 1024$ sequences in the sequence space). In each simulation, we (i) sampled a sparse latent fitness function from the NK model, (ii) produced an observed fitness function by applying one of three nonlinearities to the latent fitness function: exponential, $g(f) = \exp(a \cdot f)$, sigmoid, $g(f) = (1 + e^{-a \cdot f})^{-1}$, or a cubic polynomial $g(f) = x^3 + ax$ with settings of the parameter a that ensured the nonlinearity was monotonically increasing, (iii) sampled 256 sequence/fitness pairs uniformly at random from the observed fitness function to be used as training data, and (iv) trained models with this data by performing SGD on the MSE and BT losses. We ran this simulation for 50 replicates of each of 20 settings of the a parameter for each of the three nonlinearities. In every case, the models were fully-connected neural networks with two hidden layers and the optimization was terminated using early stopping with 20 percent of the training data used as a validation set. After training, we measured the extent to which the models estimated the fitness function by calculating Spearman correlation between the model predictions and true fitness values on all sequences in the sequence space. Spearman correlation is commonly used to benchmark fitness prediction methods [14, 27].

The results of each of these simulations are shown in Fig. 2b. We see that the predictive performance of models trained with the MSE loss degrades as the entropy of the fitness function in the epistatic domain increases, regardless of the type of nonlinearity that is applied to the latent fitness function. This is in contrast to the models trained with the BT loss, which often achieve nearly perfect estimation of the fitness function even when the entropy of the fitness function in the epistatic domain approaches its maximum possible value of $L \log 2$. This demonstrates the key result that the MSE loss is sensitive to the density of the epistatic representation resulting from global epistasis (as implied by the analogy to CS), while the BT loss is robust to these effects. We additionally find that these results are maintained when the degree of epistatic interactions in the latent fitness function is changed (Appendix E.1) and when noise is added to the observed fitness functions (Appendix E.2).

Next, we tested how training set size effects the predictive performance of models trained with MSE and BT losses on a fitness function affected by global epistasis. In order to do so, we sampled a single

		SPEARMAN		TOP 10% RECALL	
Data set	SPLIT	MSE Loss	BRADLEY-TERRY	MSE Loss	BRADLEY-TERRY
GB1	1-VS-REST	0.133 ± 0.150	0.091 ± 0.093	0.097 ± 0.030	0.138 ± 0.051
	2-VS-REST	0.564 ± 0.026	$\boldsymbol{0.607 \pm 0.009}$	0.250 ± 0.030	$\textbf{0.282} \pm \textbf{0.008}$
	3-VS-REST*	0.814 ± 0.049	$\textbf{0.880} \pm \textbf{0.003}$	0.539 ± 0.084	$\textbf{0.664} \pm \textbf{0.014}$
	Low-vs-High*	0.499 ± 0.010	$\textbf{0.567} \pm \textbf{0.013}$	0.381 ± 0.028	$\textbf{0.443} \pm \textbf{0.024}$
	SAMPLED	0.930 ± 0.002	$\textbf{0.951} \pm \textbf{0.002}$	0.823 ± 0.009	0.816 ± 0.010
AAV	Mut-Des	0.751 ± 0.006	0.757 ± 0.007	0.288 ± 0.004	0.307 ± 0.005
	DES-MUT	0.806 ± 0.006	$\textbf{0.832} \pm \textbf{0.002}$	0.318 ± 0.013	$\textbf{0.387} \pm \textbf{0.008}$
	1-VS-REST*	0.335 ± 0.117	$\textbf{0.485} \pm \textbf{0.078}$	0.052 ± 0.053	$\textbf{0.143} \pm \textbf{0.049}$
	2-vs-rest	0.748 ± 0.010	$\boldsymbol{0.798 \pm 0.003}$	$\textbf{0.490} \pm \textbf{0.011}$	0.457 ± 0.010
	7-vs-rest	0.732 ± 0.003	$\textbf{0.742} \pm \textbf{0.003}$	0.694 ± 0.006	0.695 ± 0.006
	Low-vs-High	0.401 ± 0.006	$\textbf{0.410} \pm \textbf{0.009}$	$\textbf{0.180} \pm \textbf{0.009}$	0.170 ± 0.006
	SAMPLED	0.927 ± 0.001	$\textbf{0.933} \pm \textbf{0.000}$	0.650 ± 0.005	0.652 ± 0.010
THERMOSTABILITY	MIXED	0.349 ± 0.011	$\textbf{0.453} \pm \textbf{0.007}$	0.636 ± 0.013	0.616 ± 0.011
	Human	0.511 ± 0.016	$\textbf{0.589} \pm \textbf{0.002}$	0.382 ± 0.022	$\textbf{0.405} \pm \textbf{0.015}$
	HUMAN-CELL	0.490 ± 0.021	$\textbf{0.570} \pm \textbf{0.004}$	0.316 ± 0.018	$\textbf{0.355} \pm \textbf{0.012}$

Table 1: Comparison between MSE and Bradley-Terry losses on FLIP benchmark tasks using the CNN baseline model. Each row represents a data set and split combination. Numerical columns indicate the mean and standard deviation of test set metrics over 10 random initializations of the model. Asterisks indicate that unmodified portions of sequences were used in training data. Bold values indicate that a loss has significantly improved performance over all other tested losses (p < 0.05). Additional benchmark results are shown in Appendix G.

L=10, K=2 fitness function from the NK model and applied the nonlinearity $g(f)=\exp(10\cdot f)$ to produce an observed fitness function. Then, for each of a range of training set sizes between 25 and 1000, we randomly sampled a training set and fit models with MSE and BT losses using the same models and procedure as in the previous simulations. We repeated this process for 200 training set replicates of each size, and calculated both the Spearman and Pearson correlations between the resulting model predictions and true observed fitness values for all sequences in the sequence space.

Fig. 2c shows the mean correlation values across all 200 replicates of each training set size. There are two important takeaways from this plot. First, the BT loss achieves higher Spearman correlations than the MSE loss in all data regimes. This demonstrates the general effectiveness of this loss to estimate fitness functions affected by global epistasis. Next, we see that models trained with BT loss converge to a maximum Spearman correlation faster than models trained with MSE loss do to a maximum Pearson correlation, which demonstrates that the difference in predictive performance between models trained with MSE and BT losses is not simply due to a result of the evaluation metric being more tailored to one loss than the other. This result also reinforces our claim that fitness functions affected by global epistasis require more data to learn effectively with MSE loss, as would be predicted by CS scaling laws. The BT loss on the other hand, while not performant with the Pearson metric as expected by a ranking loss, seems to overcome this barrier and can be used to estimate a fitness function from a small amount of data, despite the effects of global epistasis.

