

# Inference of fitness landscapes with heterogeneous patterns of epistasis across sites

Carlos Martí-Gómez<sup>1,\*</sup> and David M. McCandlish<sup>1,\*</sup>

<sup>1</sup>Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 11724

\*Correspondence: [martigo@cshl.edu](mailto:martigo@cshl.edu), [mccandlish@cshl.edu](mailto:mccandlish@cshl.edu)

## 1 Abstract

2 Fitness landscapes provide a framework for understanding how genetic variation shapes evolutionary outcomes. Although these landscapes  
3 were long treated as abstract conceptual objects, recent advances in genetic engineering and high-throughput phenotyping have enabled the  
4 empirical measurement of phenotypic values across large combinatorial sequence spaces. These developments create a need for statistical  
5 frameworks that can summarize, infer, and interpret fitness landscapes in the presence of complex genetic interactions. Here, we introduce a  
6 framework for summarizing the structure of genetic interactions across sites based on the average squared local  $k$ -way epistatic coefficients  
7 between mutations at different subsets of sites, and derive the precise manner in which the variance in these local  $k$ -way epistatic coefficients  
8 across backgrounds relates to epistasis of orders higher than  $k$ . These statistics can be computed exactly for complete combinatorial landscapes  
9 and are related to classical statistics in the fitness landscape literature. Moreover, they can be estimated from empirical correlations when data  
10 are incomplete or noisy, and used to define an empirical Bayes prior for fitness landscape inference that differentially penalizes interactions  
11 involving different subsets of sites. We apply this inference method to diverse high-throughput protein and RNA combinatorial mutagenesis  
12 datasets and find that fitness landscapes often show highly structured patterns of genetic interactions across positions. Finally, we use this  
13 model to infer a fitness landscape for a dynamic self-splicing intron comprising 65,536 genotypes, and describe in detail the main genetic  
14 interactions that shape the structure of this landscape and how they relate to the underlying molecular mechanism. Together, these results  
15 provide new tools for summarizing and modeling complex fitness landscapes, and for linking large-scale empirical data to the mathematical  
16 theory of fitness landscapes.

17 **Keywords:** fitness landscape; epistasis; Gaussian process

## 1 Introduction

2 The fitness landscape is a fundamental concept in evolutionary  
3 biology and genetics. First introduced by Wright (1932), it de-  
4 scribes the mapping between genotypes and their associated fit-  
5 ness. The structure of a fitness landscape, including the number,  
6 distribution, and connectivity of fitness peaks, plays a crucial  
7 role in how populations evolve and diverge over time (Kauff-  
8 man and Levin 1987; Kondrashov *et al.* 2002; Gavrilets 2004;  
9 Weinreich *et al.* 2006; McCandlish 2011; De Visser and Krug  
10 2014; Fragata *et al.* 2019; Manrubia *et al.* 2021; Bank 2022; John-  
11 son *et al.* 2023). Understanding the mapping from genotype  
12 to fitness is not only important for explaining and predicting  
13 evolution, but also has critical applications in cancer and human  
14 disease (Moore and Williams 2009; Dasari *et al.* 2021) as well as  
15 plant and animal breeding (De Los Campos *et al.* 2013; Sackton  
16 and Hartl 2016; Soyk *et al.* 2020; Dwivedi *et al.* 2024). However,  
17 despite its importance, characterizing this mapping is inherently  
18 challenging due to the high dimensionality of sequence space.  
19 Because the number of possible sequences grows exponentially  
20 with sequence length, such high-dimensional fitness landscapes  
21 are often summarized by computing low-dimensional summary  
22 statistics that characterize the ruggedness or structure of the  
23 landscape, e.g. the number of local optima, lengths of adaptive  
24 walks, and the number of alternative local optima accessible

from starting genotype (Kauffman and Levin 1987; Szendro  
*et al.* 2013; Ferretti *et al.* 2018). Another common approach is  
to characterize how mutational effects change across genetic  
backgrounds, for instance by quantifying the average magni-  
tude of local epistatic coefficients (Zhou and McCandlish 2020)  
or by measuring the correlation of mutational effects between  
genotypes separated by increasing numbers of mutations (Wein-  
berger 1990; Stadler 1996; Neidhart *et al.* 2013; Bank *et al.* 2016;  
Ferretti *et al.* 2016).

Historically, the scarcity of comprehensive experimental data  
has motivated the development of theoretical and computational  
models of fitness landscapes, which have provided a framework  
for understanding how summary statistics behave across differ-  
ent classes of landscapes. One approach is to consider families  
of fitness landscapes drawn from a probability distribution, gen-  
erally known as random field models (Kauffman and Levin 1987;  
Stadler and Happel 1999). Classical examples include the House  
of Cards model (Kingman 1978), the NK model (Kauffman and  
Levin 1987), and the Rough Mount Fuji landscape (Aita and  
Husimi 1998). Within these frameworks, one can compute the  
expected values of summary statistics and analyze how they  
depend on parameters that control the smoothness or rugged-  
ness of the landscape (Kauffman and Levin 1987; Schmiegelt  
and Krug 2014; Neidhart *et al.* 2014; Hwang *et al.* 2018; Reddy

1 and Desai 2021).

2 More recent efforts have increasingly turned toward the ex-  
3 perimental characterization of empirical fitness landscapes by  
4 measuring growth rates or other measures of biological function-  
5 ality for many different combinations of mutations (De Visser  
6 and Krug 2014; Fragata *et al.* 2019). The size of the first em-  
7 pirically reconstructed fitness landscapes was limited by the  
8 difficulty of engineering large numbers of genotypes and mea-  
9 suring their fitness experimentally (Khan *et al.* 2011; Chou *et al.*  
10 2011; Flynn *et al.* 2013; Szendro *et al.* 2013; Ogbunugafor *et al.*  
11 2016; Weinreich *et al.* 2018; Gao *et al.* 2022; Aguirre *et al.* 2023;  
12 Zebell *et al.* 2025). However, recent advances in high-throughput  
13 assays (Kinney *et al.* 2010; Fowler and Fields 2014; Kinney and  
14 McCandlish 2019) have substantially expanded this scope, en-  
15 abling the parallel measurement of thousands to millions of  
16 genotypes. These techniques have been used to characterize  
17 fitness landscapes across a range of biological systems, includ-  
18 ing regulatory sequences (Noderer *et al.* 2014; Rosenberg *et al.*  
19 2015; Bonde *et al.* 2016; Evfratov *et al.* 2017; Rabani *et al.* 2017;  
20 Wong *et al.* 2018; Baeza-Centurion *et al.* 2019; Kuo *et al.* 2020;  
21 Komarova *et al.* 2020; de Boer *et al.* 2020; Vaishnav *et al.* 2022;  
22 Liao *et al.* 2023; Westmann *et al.* 2024b,a; Kuo *et al.* 2025; Chat-  
23 topadhyay *et al.* 2025; Agarwal *et al.* 2025), RNAs (Domingo *et al.*  
24 2018; Bendixsen *et al.* 2019; Soo *et al.* 2021; Rotrattanadumrong  
25 and Yokobayashi 2022), proteins (O’Maille *et al.* 2008; Bank *et al.*  
26 2016; Wu *et al.* 2016; Starr *et al.* 2017; Poelwijk *et al.* 2019; Lite  
27 *et al.* 2020; Jalal *et al.* 2020; Bryant *et al.* 2021; Somermeyer *et al.*  
28 2022; Moulana *et al.* 2023; Papkou *et al.* 2023; Sundar *et al.* 2024;  
29 Zarin and Lehner 2024; Johnston *et al.* 2024; Faure *et al.* 2024b;  
30 Escobedo *et al.* 2025; Herrera-Álvarez *et al.* 2025), and genome-  
31 wide gene interactions (Bakerlee *et al.* 2022; Nguyen Ba *et al.* 2022;  
32 Matsui *et al.* 2022; N’Guessan *et al.* 2025). Despite these advances,  
33 contemporary datasets are often noisy and typically do not cover  
34 the full sequence space, so that a key challenge is to develop  
35 flexible statistical methods for inferring full fitness landscapes  
36 from empirical data without distorting the rich fitness landscape  
37 geometry revealed by these high-throughput measurements.

38 One powerful approach to address this problem is to combine  
39 theoretical models of fitness landscapes with empirical measure-  
40 ments by recasting these theoretical models as Bayesian priors  
41 for reconstructing complete landscapes from incomplete and  
42 noisy data (Zhou and McCandlish 2020; Chen *et al.* 2021; Zhou  
43 *et al.* 2022, 2025; Petti *et al.* 2025; Martí-Gómez *et al.* 2026b). Such  
44 an approach can leverage our mathematical understanding of  
45 these models to define prior distributions that confer the over-  
46 all inference procedure with desirable properties. For example,  
47 Minimum Epistasis Interpolation defines a prior that depends  
48 on the average squared epistatic coefficients between all pairs  
49 of mutations (Zhou and McCandlish 2020), favoring reconstruc-  
50 tions that are locally approximately additive, which results in  
51 reconstructions that can capture genetic interactions of all or-  
52 ders where data is abundant but extrapolates additively far from  
53 the data. As another example, Empirical Variance Component  
54 regression (Zhou *et al.* 2022) constructs a prior parametrized  
55 by the variance explained by genetic interactions of each possi-  
56 ble order, resulting in reconstructions that accurately reflect  
57 how quickly the predictability of mutational effects decay in  
58 increasingly distant genetic backgrounds.

59 While Minimum Epistasis Interpolation and Empirical Vari-  
60 ance Component regression can incorporate epistatic interac-  
61 tions of all orders, the corresponding priors are still only weakly  
62 informative in the sense that they are “isotropic” (Stadler 1996,

2002), i.e. the prior treats all sites and all mutations equally. 63  
64 However, in reality some sites and alleles are more influential  
65 and more likely to be involved in epistatic interactions than  
66 others (Weinreich *et al.* 2005; Kvitek and Sherlock 2011; Ferretti  
67 *et al.* 2016; Bank *et al.* 2016; Pokusaeva *et al.* 2019; Reddy and  
68 Desai 2021). Reddy and Desai (2021) recently proposed a new  
69 family of theoretical random field models whose parameters  
70 control the site-specific probability that mutations participate in  
71 genetic interactions, and these models were then extended by  
72 Zhou *et al.* (2025) to include allele-specific and mutation-specific  
73 parameters. By treating these site-, allele-, or mutation-specific  
74 parameters as hyperparameters of an informative Bayesian prior,  
75 the resulting model learns which mutations most strongly in-  
76 fluence the predictability of other mutations, and incorporates  
77 this information when inferring the fitness landscape from data,  
78 achieving state-of-the-art predictive performance (Zhou *et al.*  
79 2025). Nonetheless, these models still implicitly assume that the  
80 the propensity for a set of sites to interact is determined solely by  
81 these site-specific parameters, while empirical observations from  
82 both pairwise interaction models (Marks *et al.* 2012; Haldane *et al.*  
83 2016, 2018) and the posterior distributions of models containing  
84 interactions of all orders (Chen *et al.* 2021; Martí-Gómez *et al.*  
85 2026b) suggest that patterns of epistatic interaction are often  
86 sparse (Poelwijk *et al.* 2019) or modular (Rojas Echenique *et al.*  
87 2019; Hwang *et al.* 2018).