3.4 Empirical benchmark results

In the previous sections, we used noiseless simulated data to explore the interaction between global epistasis and loss functions. Now we present results demonstrating the practical utility of our insights by comparing the predictive performance of models trained with MSE and BT losses on experimentally determined protein fitness data. We particularly focus on the FLIP benchmark [14], which comprises a total of 15 fitness prediction tasks derived from three empirical fitness datasets. These three datasets explore multiple types of proteins, definitions of protein fitness, and experimental assays. In particular, one is a combinatorially complete dataset that contains the binding fitness of all combinations of mutations at 4 positions to the GB1 protein [41], another contains data about the viability of Adeno-associated virus (AAV) capsids for many different sets of mutations to the wild-type capsid sequence [8], and another contains data about the thermostability of many distantly related proteins [23]. Although there are a number of protein fitness benchmarks, we chose to primarily focus on FLIP because the FLIP datasets all contain a large number of sequences with three or more mutations, which is the regime where the effects of global epistasis are most apparent. In contrast, other benchmarks such as ProteinGym [28] mostly contain datasets with only single and double mutants. Below we describe results on the FLIP benchmark; however, we also tested on two datasets in ProteinGym that contained a large number of higher-order mutations and show these corroborating results in Appendix G.

For each of the three FLIP datasets, the benchmark provides multiple train/test splits that are relevant for protein engineering scenarios. For example, in the GB1 and AAV datasets, there are training sets that contain only single and double mutations to the protein, while the associated test sets contain sequences with more than two mutations. This represents a typical situation in protein engineering where data can easily be collected for single mutations (and some double mutations) and the goal is then to design sequences that combine these mutations to produce a sequence with high fitness. In all of the FLIP tasks the evaluation metric is Spearman correlation between model predictions and fitness labels in the test set, since ranking sequences by fitness is the primary task that models are used for in data-driven protein engineering.

In the FLIP benchmark paper, the authors apply a number of different modeling strategies to these splits, including Ridge regression, training a CNN, and a number of variations on fine-tuning the ESM language models for protein sequences [35]. All of these models use a MSE loss to fit the model to the data, along with any model-specific regularization losses. In our tests, we consider only the CNN model as it balances consistently high performance in the benchmark tasks with relatively straightforward training protocols, enabling fast replication with random restarts.

We trained the CNN model on each split using the standard MSE loss and BT contrastive losses. The mean and standard deviation of Spearman correlations between the model predictions and test set labels over 10 random restarts are shown in Table 1, third and fourth columns. By default, the FLIP datasets contain portions of sequences that are never mutated in any of the data (e.g., only 4 positions are mutated in the GB1 data, but the splits contain the full GB1 sequence of length 56). We found that including these unmodified portions of the sequence often did not improve, and sometimes hurt, the predictive performance of the CNN models while requiring significantly increased computational complexity. Therefore most of our results are reported using inputs that contain only sequence positions that are mutated in at least one train or test data point. We found that including the unmodified portions of sequences improved the performance for the Low-vs-High and 3-vs-Rest GB1 well splits, as well as the 1-vs-rest AAV split and so these results are reported in Table 1; in these cases we found both models trained with MSE and contrastive losses had improved performance.

Although Spearman correlation is commonly used to benchmark models of fitness functions, in practical protein engineering settings it also important to consider the ability of the model to classify the sequences with the highest fitness. To test this, we calculated the "top 10% recall" for models trained with the BT and MSE losses on the FLIP benchmark data, which measures the ability of the models to correctly classify the 10% of test sequences with the highest fitness in the test set [18]. These results are shown in the fifth and sixth columns of Table 1. It may be expected that the BT loss improves performance based on Spearman correlation because both are measures of ranking performance; however, the consistently improved performance of the BT loss based on top 10% recall demonstrates that this loss will lead to improvements in fitness prediction that are practical to protein engineering.

The results in Table1 show that using contrastive losses (and particularly the BT loss) consistently results in improved predictive performance across a variety of practically relevant fitness prediction tasks. Further, in no case does the BT loss result in worse performance than MSE. The reasons for this result may be manifold; however, we hypothesize that it is partially a result of sparse latent fitness functions being corrupted by global epistasis. Indeed, it is shown in Otwinowski et al. [31] that a GB1 landscape closely associated with that in the FLIP benchmark is strongly affected by global epistasis. Further, many of the FLIP training sets are severely undersampled in the sense of CS scaling laws, which is the regime in which differences between MSE and contrastive losses are most apparent when global epistasis is present, as shown in Fig. 2.

4 Discussion

Our results leave open a few avenues for future exploration. First, it is not immediately clear in what situations we can expect to observe the nearly-perfect recovery of a latent fitness function as seen in Fig. 1. A theoretical understanding of this result may either cement the promise of the BT loss, or provide motivation for the development of techniques that can be applied in different scenarios. Next, we have made a couple of logical steps in our interpretations of these results that are intuitive, but not fully supported by any theory. In particular, we have drawn an analogy to CS scaling laws

to explain why neural networks trained with MSE loss struggle to learn a fitness function that has a dense representation in the epistatic domain. However, these scaling laws only strictly apply for a specific set of methods that use an orthogonal basis as the representation of the signal; there is no theoretical justification for using them to understanding the training of neural networks (although applying certain regularizations to neural network training can provide similar guarantees [2]). It is also not clear from a theoretical perspective why the BT loss seems to be robust to the dense representations produced by global epistasis. A deeper understanding of these phenomena could be useful for developing improved techniques.

Additionally, our simulations largely do not consider how models trained with contrastive losses may be affected by the complex measurement noise commonly seen in experimental fitness assays based on sequencing counts [10]. Although our simulations mostly do not consider the effects of noise (except for the simple Gaussian noise added in Appendix E.2), our results on multiple empirical benchmarks demonstrate that contrastive losses can be robust to the effects of noise in practical scenarios. Further, we show in Appendix F that the BT loss can be robust to noise in a potentially pathological scenario. A more complete analysis of the effects of noise on contrastive losses would complement these results.

The code for running the simulations used herein is available at https://github.com/dhbrookes/Contrastive-Losses-Global-Epistasis.git.

References

- [1] Amirali Aghazadeh, Hunter Nisonoff, Orhan Ocal, David H. Brookes, Yijie Huang, O. Ozan Koyluoglu, Jennifer Listgarten, and Kannan Ramchandran. Epistatic Net allows the sparse spectral regularization of deep neural networks for inferring fitness functions. *Nature Communications*, 12(1):5225, 2021.
- [2] Amirali Aghazadeh, Nived Rajaraman, Tony Tu, and Kannan Ramchandran. Spectral regularization allows data-frugal learning over combinatorial spaces, 2022.
- [3] Pablo Baeza-Centurion, Belén Miñana, Jörn M Schmiedel, Juan Valcárcel, and Ben Lehner. Combinatorial Genetics Reveals a Scaling Law for the Effects of Mutations on Splicing. *Cell*, 176(3):549—-563.e23, 2019.
- [4] Christopher W. Bakerlee, Alex N. Nguyen Ba, Yekaterina Shulgina, Jose I. Rojas Echenique, and Michael M. Desai. Idiosyncratic epistasis leads to global fitness–correlated trends. *Science*, 376(6593):630–635, 2022.
- [5] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39:324, 1952.
- [6] David H. Brookes, Hahnbeom Park, and Jennifer Listgarten. Conditioning by adaptive sampling for robust design. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 773–782. PMLR, 09-15 Jun 2019.
- [7] David H. Brookes, Amirali Aghazadeh, and Jennifer Listgarten. On the sparsity of fitness functions and implications for learning. *Proceedings of the National Academy of Sciences of the United States of America*, 119(1):e2109649118, 2022.
- [8] Drew H. Bryant, Ali Bashir, Sam Sinai, Nina K. Jain, Pierce J. Ogden, Patrick F. Riley, George M. Church, Lucy J. Colwell, and Eric D. Kelsic. Deep diversification of an AAV capsid protein by machine learning. *Nature Biotechnology*, pages 1–6, 2021.
- [9] Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. Learning to rank using gradient descent. In *Proceedings of the 22nd International Conference on Machine Learning*, ICML '05, page 89–96, New York, NY, USA, 2005. Association for Computing Machinery.
- [10] Akosua Busia and Jennifer Listgarten. MBE: model-based enrichment estimation and prediction for differential sequencing data. *Genome Biology*, 24(1):218, 2023.