88 Here, we present a method for fitness landscape inference  
89 incorporating a prior that can encode this type of highly struc-  
90 tured tendency for specific sets of sites to interact with each  
91 other. We begin by proposing a simple approach to summarize  
92 the structure of genetic interactions of all orders by comput-  
93 ing, for each pair of sites, the average squared local epistatic  
94 coefficient between mutation at those sites, which enables the  
95 identification of sets of sites involved in epistatic interactions  
96 of arbitrary order in complete combinatorial fitness landscapes.  
97 Based on these summary statistics, we define a new family of  
98 prior distributions that differentially penalize interactions in-  
99 volving different pairs of sites. The parameters of these priors  
100 can be estimated from empirical correlations in fitness values  
101 between sequences that differ at specific subsets of sites, and we  
102 then use the resulting priors to infer complete fitness landscapes  
103 from several empirical datasets. Finally, we apply our method  
104 to the fitness landscape of a dynamic self-splicing intron (Soo  
105 *et al.* 2021), and show how the higher-order interactions in this  
106 system have an interpretable structure wherein many aspects of  
107 the genetic architecture are systematically rewired depending  
108 on the nucleotide identities at one specific pair of sites.

## 109 Results

### 110 Epistatic coefficients for subsets of sites

111 In this section, our aim is to quantify the overall amount of  
112 epistasis between a specific pair of sites in an arbitrary fitness  
113 landscape  $f$ . We start by considering a simple two-locus bial-  
114 lelic fitness landscape and assuming that we know the fitness  
115 values of the four possible genotypes  $f_{AB}$ ,  $f_{Ab}$ ,  $f_{aB}$  and  $f_{ab}$ . We  
116 can measure how much the effect of a mutation  $A \rightarrow a$  in one  
117 locus changes in the presence of an additional mutation  $B \rightarrow b$   
118 at the other locus via the traditional double-mutant epistatic  
119 coefficient:

$$\epsilon = (f_{Ab} - f_{ab}) - (f_{AB} - f_{aB}), \quad (1)$$

120 and we can use the squared epistatic coefficient  $\epsilon^2$  to quantify the  
121 magnitude of epistasis in a manner that is independent of which

1 physical allele is encoded as e.g. “A” or “a”. As we add a third  
 2 locus, the epistatic coefficient between mutations  $A \rightarrow a$  and  
 3  $B \rightarrow b$  can be calculated in the presence of alleles C and c at this  
 4 new locus, where these two different genetic backgrounds can  
 5 be viewed as defining parallel faces in sequence space (Figure 1).  
 6 Thus, one way of quantifying the extent of epistasis between  
 7 sites 1 and 2 is by taking the average of the squared local epistatic  
 8 coefficients defined on each of these two parallel faces. Similarly,  
 9 we can quantify epistasis between sites 1 and 3 by averaging the  
 10 local squared epistatic coefficients found on the corresponding  
 11 pair of parallel faces, and we can quantify epistasis between sites  
 12 2 and 3 by averaging the squared local epistatic coefficients from  
 13 the final remaining pair of faces (Figure 1).

14 We can generalize these statistics for larger fitness landscapes  
 15 composed of  $\ell$  sites with  $\alpha$  alleles per site represented by a  
 16 vector  $f$  of length  $\alpha^\ell$ . In particular, epistasis between any pair  
 17 of sites  $i, j$  can be quantified by computing the average squared  
 18 local epistatic coefficient  $\overline{\epsilon_{ij}^2}(f)$ , where the average is taken over  
 19 every possible pair of mutations at those sites and across every  
 20 possible genetic background in which they can be introduced. It  
 21 is easy to see that the average of these quantities  $\overline{\epsilon_{ij}^2}(f)$  across all  
 22 possible pairs of sites corresponds to the previously proposed  
 23 average squared epistatic coefficient  $\overline{\epsilon^2}(f)$  describing the overall  
 24 local smoothness of a fitness landscape (Zhou and McCandlish  
 25 2020) (see Supplement). However, separately averaging these  
 26 squared epistatic coefficients for each pair of sites provides a  
 27 more granular description of epistasis, showing not only how  
 28 much epistasis there is, but whether it is equally distributed  
 29 across different pairs of sites or concentrated within specific  
 30 subsets of sites.

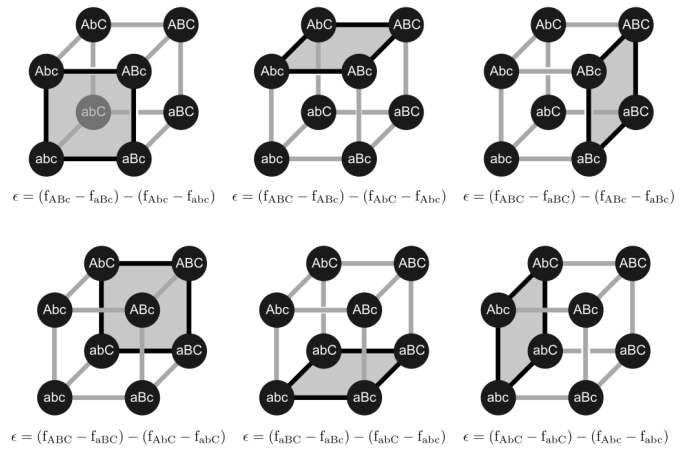
31 For the special case of bi-allelic fitness landscapes, Ferretti  
 32 *et al.* (2016) previously discussed the  $\overline{\epsilon_{ij}^2}$  as part of their proposal  
 33 of the statistic  $\gamma_{i \rightarrow j}$ , which measures the correlation in the ef-  
 34 fects of mutations at site  $j$  before and after mutating site  $i$ . Here,  
 35 we generalize the  $\overline{\epsilon_{ij}^2}$  statistics to multi-allelic landscapes and to  
 36 higher-order epistatic coefficients between any subset of sites  $U$   
 37 (see Supplement). Moreover, we show the relationship between  
 38 these statistics and the total variance explained by interactions  
 39 involving the sites in  $U$  of order equal or higher than  $|U|$ . Specif-  
 40 ically, let  $\text{Var}^{(U)}[f]$  be the total variance explained by  $|U|$ -way  
 41 interactions among the sites in  $U$  as well as interactions of order  
 42 greater than  $|U|$  that also involve all the sites in  $U$  (Reddy and  
 43 Desai 2021; Martí-Gómez *et al.* 2026b). In the supplement, we  
 44 show that this variance  $\text{Var}^{(U)}[f]$  is proportional to the average  
 45 squared local  $|U|$ -way epistatic coefficient  $\overline{\epsilon_{ij}^2}$ :

$$\text{Var}^{(U)}[f] = \left( \frac{\alpha^\ell}{2^{|U|}} \left( \frac{\alpha - 1}{\alpha} \right)^{|U|} \right) \overline{\epsilon_{ij}^2}. \quad (2)$$

46 In addition, we show that the portion of this variance explained  
 47 by genetic interactions of order strictly higher than  $|U|$  is pro-  
 48 portional to the variance across backgrounds in  $|U|$ -way local  
 49 epistatic coefficients between mutations at sites  $U$ :

$$\text{Var}_{k>|U|}^{(U)}[f] = \left( \frac{\alpha^\ell}{2^{|U|}} \left( \frac{\alpha - 1}{\alpha} \right)^{|U|} \right) \text{Var}[\epsilon_U]. \quad (3)$$

50 These results provide a quantitative link between the size and  
 51 variance of local  $|U|$ -way epistatic coefficients involving a set of  
 52 sites  $U$  and the amount of  $|U|$ -way and higher epistatic variance  
 53 explained by that set of sites.



**Figure 1** Local epistatic coefficients for mutations at different pairs of sites in a three-site bi-allelic fitness landscape. Local epistatic coefficients can be grouped into sets corresponding to the interaction of the same pair of mutations on different genetic backgrounds. Geometrically this corresponds to grouping epistatic coefficients on parallel sets of “faces”; here we show each set of parallel faces in its own column, with those corresponding to interactions between sites 1 and 2 in the first column, those corresponding to interactions between sites 2 and 3 in the middle column, and those corresponding to interactions between sites 1 and 3 in the last column.

### Local epistasis regression

54 In the previous section, we defined a set of descriptive statis-  
 55 tics  $\overline{\epsilon_{ij}^2}$  quantifying the magnitude of the local double-mutant  
 56 epistatic coefficients when introducing mutations at a specific  
 57 pair of sites  $i$  and  $j$ . While these statistics can only be evaluated  
 58 for complete fitness landscapes, we can easily use them to con-  
 59 struct a prior distribution for use in inferring complete fitness  
 60 landscapes from noisy and incomplete data. We call the resulting  
 61 Bayesian regression method Local Epistasis Regression.  
 62

In particular, we consider the prior

$$p(f) \propto e^{-\sum_{i<j} a_{ij} \overline{\epsilon_{ij}^2}(f)}, \quad (4)$$

63 where the  $\binom{\ell}{2}$  hyperparameters  $a_{ij} > 0$  penalize the size of the  
 64 local epistatic coefficients between each specific pair of sites  $i$   
 65 and  $j$ . This prior turns out to be an improper Gaussian prior  
 66 that penalizes the size of local genetic interactions as specified  
 67 by the  $a_{ij}$  but which places a flat prior over the non-epistatic  
 68 component of the landscape.  
 69

70 In order to better understand the properties of this prior, and  
 71 in particular why this prior allows epistatic interactions of all  
 72 orders, it is helpful to re-express this new prior within a more  
 73 general class of priors previously described by Zhou *et al.* (2022)  
 74 and Martí-Gómez *et al.* (2026b). This more general prior is a  
 75 Gaussian random field model parametrized by the variance  
 76 explained by epistatic interactions between each possible subset  
 77 of sites  $U$ . The prior as a whole is given by

$$f \sim \mathcal{N} \left( 0, \sum_U \lambda_U P_U \right), \quad (5)$$

78 where  $P_U$  is the projection matrix into the subspace of landscapes  
 79 composed solely of  $|U|$ -way interactions between the sites  $U$

#### 4 Local epistasis regression

1 and we constrain the  $2^\ell$  hyperparameters  $\lambda_U$  to be non-negative,  
 2  $\lambda_U \geq 0$ . This is the most general model in which the covariance  
 3 between two sequences  $x, x'$  only depends on the set of sites  
 4 at which they differ, and the expected variance explained by  
 5  $|U|$ -way genetic interactions between a set of sites  $U$  is given by  
 6  $(\alpha - 1)^{|U|} \lambda_U$  (see Supplement).