- [11] E.J. Candes, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52 (2):489–509, 2006. doi: 10.1109/TIT.2005.862083.
- [12] Alvin Chan, Ali Madani, Ben Krause, and Nikhil Naik. Deep extrapolation for attribute-enhanced generation. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [13] Wei Chen, Tie-yan Liu, Yanyan Lan, Zhi-ming Ma, and Hang Li. Ranking Measures and Loss Functions in Learning to Rank. In Y Bengio, D Schuurmans, J Lafferty, C Williams, and A Culotta, editors, *Advances in Neural Information Processing Systems*, volume 22. Curran Associates, Inc., 2009.
- [14] Christian Dallago, Jody Mou, Kadina E. Johnston, Bruce Wittmann, Nick Bhattacharya, Samuel Goldman, Ali Madani, and Kevin K. Yang. FLIP: Benchmark tasks in fitness landscape inference for proteins. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [15] Amir Dembo, Thomas M. Cover, and Joy A. Thomas. Information Theoretic Inequalities. *IEEE Transactions on Information Theory*, 37(6):1501–1518, 1991.
- [16] David Ding, Anna G. Green, Boyuan Wang, Thuy-Lan Vo Lite, Eli N. Weinstein, Debora S. Marks, and Michael T. Laub. Co-evolution of interacting proteins through non-contacting and non-specific mutations. *Nature Ecology & Evolution*, 6(5):590–603, 2022.
- [17] Holger Eble, Michael Joswig, Lisa Lamberti, and Will Ludington. Master regulators of evolution and the microbiome in higher dimensions, 2020. URL https://arxiv.org/abs/ 2009.12277.
- [18] Sam Gelman, Sarah A. Fahlberg, Pete Heinzelman, Philip A. Romero, and Anthony Gitter. Neural networks to learn protein sequence–function relationships from deep mutational scanning data. *Proceedings of the National Academy of Sciences*, 118(48):e2104878118, 2021.
- [19] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), volume 2, pages 1735–1742. IEEE, 2006.
- [20] Ralf Herbrich, Thore Graepel, and Klause Obermayer. Large Margin Rank Boundaries for Ordinal Regression. In Advances in Large Margin Classifiers, chapter 7, pages 115–132. The MIT Press, 1999.
- [21] Chloe Hsu, Hunter Nisonoff, Clara Fannjiang, and Jennifer Listgarten. Learning protein fitness models from evolutionary and assay-labeled data. *Nature Biotechnology*, 40(7):1114–1122, 2022.
- [22] Kabir Husain and Arvind Murugan. Physical Constraints on Epistasis. *Molecular Biology and Evolution*, 37(10):2865–2874, 2020.
- [23] Anna Jarzab, Nils Kurzawa, Thomas Hopf, Matthias Moerch, Jana Zecha, Niels Leijten, Yangyang Bian, Eva Musiol, Melanie Maschberger, Gabriele Stoehr, Isabelle Becher, Charlotte Daly, Patroklos Samaras, Julia Mergner, Britta Spanier, Angel Angelov, Thilo Werner, Marcus Bantscheff, Mathias Wilhelm, Martin Klingenspor, Simone Lemeer, Wolfgang Liebl, Hannes Hahne, Mikhail M. Savitski, and Bernhard Kuster. Meltome atlas—thermal proteome stability across the tree of life. *Nature Methods*, 17(5):495–503, 2020.
- [24] Stuart A. Kauffman and Edward D. Weinberger. The NK model of rugged fitness landscapes and its application to maturation of the immune response. *Journal of Theoretical Biology*, 141 (2):211–245, 1989.
- [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.

- [26] Sergey Kryazhimskiy, Daniel P. Rice, Elizabeth R. Jerison, and Michael M. Desai. Global epistasis makes adaptation predictable despite sequence-level stochasticity. *Science*, 344(6191): 1519–1522, 2014.
- [27] Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena Hurtado, Aidan N Gomez, Debora Marks, and Yarin Gal. Tranception: Protein fitness prediction with autoregressive transformers and inference-time retrieval. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 16990–17017. PMLR, 17–23 Jul 2022.
- [28] Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood van Niekerk, Steffanie Paul, Han Spinner, Nathan Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, Jonathan Frazer, Mafalda Dias, Dinko Franceschi, Yarin Gal, and Debora Marks. Proteingym: Large-scale benchmarks for protein fitness prediction and design. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, Advances in Neural Information Processing Systems, volume 36, pages 64331–64379. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/cac723e5ff29f65e3fcbb0739ae91bee-Paper-Datasets_and_Benchmarks.pdf.
- [29] Jakub Otwinowski. Biophysical Inference of Epistasis and the Effects of Mutations on Protein Stability and Function. *Molecular Biology and Evolution*, 35(10):2345–2354, 2018.
- [30] Jakub Otwinowski and Joshua B. Plotkin. Inferring fitness landscapes by regression produces biased estimates of epistasis. *Proceedings of the National Academy of Sciences*, 111(22): E2301–E2309, 2014.
- [31] Jakub Otwinowski, David M. McCandlish, and Joshua B. Plotkin. Inferring the shape of global epistasis. *Proceedings of the National Academy of Sciences*, 115(32):E7550–E7558, 2018.
- [32] Frank J. Poelwijk, Michael Socolich, and Rama Ranganathan. Learning the pattern of epistasis linking genotype and phenotype in a protein. *Nature Communications*, 10(1):4213, 2019.
- [33] Victoria O. Pokusaeva, Dinara R. Usmanova, Ekaterina V. Putintseva, Lorena Espinar, Karen S. Sarkisyan, Alexander S. Mishin, Natalya S. Bogatyreva, Dmitry N. Ivankov, Arseniy V. Akopyan, Sergey Y Avvakumov, Inna S. Povolotskaya, Guillaume J. Filion, Lucas B. Carey, and Fyodor A. Kondrashov. An experimental assay of the interactions of amino acids from orthologous sequences shaping a complex fitness landscape. *PLOS Genetics*, 15(4):1–30, 2019.
- [34] Gautam Reddy and Michael M. Desai. Global epistasis emerges from a generic model of a complex trait. *Elife*, 10:e64740, 2021.
- [35] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021. doi: 10.1073/pnas.2016239118.
- [36] Zachary R. Sailer and Michael J. Harms. Detecting High-Order Epistasis in Nonlinear Genotype-Phenotype Maps. *Genetics*, 205(3):1079–1088, 2017.
- [37] Karen S. Sarkisyan, Dmitry A. Bolotin, Margarita V. Meer, Dinara R. Usmanova, Alexander S. Mishin, George V. Sharonov, Dmitry N. Ivankov, Nina G. Bozhanova, Mikhail S. Baranov, Onuralp Soylemez, Natalya S. Bogatyreva, Peter K. Vlasov, Evgeny S. Egorov, Maria D. Logacheva, Alexey S. Kondrashov, Dmitry M. Chudakov, Ekaterina V. Putintseva, Ilgar Z. Mamedov, Dan S. Tawfik, Konstantin A. Lukyanov, and Fyodor A. Kondrashov. Local fitness landscape of the green fluorescent protein. *Nature*, 533(7603):397–401, 2016. ISSN 1476-4687.
- [38] Peter F. Stadler. Towards a theory of landscapes. In Ramón López-Peña, Henri Waelbroeck, Riccardo Capovilla, Ricardo García-Pelayo, and Federico Zertuche, editors, *Complex Systems and Binary Networks*, pages 78–163, Berlin, Heidelberg, 1995. Springer Berlin Heidelberg.
- [39] Ammar Tareen, Mahdi Kooshkbaghi, Anna Posfai, William T. Ireland, David M. McCandlish, and Justin B. Kinney. MAVE-NN: learning genotype-phenotype maps from multiplex assays of variant effect. *Genome Biology*, 23(1), 2022.

- [40] E. D. Weinberger. Fourier and Taylor series on fitness landscapes. *Biological Cybernetics*, 65 (5):321–330, 1991.
- [41] Nicholas C. Wu, Lei Dai, C. Anders Olson, James O. Lloyd-Smith, and Ren Sun. Adaptation in protein fitness landscapes is facilitated by indirect paths. *eLife*, 5:e16965, 2016.
- [42] Zachary Wu, S. B. Jennifer Kan, Russell D. Lewis, Bruce J. Wittmann, and Frances H. Arnold. Machine learning-assisted directed protein evolution with combinatorial libraries. *Proceedings of the National Academy of Sciences*, 116(18):8852–8858, 2019.
- [43] Kevin K Yang, Zachary Wu, and Frances H. Arnold. Machine-learning-guided directed evolution for protein engineering. *Nature Methods*, 16(8):687–694, 2019.
- [44] Juannan Zhou and David M. McCandlish. Minimum epistasis interpolation for sequence-function relationships. *Nature Communications*, 11(1):1782, 2020.
- [45] Juannan Zhou, Mandy S. Wong, Wei-Chia Chen, Adrian R. Krainer, Justin B. Kinney, and David M. McCandlish. Higher-order epistasis and phenotypic prediction. *Proceedings of the National Academy of Sciences*, 119(39):e2204233119, 2022.

Appendix A Simulation details

Here we provide more specific details about the simulations used to generate Figures 1 and 2 in the main text.

A.1 Complete data simulation

This section provides more details on the simulation that is described in Section 3.1 and Fig. 1 in the main text, in which the goal was to recover a latent fitness function given complete fitness data. For this task, we first sampled a complete latent fitness function, $f(\mathbf{x}_i)$ for $i=1,2,...2^L$, from the NK model, using the parameters L=8, K=2 and q=2. The NK model was implemented as in [7], using the random neighborhood scheme (described in detail in Section A.3, below) We then applied a monotonic nonlinearity to the fitness function to produce a complete set of observed data, $y_i=g(f(\mathbf{x}_i))$ for $i=1,2,...2^L$. We then constructed a neural network model, f_{θ} , in which input binary sequences of length L were transformed by two hidden layers with 100 nodes each and ReLU activation functions and a final linear layer that produced a single fitness output. To fit the parameters of this model, we performed stochastic gradient descent using the Adam method [25] on the Bradley-Terry (BT) loss with all (\mathbf{x}_i, y_i) pairs as training data, a learning rate of 0.001, a batch size of 256. The optimization procedure was terminated when the Spearman between the model's predictions and the observed data y failed to improve over 100 epochs. Letting $\hat{\theta}$ be the parameters of the model at the end of the optimization, we denote the estimated fitness function as $\hat{f} := f_{\hat{\theta}}$.

In order to calculate the epistatic representations of the latent, observed and estimated fitness functions, we arranged each of these fitness functions into appropriate vectors: \mathbf{f} , \mathbf{y} , and $\hat{\mathbf{f}}$, respectively. These vectors are arranged such that the i^{th} element corresponds to the sequence represented by the i^{th} row of the Graph Fourier basis. In this case of binary sequences, the Graph Fourier basis is known as the Walsh-Hadamard basis and is easily constructed, as in [1] and [7]. The epistatic representations of of the latent, observed and estimated fitness functions were then calculated as $\boldsymbol{\beta} = \boldsymbol{\Phi}^T \mathbf{f}$, $\boldsymbol{\beta}_y = \boldsymbol{\Phi}^T \mathbf{y}$, and $\hat{\boldsymbol{\beta}} = \boldsymbol{\Phi}^T \hat{\mathbf{f}}$, respectively.

A.2 Incomplete data simulations

This section provides further details on the simulations in Section 3.3 and Fig. 2 where we tested the ability of models trained with MSE and BT losses to estimate fitness functions given incomplete data corrupted by global epistasis. Many of the details of these simulations are provided in the main text. We used a different set of settings of the a for each nonlinearity used in the simulations shown in Fig. 2b. In particular, for the exponential, sigmoid, and cubic functions we used 20 logarithmically spaced values of a in the ranges [0.1, 50], [0.5, 25], and [0.001, 10], respectively. In all cases where models were trained in these simulations, we used the same neural network model described in the previous section. In all of these cases, we performed SGD on either the MSE or BT loss using the Adam method with a learning rate of 0.001 and batch size equal to 64. In cases where the size of the training set was less than 64, we set the batch size to be equal to the size of the training set. In all cases, the optimization was terminated when the validation metric failed to improve after 20 epochs. The validation metrics for models trained with the BT and MSE losses were the Spearman and Pearson correlations, respectively, between the model predictions on the validation set and the corresponding labels. In order to calculate the yellow curve in Fig. 2a and the values on the horizontal axis in Fig. 2b, the epistatic representations of the observed fitness functions were calculated as described in the previous section.