7 Turning back to our prior in terms of the  $a_{ij}$  and  $\overline{\epsilon_{ij}^2}(f)$  given  
 8 in Equation 4, in the Supplement we show that this prior is  
 9 equivalent to choosing

$$\lambda_U = \begin{cases} \tilde{\lambda}_\emptyset & \text{if } U = \emptyset \\ \tilde{\lambda}_i & \text{if } U = \{i\} \\ \frac{1}{\alpha^2 \sum_{i < j \in U} a_{ij}} & \text{otherwise,} \end{cases} \quad (6)$$

10 in the limit where  $\tilde{\lambda}_\emptyset = \tilde{\lambda}_1 = \dots = \tilde{\lambda}_\ell \rightarrow \infty$ . Since we assume  
 11  $a_{ij} > 0$ , we see that  $\lambda_U > 0$  for  $|U| \geq 2$  indicating that our  
 12 prior has positive variance for every  $U$  with  $|U| \geq 2$  and can  
 13 thus fit genetic interactions of all orders. In practice, instead of  
 14 taking the limit and working with an improper prior, we retain  
 15  $\tilde{\lambda}_\emptyset$  and  $\tilde{\lambda}_1, \dots, \tilde{\lambda}_\ell$  as additional hyper-parameters that provide  
 16 control over the extent to which the non-epistatic component of  
 17 the model is regularized. Finally, although  $a_{ij}$  penalizes the size  
 18 of  $\overline{\epsilon_{ij}^2}(f)$ , the expected value of  $\overline{\epsilon_{ij}^2}(f)$  under the prior depends  
 19 not only on  $a_{ij}$ , but also on the  $a_{ik}$  and  $a_{jk}$  for all the  $\ell - 2$  other  
 20 sites  $k$  (see Supplement).

21 In order to perform inference under this prior, we further  
 22 assume that the experimental measurements  $y$  for a series of  
 23 sequences  $X$  have normally distributed errors with known ex-  
 24 perimental variance  $y_{var}$  so that we can use standard Gaussian  
 25 process results (Rasmussen and Williams 2008) to compute the  
 26 closed form Gaussian posterior distribution over the complete  
 27 fitness landscape  $f$  with mean

$$\hat{f} = K_{*X}(K_{XX} + D_{var})^{-1}y \quad (7)$$

28 and covariance matrix

$$K - K_{*X}(K_{XX} + D_{var})^{-1}K_{X*}, \quad (8)$$

29 where  $K_{XX}, K_{*X}, K_{X*}$  are submatrices of  $K = \sum_U \lambda_U P_U$  indexed  
 30 by sequences  $X$  and  $*$ , where  $*$  represents all possible sequences,  
 31 and  $D_{var}$  is a diagonal matrix with the known experimental  
 32 variances  $y_{var}$  along the diagonal. Naively evaluating and com-  
 33 puting with these matrices becomes quickly impractical beyond  
 34 a few thousand measurements. However, we can leverage the  
 35 mathematical properties of the covariance matrix to derive ef-  
 36 ficient routines for computing matrix-vector products without  
 37 explicitly constructing the matrices, which can be implemented  
 38 as linear operators and used for solving large linear systems  
 39 using iterative methods in *gpmmap-tools* (see Supplement) (Martí-  
 40 Gómez et al. 2026b).

41 The final remaining issue is how to choose the values for  
 42 the  $1 + \ell + \binom{\ell}{2}$  hyperparameters (specifically,  $\tilde{\lambda}_\emptyset, \tilde{\lambda}_1, \dots, \tilde{\lambda}_\ell$   
 43 and the  $\binom{\ell}{2} a_{ij}$ ). In principle, these hyperparameters can be speci-  
 44 fied based on prior knowledge, for example if certain sites are  
 45 known to interact more strongly with each other such as con-  
 46 tacting residues within a known structure. Here, however, we  
 47 adopt an empirical Bayes approach and assume that the sta-  
 48 tistical properties of the measured sequences generalize to the  
 49 full sequence space, allowing us to estimate these parameters  
 50 directly from the data without needing any a priori knowledge

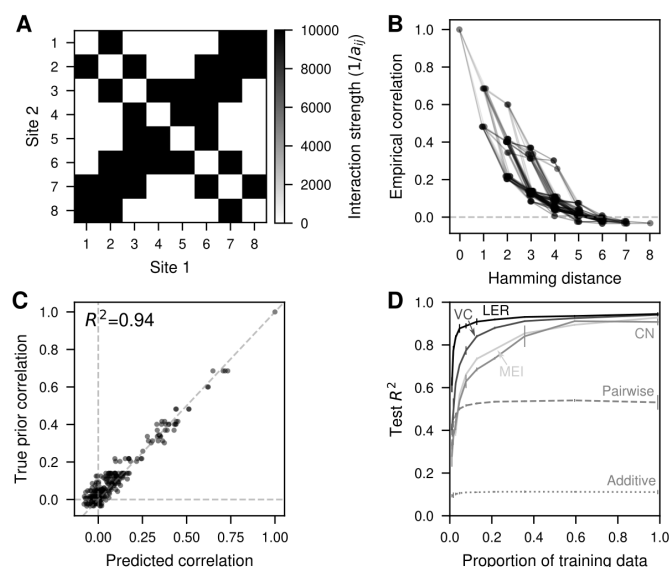
51 about the structure of genetic interactions between sites. Practi-  
 52 cally speaking, because the covariance in fitness between two  
 53 sequences under the prior depends only on the set of sites  $D$  at  
 54 which they differ, we can compute the empirical correlation in  
 55 measured fitness between all pairs of sequences differing exactly  
 56 at a set of sites  $D$ , for each possible subset  $D$ . We then choose the  
 57 hyperparameters so that the correlations implied by the prior  
 58 match these empirical correlations as closely as possible, using a  
 59 procedure known as kernel alignment (Wang et al. 2015; Zhou  
 60 et al. 2022). This structure reduces the kernel alignment problem  
 61 to a  $2^\ell$ -dimensional weighted least squares problem, which can  
 62 be solved efficiently under the constraint that the  $a_{ij}$ , the  $\tilde{\lambda}_i$  and  
 63  $\tilde{\lambda}_\emptyset$  are all positive (see Supplement).

#### Validation on simulated data

64 We first evaluate the performance of Local Epistasis Regression  
 65 using simulated data to illustrate the advantages of this model in  
 66 comparison with previously proposed approaches. Specifically,  
 67 we define a random field model inspired by the expected inter-  
 68 actions between 8 positions forming an RNA helix, where sites  
 69 interact more strongly with their neighboring positions as well  
 70 as with the sites with which they form base-pairs (Figure 2A).  
 71 Figure 2B shows the correlation under this prior between pairs  
 72 of sequences differing at every possible subset of sites  $D$  within  
 73 each Hamming distance class and where lines represent distance  
 74 classes that differ from each other at a single position e.g.  
 75  $D_1 = \{1, 2\}$  and  $D_2 = \{1, 2, 3\}$ . We see that while the correla-  
 76 tion under the prior generally decays with increasing Hamming  
 77 distance  $d = |D|$ , the correlation is not a strict function of  $d$   
 78 and instead varies based on the specific positions at which two  
 79 sequences differ.

80 Having specified a family of random fitness landscapes, we  
 81 then drew a specific fitness landscape from this prior distribution  
 82 and evaluated the performance of the hyperparameter estima-  
 83 tion procedure in recovering the ground truth hyperparameters  
 84 in the presence of limited amount of data. In particular, we  
 85 kept fitness values for only 15% of the sequences, computed the  
 86 empirical fitness correlation between pairs of sequences for each  
 87 class of differences  $D$  and estimated the hyperparameters for  
 88 Local Epistasis Regression using kernel alignment. Figure 2C  
 89 shows that we can accurately infer a prior that approximates the  
 90 true correlation structure of the landscape using only a limited  
 91 amount of data ( $R^2 = 0.94$ ).