A.3 Definition of the NK model

The NK model is an extensively studied random field model of fitness functions introduced by Kauffman and Weinberger [24]. In order to sample a fitness function from the NK model, first one chooses values of the parameters L, K, and q, which correspond to the sequence length, maximum degree of epistatic interaction, and size of the sequence alphabet, respectively. Next, one samples a "neighborhood" $\mathcal{V}^{[j]}$ for each position j in the sequence, which represents the K positions that interact with position j. Concretely, for each position j=1,...,L, $\mathcal{V}^{[j]}$ is constructed by sampling K values from the set of positions $\{1,...,L\}\setminus j$ uniformly at random, without replacement. Now let $\mathcal{S}^{(L,q)}$

be the set of all q^L sequences of length L and alphabet size q. Given the sampled neighborhoods $\{\mathcal{V}^{[j]}\}_{j=1}^L$, the NK model assigns a fitness to every sequence in $\mathcal{S}^{(L,q)}$ through the following two steps:

- 1. Let $\mathbf{s}^{[j]} \coloneqq (s_k)_{k \in \mathcal{V}^{[j]}}$ be the subsequence of \mathbf{s} corresponding to the indices in the neighborhood $\mathcal{V}^{[j]}$. Assign a 'subsequence fitness', $f_j(\mathbf{s}^{[j]})$ to every possible subsequence, $\mathbf{s}^{[j]}$, by drawing a value from the normal distribution with mean equal to zero and variance equal to $^1/L$. In other words, $f_j(\mathbf{s}^{[j]}) \sim \mathcal{N}(0, ^1/L)$ for every $\mathbf{s}^{[j]} \in \mathcal{S}^{(K,q)}$, and for every j=1,2,...,L.
- 2. For every $\mathbf{s} \in \mathcal{S}^{(L,q)}$, the subsequence fitness values are summed to produce the total fitness values $f(\mathbf{s}) = \sum_{j=1}^L f_j(\mathbf{s}^{[j]})$.

Appendix B Additional complete data results

Here we provide additional results to complement those described in Section 3.1 and shown in Fig. 1.

B.1 Additional nonlinearities and latent fitness function

In order to more fully demonstrate the ability of models trained with the BT loss to recover sparse latent fitness functions, we repeated the test shown in Fig. 1 and described in Sections 3.1 and A.1 for multiple examples of nonlinearities and different settings of the K parameter in fitness functions sampled from the NK model. In each example, a latent fitness function of L=8 binary sequences was sampled from the NK model with a chosen setting of K and a nonlinearity g(f) was applied to the latent fitness function to produce observed data y. The results of these simulations are shown in Fig. 3, below. In all cases, the model almost perfectly recovers the sparse latent fitness function given complete fitness data corrupted by global epistasis.

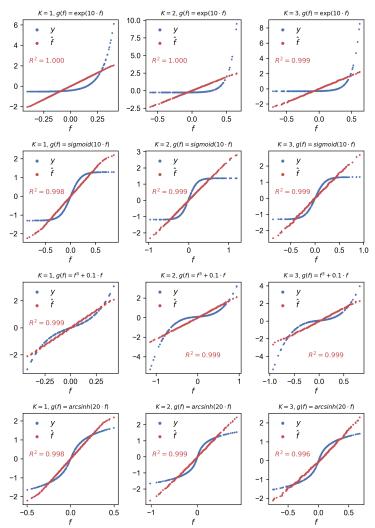


Figure 3: Results from multiple examples of the task of recovering a latent fitness function given complete observed data transformed by a global epistasis nonlinearity. Each sub-plot shows the results of one such task. The setting of K used to sample the latent fitness function from the NK model and the particular form of the nonlinearity g(f) used are indicated in each sub-plot title. The horizontal axis in each sub-plot represents the values of the latent fitness function, while the vertical axis represents the values of either the observed data (blue dots) or model predictions (red dots). For ease of plotting, all fitness functions were normalized to have an empirical mean and std. dev. of 1, respectively. The R^2 correlation between the latent fitness function and the model predictions are indicated in red text.

Appendix C Characterizing nonlinearities that decrease entropy

Here we provide a sufficient condition on a nonlinearity for the nonlinearity to reduce the entropy of a fitness function in the fitness domain. In other words, we characterize a class of g such that $H(g(\mathbf{f})) \leq H(\mathbf{f})$. First, we define probability vectors \mathbf{p} and \mathbf{q} such that $p_i \coloneqq \frac{f_i^2}{||\mathbf{f}||^2}$ and $q_i \coloneqq \frac{g(f_i)^2}{||g(\mathbf{f})||^2}$ for all i=1,...,N.

We additionally define the vector \mathbf{w} such that $w_i \coloneqq \frac{q_i}{p_i}$ and the probability vector $\mathbf{r}(\alpha)$ where $r_i(\alpha) = \frac{p_i w_i^{\alpha}}{\sum_{i=1}^N p_i w_i^{\alpha}}$ for all i=1,...,N. Note that $\mathbf{r}(0) = \mathbf{p}$ and $\mathbf{r}(1) = \mathbf{q}$.

Finally, let $H_s(\mathbf{p}) = -\sum_{i=1}^{N} p_i \log p_i$ be the Shannon entropy of a probability vector \mathbf{p} .

A note on notation: below we are primarily concerned with distributions over discrete outcomes indexed by integers. Given a probability vector \mathbf{p} representing such a distribution and a vector \mathbf{x} representing values associated with each outcome, we denote the expectation of values of \mathbf{x} over the distribution represented by \mathbf{p} as $E_p[x] = \sum_i p_i x_i$. Similarly, we write covariances as $\operatorname{Cov}_p(x,y) = E_p[xy] - E_p[x]E_p[y]$.

Theorem 1. Assume without loss of generality that \mathbf{p} is sorted such that $p_i \geq p_{i+1}$ for i=1,...,N-1. Additionally assume that if $p_i=0$ then $q_i=0$ for all i. If $w_i \geq w_{i+1}$ for all i=1,...,N-1, then $H(g(\mathbf{f})) \leq H(\mathbf{f})$.

To prove Theorem 1, we first provide a number of lemmas.

Lemma 1. The derivative of the Shannon entropy of $\mathbf{r}(\alpha)$ with respect to α is equal to the negative covariance between $\log \mathbf{r}(\alpha)$ and $\log \mathbf{w}$. Specifically, the derivative is given by

$$\frac{\mathrm{d}H_s(\mathbf{r}(\alpha))}{\mathrm{d}\alpha} = -\mathrm{Cov}_{r(\alpha)}(\log r(\alpha), \log w).$$

Lemma 2. Given a probability vector \mathbf{p} of length N and any two vectors \mathbf{x} and \mathbf{y} of length N, we have

$$Cov_p(x,y) = \sum_{j>i} p_i p_j (x_i - x_j) (y_i - y_j).$$
 (5)

Corollary 1. If two vectors x and y of length N are sorted such that $x_i \ge x_{i+1}$ and $y_i \ge y_{i+1}$ for i = 1, ..., N-1, then $Cov_p(x, y) \ge 0$ for any probability vector p.