92 Since our hyperparameter estimation procedure appeared to  
 93 be working, we then evaluated the performance of the full Local  
 94 Epistasis Regression procedure for different random samples of  
 95 training data ranging from 1% to 99% of the full landscape and  
 96 compare it with the predictive power of previously proposed  
 97 models, including classical models like additive and pairwise  
 98 interaction models, and models that allow interactions of all pos-  
 99 sible orders, including Minimum Epistasis Interpolation (Zhou  
 100 and McCandlish 2020), Empirical Variance Component regres-  
 101 sion (Zhou et al. 2022) and a new implementation of Connected-  
 102 ness regression (Zhou et al. 2025) in *gpmmap-tools* (Martí-  
 103 Gómez et al. 2026b) that uses kernel alignment for hyperparameter  
 104 inference (Figure 2D, see Supplement). Our results show that all  
 105 models that allow higher order interactions can accurately recon-  
 106 struct the true landscape given sufficient amount of data, unlike  
 107 additive and pairwise interaction models. Moreover, by encod-  
 108 ing information about which sites interact with each other and at  
 109 which sites mutations combine more additively, Local Epistasis  
 110 Regression outperformed the other methods for any amount of  
 111



**Figure 2** Validation of Local Epistasis Regression. (A) Heatmap representing the model hyperparameters as  $1/a_{ij}$  for every pair of sites  $i, j$  highlighting the patterns of genetic interactions across sites under the prior. (B) Correlation under the prior for pairs of sequences that differ at each possible subset of sites  $D$  arranged according to the Hamming distance  $d = |D|$ . Each dot represents a single distance class  $D$  and are joined by lines whenever the distance classes differ by a single position from each other. (C) Comparison of the true correlation under the prior and the estimated ones using kernel alignment on a simulated dataset comprising 15% of the sequences from a fitness landscape drawn from the prior. (D) Predictive performance evaluated by the  $R^2$  between the predicted and the true fitness values of held-out test sequences when using different amounts of training data for different models (MEI: Minimum Epistasis Interpolation, VC: Variance Component regression, CN: Connectedness Model regression, LER: Local Epistasis Regression). Predicted values are the maximum a posteriori estimate given by each method, which is equal to the posterior mean  $\hat{f}$ . Error bars represent the standard deviation across 3 different random samples for each fraction of training data.

during the splicing reaction). We start by computing the correlation between pairs of sequences differing at each possible combination of sites (Figure 3A). This analysis shows that the correlation between pairs of sequences depends, not only on the Hamming distance, but also on the specific combination of sites that are different. For example, pairs of sequences differing only at position +6 show a correlation of 0.83, whereas pairs of sequences that differ only at position +2 show a correlation of 0.12. Next, we estimated the hyperparameters of the prior for our Local Epistasis Regression model and show that the estimated prior can recapitulate the observed correlations almost perfectly (Figure 3B). Figure 3C displays the estimated strength of the penalization for local epistatic coefficients for mutations between each possible pair of sites, showing a rich structure of interactions between positions. First, we find that mutations interact more strongly with mutations at neighboring positions, as expected from its recognition mechanism by the U1 snRNA via base-pair complementarity (Wong *et al.* 2018) because the energetic contributions of a basepair to the thermodynamic stability of an RNA helix depend strongly on the adjacent basepairs (Borer *et al.* 1974). There are some exceptions to this general interaction pattern between positions: i) mutations at positions +2 and +5 tend to interact more strongly with each other than expected under the neighbor interaction model, as previously shown in the 5' splice site fitness landscape inferred from their frequencies in the human genome (Chen *et al.* 2021), and ii) the effects of mutations at position +6 are expected to combine more additively with the effects of other mutations. Interestingly, while mutations at position +2 show strong interactions with mutations at -1 and +5, mutations at -1 and +5 are expected to combine nearly additively, a pattern of interaction that cannot be captured by simpler random field models such as the Connectedness Model (Reddy and Desai 2021; Zhou *et al.* 2025). Finally, we evaluated the performance of Local Epistasis Regression when predicting the fitness of held out sequences for different amounts of training data (Figure 3D) and found that, given sufficient training data, predictions are more accurate than under all previously proposed models.

Second, we used data from a high-throughput experiment measuring the translational activity of over 250,000 9-nucleotide sequences at the Shine-Dalgarno sequence in the 5'UTR of the *dmsC* gene in *E. coli* (Kuo *et al.* 2020). As before, the correlation in the measurements between pairs of sequence depends not only on the Hamming distance between them, but also on the specific combination of sites at which they differ (Figure 3E). These empirical correlations can be accurately captured by the Local Epistasis Regression prior under the estimated hyperparameters (Figure 3F). The inferred hyperparameters recapitulate the previously characterized pattern of genetic interactions between mutations at different pairs of sites for this landscape, where sites interact more strongly with other sites within a 4-nucleotide window (Figure 3G) (Martí-Gómez *et al.* 2026b). While the Shine-Dalgarno sequence is also recognized via base-pair complementarity, in this case with the 16S rRNA (Shine and Dalgarno 1975), the observed pattern of interaction between positions can be explained by the ability of the 16S rRNA to bind the target sequence at different registers relative to the start codon (Martí-Gómez *et al.* 2026b). In contrast to results in the Snn1 dataset, Local Epistasis Regression is roughly tied as the best performing model together with Variance Component Regression, although the Variance Component Regression model performs better for very small amounts of training data

1 training data, and exhibited a particularly strong advantage for  
2 low amounts of training data.

### 3 Learning the structure of genetic interactions from empirical data

4 In the previous section, we have shown how Local Epistasis Regression can be used to learn the structure of genetic interactions  
5 from incomplete and noisy data, and then used that information to make more accurate inference of complete combinatorial  
6 fitness landscapes using simulated data. Here, we investigate the performance of Local Epistasis Regression using data from  
7 diverse empirical fitness landscapes. Here, we investigate the performance of Local Epistasis Regression using data from  
8 diverse empirical fitness landscapes. Here, we investigate the performance of Local Epistasis Regression using data from  
9 diverse empirical fitness landscapes. Here, we investigate the performance of Local Epistasis Regression using data from  
10 diverse empirical fitness landscapes. Here, we investigate the performance of Local Epistasis Regression using data from  
11 diverse empirical fitness landscapes. Here, we investigate the performance of Local Epistasis Regression using data from

12 First, we used data from a high-throughput experiment evaluating the functionality of nearly every possible 5' splice site  
13 sequence of the exon 7 in the Snn1 gene context (Wong *et al.* 2018; Zhou *et al.* 2022) (positions -3 through -1 at the exonic region  
14 and +2 through +6 at the intronic region are variable, with position +1 fixed as G, which is necessary for lariat formation  
15 and +2 through +6 at the intronic region are variable, with position +1 fixed as G, which is necessary for lariat formation  
16 and +2 through +6 at the intronic region are variable, with position +1 fixed as G, which is necessary for lariat formation  
17 and +2 through +6 at the intronic region are variable, with position +1 fixed as G, which is necessary for lariat formation

(Figure 3H).

Finally, we applied Local Epistasis Regression to estimate the structure of genetic interactions across sites in two protein fitness landscapes: (i) a complete combinatorial dataset in which nearly all possible combinations of amino acids at four positions in the binding domain of protein G were measured (GB1) (Wu *et al.* 2016), and (ii) a combinatorial landscape in which combinations of five amino acids at seven positions in the core of the SH3 domain of the FYN tyrosine kinase were quantified (FYN-SH3) (Escobedo *et al.* 2025). These datasets also show heterogeneous fitness correlations between pairs of sequences depending on the specific combination of sites at which they differ, which can be effectively captured by our estimated prior distribution (Figure S1A,B,D,E). The inferred parameters suggest that epistatic interactions are concentrated within specific subsets of sites (Figure S1C,F), consistent with patterns observed in the empirical RNA landscapes described above. For GB1, Local Epistasis Regression and Variance Component Regression are the best performing models and perform essentially identically over the whole range of training data sizes (Figure S1D), whereas for the SH3 domain Variance Component Regression performed better across the whole range of training data set sizes (Figure S1H). Thus, in summary, our analysis shows that empirical fitness landscapes often display heterogeneous correlations between pairs of sequences differing by the same number of mutations depending on the specific sites that are mutated, however, the extent to which incorporating this information into our prior distribution improves predictive performance varies across datasets relative to previous methods.

### The fitness landscape of a self-splicing intron

After validating the ability of Local Epistasis Regression to characterize the patterns of genetic interactions between pairs of sites in well-characterized fitness landscapes, we set out to study an 8-nucleotide fitness landscape of a self-splicing type I intron from *Tetrahymena thermophila* (Soo *et al.* 2021) that we have not previously analyzed. In particular, these 8 mutagenized nucleotides can form an extension of the P1 helix that is involved in recognition of the 5' splice site (Figure 4A). This P1 helix must then dissociate to form an alternative P10 helix with the 3' splice site for catalyzing the splicing reaction between the 5' and 3' splice sites (Figure 4A). Thus, we expect a particularly rich structure of genetic interactions across sites, given the heterogeneous structural roles played by these nucleotides at different stages of the splicing reaction.

We start again by computing the correlation in the measured fitness values for pairs of sequences differing at each possible combination of sites. The correlations, as before, depend not only on the number of sites at which two sequence differ, but more specifically on the combination of sites at which they do so, suggesting that different mutations have different effects on the predictability of other mutations (Figure 4B). We next estimate the parameters of the Local Epistasis Regression prior and find that the learned prior can again accurately recapitulate the observed correlations in the data (Figure 4C). The estimated  $a_{ij}$  values indicate that epistatic coefficients are not identically distributed across pairs of mutations at different combinations of sites, but that pairs of mutations at different sites tend to interact more strongly with each other (Figure 4D). The strongest signal comes from mutations at sites 2 and 21 which are by far the most strongly interacting positions, consistent with the P1 helix extension mechanism wherein position 2 basepairs with position

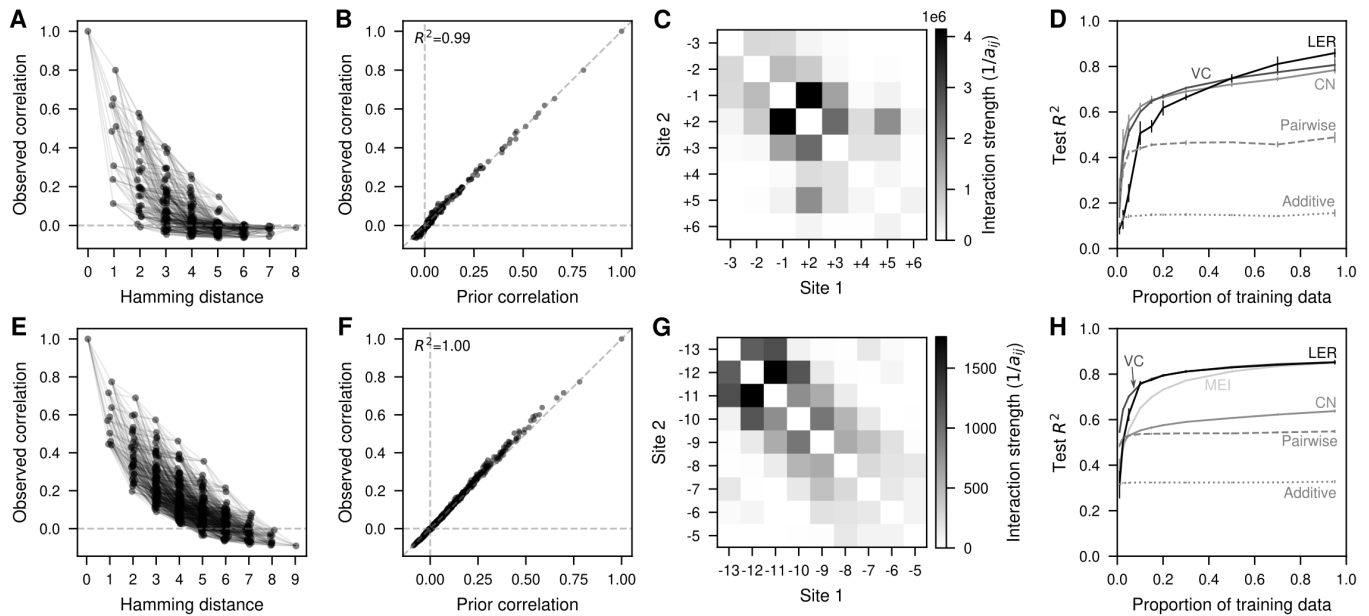
21 (Figure 4A).