Proof of Theorem 1. First we note that by the definition of entropy given in the main text, we have $H(\mathbf{f}) = H_s(\mathbf{p})$ and $H(g(\mathbf{f})) = H_s(\mathbf{q})$. Therefore we need only prove that $H_s(\mathbf{q}) \leq H_s(\mathbf{p})$.

By Lemma 1, we have that

$$H_s(\mathbf{q}) - H_s(\mathbf{p}) = -\int_0^1 \text{Cov}_{r(\alpha)}(\log r(\alpha), \log w) d\alpha.$$
 (6)

Both \mathbf{p} and \mathbf{w} are sorted by assumption and therefore it is clearly true that $\mathbf{r}(\alpha)$ is also sorted such that $r_i(\alpha) \geq r_{i+1}(\alpha)$ for i=1,...,N-1 and all $\alpha \in [0,1]$. Since the logarithm is a monotonic function, both $\log \mathbf{r}(\alpha)$ and $\log \mathbf{w}$ are similarly sorted. Therefore, by Corollary 1, we have that the integrand in Eq. 6 is positive for all $\alpha \in [0,1]$. Therefore, the integral must be positive and $H_s(\mathbf{q}) - H_s(\mathbf{p}) < 0$.

Proof of Lemma 1.

$$\begin{split} \frac{\mathrm{d}H_s(\mathbf{r}(\alpha))}{\mathrm{d}\alpha} &= -\sum_{i=1}^N (1 + \log r_i(\alpha)) \frac{\mathrm{d}r_i(\alpha)}{\mathrm{d}\alpha} \\ &= -\sum_{i=1}^N (1 + \log r_i(\alpha)) \frac{p_i w_i^{\alpha}}{\sum_{j=1}^N p_i w_i^{\alpha}} \left(\log w_i - \frac{\sum_{k=1}^N p_k w_k^{\alpha} \log w_k}{\sum_{j=1}^N p_i w_i^{\alpha}} \right) \\ &= -\sum_{i=1}^N r_i(\alpha) (1 + \log r_i(\alpha)) \left(\log w_i - \sum_k r_k(\alpha) \log w_k \right) \\ &= -\sum_{i=1}^N r_i(\alpha) \left(\log w_i + \log r_i(\alpha) \log w_i - E_{r(\alpha)} [\log w] - \log r_i(\alpha) E_{r(\alpha)} [\log w] \right) \\ &= -E_{r(\alpha)} [\log w] - E_{r(\alpha)} [\log r(\alpha) \log w] + E_{r(\alpha)} [\log w] + E_{r(\alpha)} [\log r(\alpha)] E_{r(\alpha)} [\log w] \\ &= -E_{r(\alpha)} [\log r(\alpha) \log w] + E_{r(\alpha)} [\log r(\alpha)] E_{r(\alpha)} [\log w] \\ &= -\mathrm{Cov}_{r(\alpha)} (\log r(\alpha), \log w) \end{split}$$

Proof of Lemma 2.

$$\begin{split} &\sum_{i=1}^{N} \sum_{j=i}^{N} p_{i} p_{j} (x_{i} - x_{j}) (y_{i} - y_{j}) = \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} p_{i} p_{j} (x_{i} - x_{j}) (y_{i} - y_{j}) \\ &= \frac{1}{2} \sum_{i=1}^{N} p_{i} x_{i} y_{i} \sum_{j=1}^{N} p_{j} - \frac{1}{2} \sum_{i=1}^{N} p_{i} x_{i} \sum_{j=1}^{N} p_{j} y_{j} - \frac{1}{2} \sum_{i=1}^{N} p_{i} y_{i} \sum_{j=1}^{N} p_{j} x_{j} + \frac{1}{2} \sum_{j=1}^{N} p_{j} x_{j} y_{j} \sum_{i=1}^{N} p_{i} x_{j} \\ &= E_{p}[xy] - E_{p}[x] E_{p}[y] \\ &= \operatorname{Cov}_{p}(x, y) \end{split}$$

where the first line follows from the summand equaling zero when i = j.

Proof of Corollary 1. This is a straightforward consequence of Lemma 2. For every term in the sum of Eq. 5, we have that j>i, and therefore $x_i-x_j\geq 0$ and $y_i-y_j\geq 0$ by assumption. Further ${\bf p}$ is a probability vector and thus $p_i>0$ for all i. Therefore every element of the sum is positive and the covariance is positive.

Notably, the condition of Theorem 1 is not necessary in order decrease entropy in the fitness domain. Indeed, many of the nonlinearities tested herein do not satisfy this condition. Eq. 6 is an illuminating relationship that may provide further insight into additional classes of nonlinearities that decrease entropy.

Appendix D Additional examples of uncertainty principle

Here we provide additional examples showing that nonlinearities tend to decrease the entropy in the fitness domain of sparse fitness, which causes corresponding increases in the entropy in the epistatic domain due to the uncertainty principle of Eq. 4. We show this for four nonlinearities:exponential, $g(f) = \exp(a \cdot f)$, sigmoid, $g(f) = (1 + e^{-a \cdot f})^{-1}$, cubic polynomial $g(f) = x^3 + ax$ with a > 0, and a hinge function that represents the left-censoring of data, $g(f) = \max(0, f - a)$. Left-censoring is common in fitness data, as biological sequences will often have fitness that falls below the detection threshold of an assay. For each of these nonlinearities, we chose a range of settings of the a parameter, and for each setting of the a parameter, we sampled 200 fitness functions from the NK model with parameters L = 8, K = 2 and q = 2 and transformed each function by the nonlinearity with the chosen setting of a. For each of these transformed fitness functions, we calculated the entropy of the function in the fitness and epistatic domains. The mean and standard deviation of these entropies across the fitness function replicates are shown in Fig. 4, below. In each case, we can see that the nonlinearities cause the fitness domain to become increasingly concentrated, and the epistatic domain to become increasingly sparse.

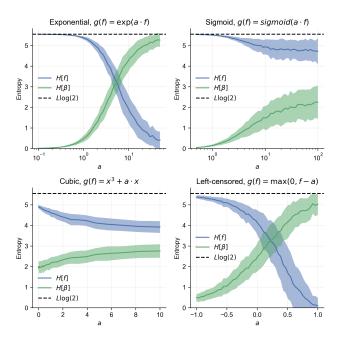


Figure 4: Demonstration of the fitness-epistasis uncertainty principle for multiple examples of nonlinearities. The title of the subplot indicates the nonlinearity used to produce the results in that subplot. The lines and shaded regions represent the mean and std. dev. of entropies, respectively, across 200 replicates of latent fitness functions sampled from the NK model. The black dotted line indicates the lower bound on the sum of the entropies in Eq. 4.

Appendix E Additional incomplete data results

Here we provide additional results to complement those described in Section 3.3 and shown in Fig. 2.