Next, we perform inference of the complete fitness landscape under the learned prior while taking into account the estimated experimental error. These estimates accurately recapitulate the measured fitness in 0.5% (327) random held-out sequences (Figure 4E). Importantly, our Gaussian process formulation allows us to obtain uncertainty estimates for the fitness values of the held-out sequences showing good coverage properties, as the 95% posterior credible interval included the measured fitness for 91.1% of the held-out sequences. We found that all the regression models allowing higher-order interactions made similar predictions and had similar predictive performance (Figure S2A), and that in particular Local Epistasis Regression and Variance Component Regression both produced very similar fits (Figure S2B) and predictions for the set of held-out sequences (Figure S2C).

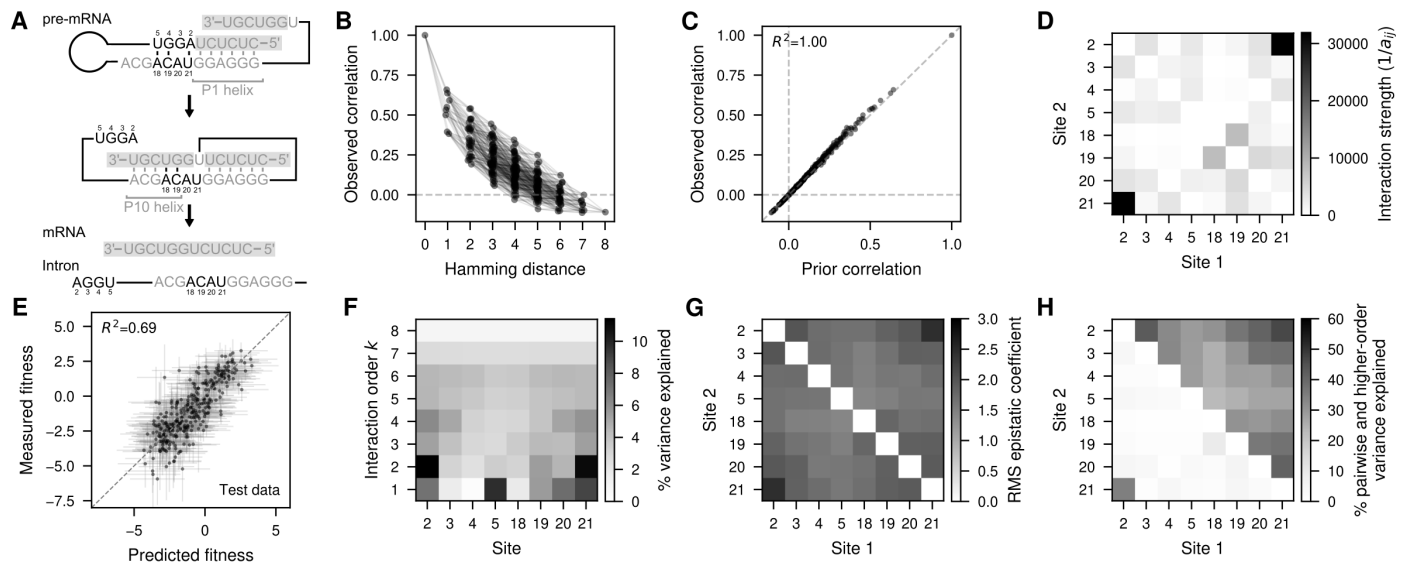
Next, we further explore the structure of epistasis in this dataset by computing a series of informative summary statistics from the maximum a posteriori (MAP) reconstruction of the complete landscape. First, we compute the percentage of variance explained by epistatic interactions of every possible order (Zhou *et al.* 2022) and find that the reconstructed landscape is highly epistatic, with only 42.3% of the variance explained by the additive component, 22.2% by pairwise interactions, and the remaining 35.4% explained by higher-order genetic interactions. To investigate how much each site contributes to interactions of different orders, we compute the percentage of variance explained by genetic interactions of order  $k$  that involve each site  $i$  (Figure 4F) (Martí-Gómez *et al.* 2026b). These statistics reveal substantial heterogeneity across sites. Some sites contribute little across all interaction orders, such as sites 4 and 18. Other sites contribute more strongly through their additive effects, such as position 5, whereas sites like 2 and 21 contribute primarily through pairwise and higher-order interactions. Next, we characterize the structure of genetic interactions between sites by computing the root mean square local double-mutant epistatic coefficient  $\sqrt{\epsilon_{ij}^2}$  for every pair of sites (Figure 4G). These statistics reveal widespread epistatic interactions across pairs of sites, particularly among sites 2, 3, 20 and 21, and within the groups of positions 2–5 and 18–21, consistent with the regularization parameters (Figure 4D). However, as shown by Eq. 2, these quantities reflect contributions from both pairwise and higher-order interactions. To disentangle these contributions, we compute the percentage of pairwise ( $\text{Var}_{k=2}^{\{(i,j)\}}[\hat{f}] = \text{Var}[P_{\{(i,j)\}}\hat{f}]$ ) and higher-order variance ( $\text{Var}_{k>2}^{\{(i,j)\}}[\hat{f}]$ ) explained by interactions involving each pair of sites (Figure 4H). This analysis shows that a large fraction of pairwise epistatic variance is explained by interactions between positions 2 and 21 alone (Figure 4A, lower triangular part), whereas higher-order epistatic variance is more broadly distributed across sites, indicating that local double mutant epistatic coefficients for essentially all pairs of sites vary substantially across genetic backgrounds (Eq.3).

### Structure of the self-spliced intron fitness landscape

In order to better understand the qualitative structure of this fitness landscape, we apply a visualization method for producing low-dimensional representations of fitness landscapes where distances between genotypes in the representation reflect the expected waiting time to evolve from one genotype to another under a model of molecular evolution that includes mutation, selection, and drift (McCandlish 2011). Applying this method as implemented in the software package *gpmmap-tools* resulted



**Figure 3** Application of Local Epistasis Regression to estimate the structure of epistatic interactions across sites in empirical fitness landscapes. (A,E) Correlation in the measured fitness values for pairs of sequences differing at each possible subset of sites  $D$  arranged according to the Hamming distance  $d = |D|$ . Each dot represents a single distance class  $D$  and are joined by lines whenever the distance classes differ by a single position from each other. (B,F) Comparison of the observed correlation values in the data and the values under the estimated prior ones using Local Epistasis Regression for every possible distance class  $D$  (each dot represents a different  $D$ ). Correlations were estimated using 80% of the data for training. (C,G) Heatmap representing the inferred model hyperparameters as  $1/a_{ij}$  for every pair of sites  $i, j$  highlighting the patterns of genetic interactions across sites under the prior. (D,H) Predictive performance evaluated by the  $R^2$  between the predicted and the measured fitness of held-out test sequences when using different amounts of training data for different models (MEI: Minimum Epistasis Interpolation, VC: Variance Component regression, CN: Connectedness Model regression, LER: Local Epistasis Regression). Predicted values are the maximum a posteriori estimate given by each method, which is equal to the posterior mean  $\hat{f}$ . Error bars represent the standard deviation across 3 different random samples for each fraction of training data. Each row represents a fitness landscape: Snn1 exon 7 5' splice site (A,B,C,D); dmsC Shine-Dalgarno sequence (E,F,G,H).



**Figure 4** Inference and summary statistics of the fitness landscape of a self-spliced intron. (A) Schematic representation of the molecular mechanism of self-splicing in the model intron from *Tetrahymena thermophila* (Soo *et al.* 2021). Positions considered in the fitness landscape are highlighted in black and are numbered by their relative position in the intron. Other relevant sequences that are fixed in the background are shown in grey. Exonic sequences are shown with gray background. (B) Correlation in the measured fitness values for pairs of sequences differing at each possible subset of sites  $D$  arranged according to the Hamming distance  $d = |D|$ . Each dot represents a single distance class  $D$  and are joined by lines whenever the distance classes differ by a single position from each other. (C) Comparison of the observed correlation values in the data and the values under the estimated prior using Local Epistasis Regression. One correlation is shown for each possible set of positions  $D$  where two sequences may differ (each dot represents a different  $D$ ). (D) Heatmap representing the inferred model hyperparameters as  $1/a_{ij}$  for every pair of sites  $i, j$  highlighting the patterns of genetic interactions across sites under the prior. (E) Measured values for held-out test sequences versus Local Epistasis Regression predictions. Horizontal error bars represent the 95% credible interval, whereas vertical error bars correspond to the 95% confidence interval under each measurement's variance. (F) Heatmap representing the percentage of variance in the maximum a posteriori Local Epistasis Regression reconstruction explained by interactions of order  $k$  involving each position. (G) Root mean squared local double-mutant epistatic coefficient magnitude between mutations at each possible pair of predictions for the maximum a posteriori reconstruction. The plot indicates that relatively large local epistatic coefficients occur between mutations at essentially all pairs of positions. (H) Heatmap representing the percentage of variance in the maximum a posteriori reconstruction explained by pairwise (lower triangle) and higher-order (upper triangle) interactions that is explained by interactions involving pairs of positions.

1 in Figure 5A (see Methods, [Martí-Gómez et al. \(2026b\)](#)). In the  
2 previous section we noted a particularly strong pattern of both  
3 pairwise and higher-order interactions involving positions 2 and  
4 21 (Figure 4F), and we see that the visualization in Figure 5A  
5 largely separates sequences based on the nucleotides present at  
6 this pair of positions. Because the axes in such visualizations  
7 tend to highlight key barriers that make it difficult for a popula-  
8 tion to diffuse from one area of sequence space to another, we  
9 examined the mean fitness conferred by different combinations  
10 of alleles at these two sites for potential allelic incompatibilities  
11 (Figure 5B), revealing a consistent pattern wherein having both  
12 positions 2 and 21 occupied by pyrimidines (or also  $C_2A_{21}$ ) re-  
13 sults in low mean fitness. This observation suggests that the  
14 P1 helix extension by base-pairing of positions 2 and 21 is not  
15 strictly necessary for functionality.