E.1 Additional nonlinearity and latent fitness function parameters

Here we show results for incomplete data simulations using latent fitness functions with varying orders of epistatic interactions and for additional parameterizations of the global epistasis nonlinearity. In particular, we repeated the simulations whose results are shown in Fig. 2b using latent fitness functions drawn from the NK model with K=1 and K=3, with all other parameters of the simulation identical to those used in the K=2 simulations described in the main text. We ran the simulations for 10 replicates of each of the 20 settings of the a parameter for each of the three nonlinearities. The results of these simulations are shown in Figure 5, below. These results are qualitatively similar to those in Figure 5, demonstrating that the simulation results are robust to the degree of epistatic interactions in the latent fitness function.

We also repeated the simulations whose results are shown in Figure 2c using different settings of the K parameter in the NK model, as well as multiple settings of the a parameter that determines the strength of the global epistasis nonlinearity. All other parameters of the simulations are identical to those used in the K=2 and a=10 simulation described in the main text. The results of these simulations, averaged over 40 replicates for each training set size, are shown in Fig 6, below. These results are qualitatively similar to those shown in Figure 2c, demonstrating the robustness of the simulation results to different forms of latent fitness functions and nonlinearities.

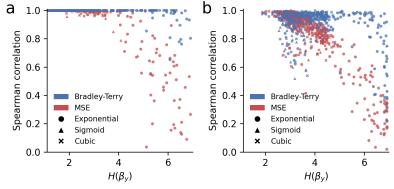


Figure 5: Results from incomplete data simulations for latent fitness functions drawn from the NK model with (a) K = 1 and (b) K = 2. Plot descriptions are as in Figure 2b.

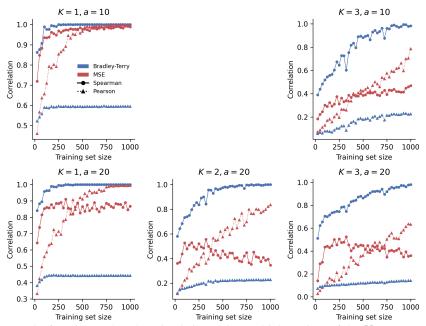


Figure 6: Results from incomplete data simulations using multiple settings of the K parameter in the NK model and the a parameter in the global epistasis nonlinearity. The settings of K and a used to generate each subplot are shown in the title of the subplots; results are not shown for the K=2 and a=10 case because these results are shown in Figure 2c of the main text. Plot descriptions are as in Figure 2c; each point in the plots represents the mean over 40 replicate simulations.

E.2 Adding noise to the observed data

Here we test how the addition of homoskedastic Gaussian noise affects the results incomplete data simulations. In particular, we repeated the set of simulations in which we test the ability of the MSE and BT losses to estimate a fitness function at varying sizes of training set, but now added Gaussian noise with standard deviation of either 0.025 or 0.1 to the observed fitness function (first normalized to have mean and variance equal to 0 and 1, respectively). The results of these tests, averaged over 40 replicates for each training set size, are shown in Figure 7, below. We can see that despite the noise drastically altering the observed ranking, the BT loss still outperforms the MSE loss in terms of Spearman correlation at both noise settings.

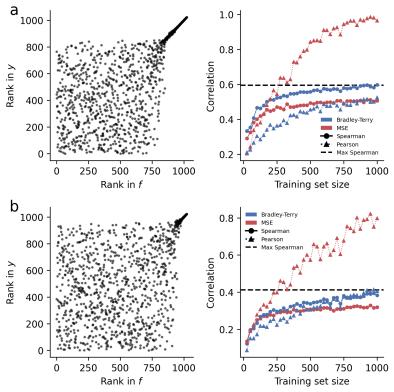


Figure 7: Incomplete data simulations with additive Gaussian noise with standard deviation equal to (a) 0.025 and (b) 0.1. Left plots show the rankings of sequences in the latent NK fitness function, \mathbf{f} , against the ranking of sequences in the noisy observed fitness function, \mathbf{y} . The right plot descriptions are as in Figure 2c, with the addition of black dashed lines that indicates the Spearman correlation between the latent fitness function, \mathbf{f} , and the noisy observed fitness function, \mathbf{y} , which is the maximum Spearman correlation that can be achieved in this test.

Appendix F Noisy toy simulation

A potential disadvantage of the BT and Margin losses used in the main text is that neither takes into consideration the size of the true observed gap between the fitness of pairs of sequences. In particular, Eq. 3 shows that any two sequences in which $y_i > y_j$ will be weighted equally in the loss, regardless of the magnitude of the difference $y_i - y_j$. This has implications for how measurement noise may affect these losses. In particular, the relative ranking between pairs of sequences with a small gap in fitness in more likely to have been swapped due to measurement noise, compared to a pair of sequences with a large gap. This suggests that contrastive losses may exhibit pathologies in the presence of certain types of noise.

Here we examine a toy scenario in which the observed fitness data has the form $y_i = I(\mathbf{x}_i) + \epsilon$, where $I(\mathbf{x}_i)$ is an indicator function that is equal to zero or one and $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is Gaussian noise. One may expect contrastive loses to exhibit pathologies when trained on this data because the relative rankings between sequences that have the same value of $I(\mathbf{x})$ is due only to noise.

In order to construct this scenario, we sampled an L=10, K=2 binary fitness function $f(\mathbf{x})$ from the NK model, and then let $I(\mathbf{x})=0$ when $f(\mathbf{x})< \mathrm{med}(f)$ and 1 otherwise, where $\mathrm{med}(f)$ is the median NK fitness of all sequences. We then added Gaussian noise with $\sigma=0.1$ to produce the observed noisy data g. We consider $I(\mathbf{x})$ to be the true binary label of a sequence and g to be the noisy label of the sequence. The relationship between the NK fitness, true labels and noisy labels is shown in Fig. 8a, below. Next, we randomly split the data into training and test sets containing 100 and 924 sequences and noisy labels, respectively. We used the training set to fit two neural network models, one using the MSE loss and the other using the BT loss. The models and training procedure were the same as described in Appendix A. The predictions of these models on test sequences, compared to noisy labels, are shown in 8b. We next constructed Receiver Operating Characteristic (ROC) curves that test the ability of each model to classify sequences according to their true labels $I(\mathbf{x})$ (Fig. 8c). These curves show that while the MSE loss does slightly outperform the BT loss in this classification task, the BT loss still results in an effective classifier and does not exhibit a major pathology in this scenario.

Notably, on empirical protein landscapes, as shown in Table 1, the performance gain by contrastive losses out-weighs the potential drawback of further sensitivity to this kind of noise, and contrastive losses ultimately result in improved predictive performance.

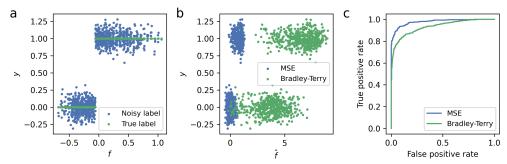


Figure 8: Toy simulation testing the effects of noise on the BT loss. (a) Noisy and true labels of training and test data (vertical axis), plotted against NK fitness (horizontal axis). True labels correspond to $I(\mathbf{x})$ values, while noisy labels correspond to $y = I(\mathbf{x}) + \epsilon$ values. (b) Observed noisy labels (vertical axis) plotted against model predictions from a models trained with the MSE and BT losses (horizontal axis). (c) ROC curves comparing the ability of the models trained with MSE and BT losses to classify test sequences according to true labels. The blue (MSE) and green (BT) curves have AUC values of 0.977 and 0.932, respectively.