16 Because the visualizations tend to spread apart high-fitness  
17 sequences that are separated by these types of incompatibil-  
18 ities, they are also useful for understanding how mutational  
19 effects and local epistatic coefficients vary across sequence space.  
20 Figures 5D and E show how  $G_2C$  and  $C_{21}G$  mutations are dele-  
21 terious in the  $G_2C_{21}$  cluster on the left-hand (negative) side of  
22 Diffusion Axis 1, but the same mutations become advantageous  
23 as one moves to the right-hand (positive) side. In fact, looking  
24 across all 8 positions, Figure S4 shows a widespread tendency  
25 for allelic preferences to change in a coherent manner across  
26 different clusters of sequences. Similarly, we can see how local  
27 double-mutant epistatic coefficients change in different regions  
28 of sequence space. For example, in the  $C_2G_{21}$  cluster on the right-  
29 hand side of Figure 5A,  $G_3C$  and  $C_{20}G$  are strongly negatively  
30 epistatic (Figure 5F) whereas in many other backgrounds such  
31 as  $U_2A_{21}$  background (Figure 5C), their interaction is strongly  
32 positive. Likewise the local double-mutant epistatic coefficient  
33 between  $A_3U$  and  $U_{20}A$  are strongly negative in the  $C_2G_{21}$  clus-  
34 ter but typically neutral or positive elsewhere (Figure 5G). Note  
35 that Figures 5F and 5G provide a simple illustration of the geo-  
36 metry of a 4th-order interaction, as we can see that the epistatic  
37 interactions at position 3 and 20 differ across the clusters of  
38 sequences defined by the base identities at positions 2 and 21.

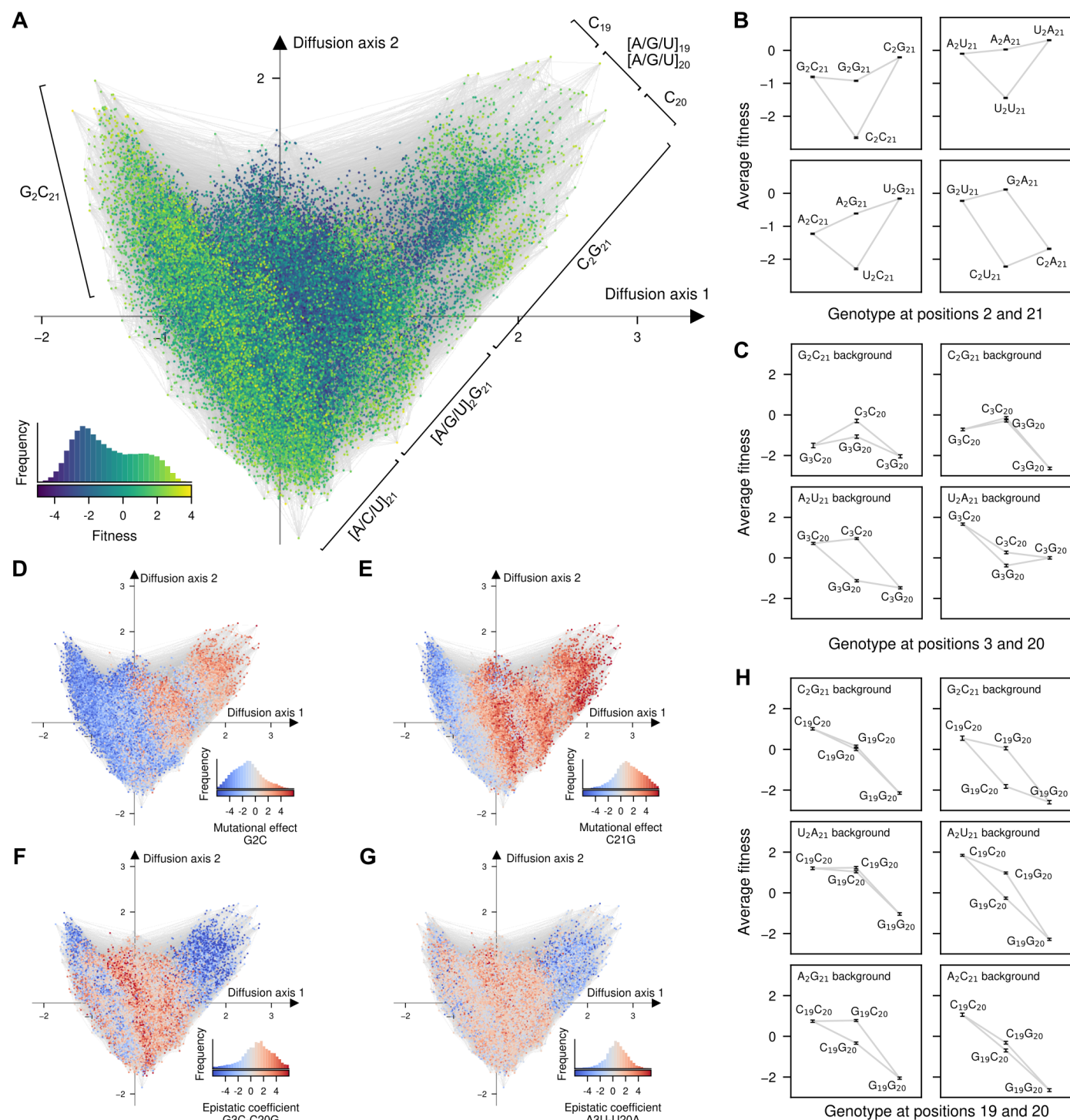
39 The visualization also separates sequences within the  $C_2G_{21}$   
40 background into three clusters, depending on the presence of C  
41 at positions 19 and 20 (Figure 5A). Moreover, Figure S4 shows  
42 that while  $C_{19}$  and  $C_{20}$  are generally preferred across the land-  
43 scape, consistent with base pairing with cognate bases in the  
44 P10 helix in the second step of the splicing reaction (Figure 4A),  
45 the effects of alleles C and G at these two positions are strongly  
46 dependent on genetic background. To better understand these  
47 dependencies, we examined the average fitness of the four com-  
48 binations of C and G alleles at positions 19 and 20 across differ-  
49 ent genetic contexts (Figure 5H). We find that in some genetic  
50 contexts, like  $C_2G_{21}$ ,  $U_2A_{21}$ , and  $C_2A_{21}$ , the base-pair-breaking  
51 mutations  $C_{19}G$  and  $C_{20}G$  are on average neutral or only mildly  
52 deleterious individually, but become substantially more dele-  
53 terious when combined together. In contrast, in other genetic  
54 contexts, such as  $G_2C_{21}$ ,  $A_2U_{21}$ , and  $A_2C_{21}$ , these mutations are  
55 more strongly deleterious individually and combine more addi-  
56 tively. These epistatic interactions limit the accessible mutational  
57 paths between  $C_{19}G_{20}$  and  $G_{19}C_{20}$  sequences, and explain why  
58 these clusters appear separated in the visualization (Figure 5A).  
59 The background dependence of this interaction provides an addi-  
60 tional illustration of how higher-order interactions in this system  
61 can be largely understood in terms of the identities at positions  
62 2 and 21, and how the alleles at these two positions modulate

epistatic interactions between the other sites.

## Discussion

63  
64  
65 In this work, we propose a new method for inferring empiri-  
66 cal fitness landscapes from experimental data. Our method is  
67 based on the idea that genetic interactions are not uniformly  
68 distributed across sites; rather, some sites tend to interact more  
69 strongly with specific other sites. This information about sets of  
70 sites exhibiting stronger epistatic interactions can be extracted  
71 from how correlations between measured fitness values depend  
72 on the specific combinations of sites at which sequences differ.  
73 We then incorporate this information into a prior distribution  
74 defined over all possible fitness landscapes for Bayesian infer-  
75 ence of complete combinatorial landscapes from incomplete  
76 and noisy data. Importantly, although our prior allows the in-  
77 corporation of different degrees of regularization for epistatic  
78 coefficients between different pairs of sites, it still allows these co-  
79 efficients to vary across genetic backgrounds, thereby preserving  
80 the benefits of other methods that contain genetic interactions  
81 of all orders ([Zhou and McCandlish 2020](#); [Zhou et al. 2022](#)). Ap-  
82 plying this method to three experimental datasets revealed that  
83 mutations do tend to interact more strongly with other specific  
84 mutations in empirical landscapes (Figures 3C,F, S1C,G and 4D),  
85 often reflecting known molecular interactions between positions  
86 (Figures 4A,D). Finally, we performed an in-depth analysis of  
87 the inferred landscape for a self-splicing intron, where we found  
88 that the nucleotide identities at positions 2 and 21, which can  
89 potentially form an extension of the P1 helix, determine the na-  
90 ture of genetic interactions at many other pairs of positions. This  
91 shows how higher-order interactions, which may at first seem  
92 mysterious, can be understood in terms of local pairwise inter-  
93 actions that change coherently as one moves from one region of  
94 the fitness landscape to another.

95 Our method relies on a simple summary statistic that quanti-  
96 fies the overall strength of local genetic interactions between a  
97 given pair of mutations across across all genetic backgrounds.  
98 Specifically, we compute the average squared epistatic coeffi-  
99 cient between mutations at every possible set of sites. While  
100 this statistic was previously derived in the context of the  $\gamma_{i \rightarrow j}$   
101 measure ([Ferretti et al. 2016](#)), here we show that it can be ex-  
102 pressed as a quadratic form with a positive semi-definite matrix  
103 that admits a Kronecker product decomposition into  $\ell$  smaller  
104 matrices. This allowed us to derive (i) the relationship between  
105 the total epistatic variance of all orders involving any pair of  
106 sites  $i$  and  $j$  ([Reddy and Desai 2021](#); [Martí-Gómez et al. 2026b](#))  
107 and the average squared size of local double-mutant epistatic  
108 coefficients defined by pairs of mutations at sites  $i$  and  $j$ , and (ii)  
109 the relationship between the higher order epistatic variance and  
110 the variance in these epistatic coefficients across genetic back-  
111 grounds, clarifying how local epistatic effects and their context  
112 dependence connect to the global variance decomposition of  
113 fitness landscapes. Interestingly, while the ability to conduct  
114 this Kronecker factorization enables efficient computation of  
115 these summary statistics in combinatorially complete fitness  
116 landscapes ([Martí-Gómez et al. 2026b,a](#)), it also opens up the  
117 possibility of calculating these statistics for astronomically large  
118 fitness landscapes for the maximum a posteriori estimate under  
119 a certain class of Gaussian process models ([Petti et al. 2025](#)),  
120 which would allow the incorporation of higher-order interac-  
121 tions in applications such as protein contact prediction ([Marks  
122 et al. 2012](#)) and 3D-structure inference from mutagenesis exper-  
123 iments ([Schmiedel and Lehner 2019](#)). Finally, the Kronecker



**Figure 5** Visualization and characterization of the fitness landscape of a self-spliced intron. (A) Visualization of the inferred fitness landscape using Local Epistasis Regression. Every dot represents one of the possible  $4^8$  possible sequences and is colored according to the predicted fitness. The inset represents the phenotypic distribution along with their corresponding color in the map. Sequences are laid out according to the first two Diffusion axes and dots are plotted in order according to Diffusion axis 3. (B,C,H) Diagrams representing average fitness for different subsets of sequences at positions 2 and 21 (B), for combinations of G and C alleles at positions 3 and 20 across different genetic contexts at positions 2 and 21 (C) and for combinations of G and C alleles at positions 19 and 20 across different genetic contexts at positions 2 and 21 (H). Indicated subsets of sequences are arranged along the x-axis according to Hamming distance from the leftmost subset. Error bars represent the 95% posterior credible intervals for these average fitness values. (D-G) Visualization of the inferred fitness landscape, as shown in (A), where nodes are colored by the mutational effect of G2C (D) and C21G (E), and the local epistatic coefficients between mutations G3C and C20G (F) and between A3U and U20A (G) when introduced in every possible genetic background throughout the landscape. The inset represents the distribution of the specific mutational effects or epistatic coefficients along with their corresponding color in the map.

1 factorization form makes it straightforward to generalize these  
2 statistics in two directions: (i) computing squared epistatic co-  
3 efficients between specific mutations within the same pair of  
4 sites (see Supplement) and (ii) averaging of the squared coeffi-  
5 cients across different genetic backgrounds or over backgrounds  
6 drawn from site-factorizable probability distributions (rather  
7 than from the uniform distribution as we have done here).

8 As in previous work (Zhou and McCandlish 2020; Chen *et al.*  
9 2021; Zhou *et al.* 2022; Reddy and Desai 2021; Zhou *et al.* 2025),  
10 we use our mathematical understanding of this summary statisti-  
11 c to define random field models of fitness landscapes, which  
12 can be used as prior distributions for Gaussian process inference  
13 of empirical landscapes from high-throughput experimental  
14 data. Our prior allows the same mutation to have different  
15 effects on the predictability of different mutations, whereas pre-  
16 vious approaches assign all mutations (Zhou and McCandlish  
17 2020; Chen *et al.* 2021) or each individual mutation (Reddy and  
18 Desai 2021; Zhou *et al.* 2025) a constant effect on the predictability  
19 of other mutations, preventing them from encoding structured  
20 interaction patterns between sites. Despite finding widespread  
21 evidence that specific pairs of sites tend to interact more than  
22 others, incorporating this realism into our prior did not always  
23 increase predictive performance. One potential reason for this is  
24 that while more flexibly encoding differences between pairs of  
25 sites, our model makes assumptions about the form of higher-  
26 order epistasis that may be more appropriate for some datasets  
27 than others. A second potential reason for this variable perfor-  
28 mance concerns our kernel alignment procedure for choosing the  
29 hyperparameters. We had previously seen a modest improve-  
30 ment in model performance for Empirical Variance Component  
31 Regression when the hyperparameters were chosen by evidence  
32 maximization rather than kernel alignment (Zhou *et al.* 2025),  
33 suggesting that the two modes of inference should generally  
34 have similar performance. However, here we see that the con-  
35 nectedness model fit by kernel alignment performed far worse  
36 on the GB1 and Smn1 datasets than our previous implementa-  
37 tion using evidence maximization (Zhou *et al.* 2025), suggesting  
38 that the method of hyperparameter optimization might have a  
39 greater impact than previously thought. Unfortunately, unlike  
40 kernel alignment, which we were able to implement efficiently  
41 using the fact that the covariance under our model only depends  
42 on the subset of sites at which two sequences differ, implement-  
43 ing evidence maximization for Local Epistasis Regression as  
44 in Zhou *et al.* (2025) is more challenging due to the need to  
45 calculate and combine  $2^\ell$  kernels.

46 A fundamental aspect of fitness landscapes is their depen-  
47 dence on the environment. Increasingly, datasets measure fitness  
48 across multiple environmental conditions, necessitating models  
49 that allow mutational effects and epistatic interactions to vary  
50 across environments (Nguyen Ba *et al.* 2022; N’Guessan *et al.*  
51 2025; Bakerlee *et al.* 2022; Soo *et al.* 2021; Ishigami *et al.* 2024).  
52 Our framework naturally extends to this setting by treating en-  
53 vironmental conditions as additional loci, enabling the inference  
54 of priors in which mutational effects depend jointly on genetic  
55 background and environment. Importantly, our prior can learn  
56 and encode that mutations tend to change by different mag-  
57 nitudes when introducing an additional mutation compared  
58 to when changing environments, effectively allowing different  
59 prior variances for gene-by-gene and gene-by-environment in-  
60 teractions.

61 Our method shares several limitations with previous Gaus-  
62 sian process approaches (Zhou and McCandlish 2020; Zhou *et al.*

2022, 2025). First, it does not explicitly model non-specific or  
global epistasis (Otwinowski *et al.* 2018; Domingo *et al.* 2019),  
but instead relies on specific epistatic interactions to fit these  
global dependencies, potentially limiting the interpretability of  
the estimated hyperparameters as the structure of specific ge-  
netic interactions between sites. Second, by taking an empirical  
Bayes approach, in which we first estimate the parameters of the  
prior and then use those parameters for inference of the fitness  
landscape, we are not taking into account potential uncertainty  
in the estimation of these parameters. Our new model may  
be more sensitive to this limitation than previous approaches  
because of the larger number of hyperparameters that need to  
be estimated (Zhou and McCandlish 2020; Zhou *et al.* 2022).  
Third, our current implementation takes advantage of the math-  
ematical structure of the covariance matrix to allow efficient  
computation of the posterior distribution (Equations 7 and 8)  
without explicitly building the covariance matrix for complete  
fitness landscapes. While this trick allows us to compute the pos-  
terior distribution for fitness landscapes containing hundreds of  
thousands of sequences, it is limited to sequences of relatively  
short length (about 5 amino acids, 12 nucleotides or 24 bi-allelic  
loci). Moreover, the need to evaluate  $2^\ell$  different kernels, even  
if the number of hyperparameters is much lower, hinders the  
applicability even under GPU-accelerated frameworks for scal-  
able Gaussian process inference (Gardner *et al.* 2018; Wang *et al.*  
2019) that facilitated the application of previous models to fitness  
landscapes defined over longer sequences (Zhou *et al.* 2025).

## Materials and methods

### Fitness landscape of a 5' splice site sequence

Data reported by Wong *et al.* (2018) was processed as previously  
reported (Zhou *et al.* 2022, 2025). Briefly, we assumed a log-  
normal distribution of enrichment ratios across 1-7 replicates,  
for each different 5' splice site sequence  $x$ . The bias corrected  
geometric mean of the enrichment ratio was used as an esti-  
mate of the median enrichment ratio when the enrichment ratio  
was strictly positive for all replicates. Otherwise, the median  
of enrichment scores was used to estimate the phenotype  $y_x$ .  
Sequence-specific variance  $y_{x,var}$  was estimated as indicated be-  
low, where  $s_x^2$  is the sample variance of the log-enrichment ratios  
if all replicates were strictly positive and were measured in at  
least two samples or the median of all  $s_x^2$  for sequences  $x$  with at  
least two replicate measurements:

$$y_{x,var}^2 = \left( e^{s_x^2} - 1 \right) e^{2y_x + s_x^2}, \quad (9)$$

see (Zhou *et al.* 2022) for more details.

### Fitness landscape of the Shine-Dalgarno sequence

Data reported by Kuo *et al.* (2020) was processed as previously re-  
ported (Martí-Gómez *et al.* 2026b). Briefly, fitness was estimated  
as the mean log(GFP) for 257,565 measured sequences with a  
common measurement variance of  $s^2 = 0.058$  using genotypes  
measured across all three experimental replicates. The squared  
standard error for each genotype  $i$  was computed by dividing  
this observed experimental variance  $s^2$  by the number of repli-  
cates  $n_x$  in which each sequence was measured ( $y_{x,var} = s^2/n_x$ ).

### Fitness landscape of protein G binding domain

Data was processed as previously described (Zhou and McCan-  
dlish 2020; Zhou *et al.* 2022). Briefly, we used the number of

1 sequencing reads for each sequence  $x$  in the input sample ( $c_x^{input}$ )  
2 and in the selected sample ( $c_x^{sel}$ ) reported in (Wu *et al.* 2016) to  
3 estimate the log-enrichment ratio relative to the wild-type se-  
4 quence  $y_x$  as a measure of the binding strength. Moreover, we  
5 estimated the error variance  $y_{x,var}$  of this estimate (Rubin *et al.*  
6 2017):

$$y_x = \log \left( \frac{c_x^{sel} + 0.5}{c_x^{input} + 0.5} \right) - \log \left( \frac{c_{wt}^{sel} + 0.5}{c_{wt}^{input} + 0.5} \right) \quad (10)$$

$$y_{x,var} = \frac{1}{c_x^{input} + 0.5} + \frac{1}{c_x^{sel} + 0.5} + \frac{1}{c_{wt}^{input} + 0.5} + \frac{1}{c_{wt}^{sel} + 0.5}. \quad (11)$$

## 7 Fitness landscape of the FYN protein SH3 domain

8 Data was used as reported by the original study (Escobedo *et al.*  
9 2025). Specifically, we downloaded the processed data from GEO  
10 (GSE266299) used the scaled relative fitness measurements and  
11 the reported experimental errors for our downstream analysis.