Appendix G Results on ProteinGym benchmark datasets

We compared the MSE vs Bradley-Terry loss on the CAPSD_AAV2S_Sinai_2021 and GFP_AEQVI_Sarkisyan_2016 datasets from the ProteinGym benchmark [28] using the same protocol as described in 3.4. These are two of the most relevant datasets in ProteinGym because they contain a large number of multi-mutants, whereas many other datasets contain only single and double mutations. The results for 10 replicates of uniform 80/20 train test splits are shown below:

	SPEARMAN		TOP 10% RECALL	
DATASET	MSE Loss	BRADLEY-TERRY	MSE Loss	BRADLEY-TERRY
CAPSD_AAV2S_SINAI_2021 GFP_AEQVI_SARKISYAN_2016	0.912 ± 0.003 0.867 ± 0.001	$\begin{array}{c} \textbf{0.920} \pm \textbf{0.003} \\ \textbf{0.873} \pm \textbf{0.001} \end{array}$	$\begin{array}{c} 0.915 \pm 0.001 \\ 0.938 \pm 0.003 \end{array}$	0.915 ± 0.001 0.945 ± 0.003

Table 2: Comparison between MSE and Bradley-Terry losses on ProteinGym benchmark tasks using the CNN baseline model. Each row represents a data set and split combination. Numerical columns indicate the mean and standard deviation of test set metrics over 10 random initializations of the model. Asterisks indicate that unmodified portions of sequences were used in training data. Bold values indicate that a loss has significantly improved performance over all other tested losses (p < 0.05).

Appendix H Derivation of fitness-epistasis uncertainty principle

Here we show how to derive the fitness-epistasis uncertainty of Eq. 4, starting from the uncertainty principle presented in Theorem 23 of [15]. When applied to the transformation $\mathbf{f} = \Phi \boldsymbol{\beta}$, this uncertainty principle is stated as:

$$H(\mathbf{f}) + H(\boldsymbol{\beta}) \ge \log\left(\frac{1}{M^2}\right)$$
 (7)

where $M = \max_{ij} |\Phi_{ij}|$. This can be further simplified by calculating the maximum absolute value of the elements in the Graph Fourier basis, Φ . As shown in [7], this matrix for sequences of length L and alphabet size q is calculated as

$$\mathbf{\Phi} = \bigotimes_{i=1}^{L} \mathbf{P}(q) \tag{8}$$

where \otimes denotes the Kronecker product, and $\mathbf{P}(q)$ is an orthonormal set of eigenvectors of the Graph Laplacian of the complete graph of size q. Based on Eq. 8, we can make the simplification to the calculation of the maximum value in Φ that

$$M = m^L (9)$$

where $m = \max_{ij} |P_{ij}(q)|$. Further, the value of m can be determined for each setting of q by considering the following calculation of $\mathbf{P}(q)$ used in [7]:

$$\mathbf{P}(q) = \mathbf{I} - \frac{2\mathbf{w}\mathbf{w}^T}{||\mathbf{w}||^2} \tag{10}$$

where **I** is the identity matrix and $\mathbf{w} = \mathbf{1} - \sqrt{q}\mathbf{e}_1$, with **1** representing a vector with all elements equal to one, and \mathbf{e}_1 is the vector equal to one in its first element and zero in all other elements. The values of each element in $\mathbf{P}(\mathbf{q})$ can be straightforwardly calculated using Eq. 10. In particular, we have:

$$P_{ij}(q) = \begin{cases} \frac{1}{\sqrt{q}} & \text{if } i = 1 \text{ or } j = 1\\ 1 - \frac{1}{q - \sqrt{q}} & \text{if } i = j \neq 1\\ \frac{1}{q - \sqrt{q}} & \text{otherwise.} \end{cases}$$

$$(11)$$

Now all that remains is determining which of these three values has the largest magnitude for each setting of q. For q=2, only the first two values appear in the matrix, and both have magnitude $\frac{1}{\sqrt{q}}=\frac{1}{\sqrt{2}}$. For q=3, we can directly calculate that all three values are positive and we have $\frac{1}{q-\sqrt{q}}>\frac{1}{\sqrt{q}}>1-\frac{1}{q-\sqrt{q}}$. For q=4, we can again directly calculate that all three values are equal to $\frac{1}{\sqrt{q}}=\frac{1}{2}$. For q>4, all three of the values are positive. Further, algebraic manipulations can be used to show that q>4 implies $1-\frac{1}{q-\sqrt{q}}>\frac{1}{q-\sqrt{q}}$ and $\frac{1}{\sqrt{q}}>\frac{1}{q-\sqrt{q}}$. Therefore, $1-\frac{1}{q-\sqrt{q}}$ is the maximum value in $\mathbf{P}(q)$ for all q>4.

Putting these results together with Eq. 9, the fitness-epistasis uncertainty principle simplifies to

$$H(\mathbf{f}) + H(\boldsymbol{\beta}) \ge L \log\left(\frac{1}{m^2}\right)$$
 (12)

where

$$m = \begin{cases} \frac{1}{\sqrt{q}} & \text{if } q = 2\\ \frac{1}{q - \sqrt{q}} & \text{if } q = 3\\ 1 - \frac{1}{q - \sqrt{q}} & \text{if } q \ge 4 \end{cases}$$

which is the form presented in the main text.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Our abstract and introduction provide accurate summaries of the contributions made by the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Our work does have limitations, which we lay out clearly in the Discussions section.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Our paper contains only one true theoretical result, which is the uncertainty principle of Eq.. 4. We provide a complete derivation of this result in Appendix H. The remainder of our claims are supported by experiments.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Our simulation experiments are clearly described, with all necessary parameters and procedures defined in either the main text or Appendix. Further, we will release the code to run these simulations upon acceptance. Our benchmarking results use the FLIP benchmark datasets and code, which have been published in the NeurIPS benchmark track and are publicly accessible.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
 well by the reviewers: Making the paper reproducible is important, regardless of
 whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We will release the code on Github upon acceptance, which will allow any researcher to fully reproduce our results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: [TODO]

Guidelines: Our experiments are clearly described with all information required to understand the result.

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We have provided error bars and measures of statistical signficance in all benchmarking experiments.

Guidelines:

• The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: None of our experiments are sufficiently compute-intensive to require us to specify this.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our paper conforms to the Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Our paper is purely scientific in nature and does not have immediate positive or negative societal impacts.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Nothing in our paper would require such safeguards.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All contributors to the paper have been properly credited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets are introduce in our paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: We did no research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Eq.uivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We did no research with human subjects.

Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.