## 12 Fitness landscape of a self-splicing intron

13 The number of reads for each variant across six-replicates in the  
14 presence and absence of Kanamycin was collected at 30° from  
15 the original study (Soo *et al.* 2021). Following Soo *et al.* (2021),  
16 fitness  $y_x$  for each sequence  $x$  was estimated as the  $\log_2$ (Fold  
17 change) between Kanamycin treated and control samples using  
18 PyDESeq2 (Muzellec *et al.* 2023) without shrinkage towards zero,  
19 and squared standard errors were kept as measure of experi-  
20 mental variance  $y_{x,var}$  for downstream analysis. Local Epistasis  
21 Regression was used to estimate the complete combinatorial  
22 fitness landscape taking into account the experimental errors for  
23 each of the sequence using all available data except 0.5% (327)  
24 sequences. We computed the posterior mean and variance for  
25 these 327 sequences to evaluate the performance of the model in  
26 unobserved sequences. Using the posterior mean for the com-  
27 plete fitness landscape, we computed the variance explained by  
28 epistatic interactions of every possible order for each site and be-  
29 tween pair of sites using *gpmmap-tools* (Martí-Gómez *et al.* 2026b).  
30 We then generated a low-dimensional representation in which  
31 distances between pairs of sequences reflect the expected time to  
32 evolve from one sequence to the other (McCandlish 2011) under  
33 an evolutionary model in the weak mutation regime as imple-  
34 mented in *gpmmap-tools* (Martí-Gómez *et al.* 2026b). We generated  
35 visualization coordinates under different strengths of selection  
36 by choosing different values for the expected fitness under long-  
37 term mutation-selection-drift (i.e. expected fitness at stationarity,  
38 Figure S3). A long-term expected fitness of 1.6 (corresponding  
39 to the 87% percentile in the distribution of fitness values) was  
40 used for the final visualization.

## 41 Data and code availability

42 The methods presented in this work have been implemented in  
43 *gpmmap-tools* v0.4.2 (Martí-Gómez *et al.* 2026b), an open-source  
44 library available at <https://github.com/cmarti/gpmmap-tools>.  
45 Code and data to reproduce the analyses presented in this paper  
46 are available at <https://github.com/cmarti/deltaU>.

## 47 Funding

48 CMG and DMM were supported by the US National Institutes  
49 of Health (NIH) grant R35GM133613 and additional funding

50 from the Simons Center for Quantitative Biology at Cold Spring  
51 Harbor Laboratory. This work was performed with assistance  
52 from the NIH Grant S10OD028632.

## 53 Conflicts of interest

54 The authors declare no conflicts of interest.

## 55 Literature cited

- 56 Agarwal V, Inoue F, Schubach M, Penzar D, Martin BK, Dash  
57 PM, Keukeleire P, Zhang Z, Sohota A, Zhao J *et al.* 2025. Mas-  
58 sively parallel characterization of transcriptional regulatory  
59 elements. *Nature*. pp. 1–10. Publisher: Nature Publishing  
60 Group.
- 61 Aguirre L, Hendelman A, Hutton SF, McCandlish DM, Lippman  
62 ZB. 2023. Idiosyncratic and dose-dependent epistasis drives  
63 variation in tomato fruit size. *Science*. 382:315–320.
- 64 Aita T, Husimi Y. 1998. Fitness Landscape of a Biopolymer Partic-  
65 ipating in a Multi-step Reaction. *Journal of Theoretical Biology*.  
66 191:377–390.
- 67 Baeza-Centurion P, Miñana B, Schmiedel JM, Valcárcel J, Lehner  
68 B. 2019. Combinatorial Genetics Reveals a Scaling Law for  
69 the Effects of Mutations on Splicing. *Cell*. 176:549–563.e23.  
70 Publisher: Cell Press.
- 71 Bakerlee CW, Nguyen Ba AN, Shulgina Y, Rojas Echenique JI,  
72 Desai MM. 2022. Idiosyncratic epistasis leads to global fitness-  
73 correlated trends. *Science*. 376:630–635.
- 74 Bank C. 2022. Epistasis and Adaptation on Fitness Landscapes.  
75 *Annual Review of Ecology, Evolution, and Systematics*. 53:457–  
76 479. Publisher: Annual Reviews.
- 77 Bank C, Matuszewski S, Hietpas RT, Jensen JD. 2016. On the  
78 (un)predictability of a large intragenic fitness landscape. *Pro-  
79 ceedings of the National Academy of Sciences*. 113:14085–  
80 14090.
- 81 Bendixsen DP, Collet J, Østman B, Hayden EJ. 2019. Genotype  
82 network intersections promote evolutionary innovation. *PLoS  
83 Biology*. 17. Publisher: Public Library of Science.
- 84 Bonde MT, Pedersen M, Klausen MS, Jensen SI, Wulff T, Har-  
85 rison S, Nielsen AT, Herrgård MJ, Sommer MO. 2016. Pre-  
86 dictable tuning of protein expression in bacteria. *Nature Meth-  
87 ods*. 13:233–236.
- 88 Borer PN, Dengler B, Tinoco Jr I, Uhlenbeck OC. 1974. Stability of  
89 ribonucleic acid double-stranded helices. *Journal of molecular  
90 biology*. 86:843–853.
- 91 Bryant DH, Bashir A, Sinai S, Jain NK, Ogden PJ, Riley PF,  
92 Church GM, Colwell LJ, Kelsic ED. 2021. Deep diversifica-  
93 tion of an aav capsid protein by machine learning. *Nature  
94 Biotechnology*. 39:691–696.
- 95 Chattopadhyay G, Papkou A, Wagner A. 2025. The fitness land-  
96 scape of the E.coli lac operator is highly rugged in two differ-  
97 ent environments. *bioRxiv*. .
- 98 Chen Wc, Zhou J, Sheltzer JM, Kinney JB, Mccandlish DM. 2021.  
99 Field-theoretic density estimation for biological sequence  
100 space with applications to 5 splice site diversity and aneu-  
101 ploidy in cancer. *Proc. Natl. Acad. Sci. USA*. .
- 102 Chou HH, Chiu HC, Delaney NF, Segrè D, Marx CJ. 2011. Dimin-  
103 ishing Returns Epistasis Among Beneficial Mutations Deceler-  
104 ates Adaptation. *Science*. 332:1190–1192. Publisher: American  
105 Association for the Advancement of Science.
- 106 Crawford L, Zeng P, Mukherjee S, Zhou X. 2017. Detecting epista-  
107 sis with the marginal epistasis test in genetic mapping studies  
108 of quantitative traits. *PLOS Genetics*. 13:e1006869. Publisher:  
109 Public Library of Science.

- 1 Dasari K, Somarelli JA, Kumar S, Townsend JP. 2021. The somatic  
2 molecular evolution of cancer: Mutation, selection, and epista-  
3 sis. *Progress in Biophysics and Molecular Biology*. 165:56–65.
- 4 de Boer CG, Vaishnav ED, Sadeh R, Abeyta EL, Friedman N,  
5 Regev A. 2020. Deciphering eukaryotic gene-regulatory logic  
6 with 100 million random promoters. *Nature Biotechnology*.  
7 38:56–65. Publisher: Nature Publishing Group.
- 8 De Los Campos G, Hickey JM, Pong-Wong R, Daetwyler HD,  
9 Calus MPL. 2013. Whole-Genome Regression and Prediction  
10 Methods Applied to Plant and Animal Breeding. *Genetics*.  
11 193:327–345.
- 12 De Visser JAG, Krug J. 2014. Empirical fitness landscapes and the  
13 predictability of evolution. *Nature Reviews Genetics*. 15:480–  
14 490. Publisher: Nature Publishing Group.
- 15 Domingo J, Baeza-Centurion P, Lehner B. 2019. The causes and  
16 consequences of genetic interactions (epistasis). *Annual Re-  
17 view of Genomics and Human Genetics*. 20:433–460. Publisher:  
18 Annual Reviews Inc.
- 19 Domingo J, Diss G, Lehner B. 2018. Pairwise and higher-order  
20 genetic interactions during the evolution of a tRNA. *Nature*.  
21 558:117–121.
- 22 Dwivedi SL, Heslop-Harrison P, Amas J, Ortiz R, Edwards D.  
23 2024. Epistasis and pleiotropy-induced variation for plant  
24 breeding. *Plant Biotechnology Journal*. 22:2788–2807. eprint:  
25 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/pbi.14405>.
- 26 Escobedo A, Voigt G, Faure AJ, Lehner B. 2025. Genetics, ener-  
27 getics, and allostery in proteins with randomized cores and  
28 surfaces. *Science*. 389:eadq3948. Publisher: American Associa-  
29 tion for the Advancement of Science.
- 30 Evfratov SA, Osterman IA, Komarova ES, Pogorelskaya AM,  
31 Rubtsova MP, Zatsepin TS, Semashko TA, Kostryukova ES,  
32 Mironov AA, Burnaev E *et al.* 2017. Application of sorting  
33 and next generation sequencing to study 5'-UTR influence  
34 on translation efficiency in *Escherichia coli*. *Nucleic Acids  
35 Research*. 45:3487–3502.
- 36 Faure AJ, Lehner B, Miró Pina V, Serrano Colome C, Weghorn  
37 D. 2024a. An extension of the walsh-hadamard transform  
38 to calculate and model epistasis in genetic landscapes of ar-  
39 bitrary shape and complexity. *PLoS Computational Biology*.  
40 20:e1012132.
- 41 Faure AJ, Martí-Aranda A, Hidalgo-Carcedo C, Beltran A,  
42 Schmiedel JM, Lehner B. 2024b. The genetic architecture of  
43 protein stability. *Nature*. 634:995–1003. Publisher: Nature Pub-  
44 lishing Group.
- 45 Ferretti L, Schmiegel B, Weinreich D, Yamauchi A, Kobayashi Y,  
46 Tajima F, Achaz G. 2016. Measuring epistasis in fitness land-  
47 scapes: The correlation of fitness effects of mutations. *Journal  
48 of Theoretical Biology*. 396:132–143.
- 49 Ferretti L, Weinreich D, Tajima F, Achaz G. 2018. Evolutionary  
50 constraints in fitness landscapes. *Heredity*. 121:466–481. Pub-  
51 lisher: Nature Publishing Group.
- 52 Flynn KM, Cooper TF, Moore FB, Cooper VS. 2013. The Environ-  
53 ment Affects Epistatic Interactions to Alter the Topology of an  
54 Empirical Fitness Landscape. *PLoS Genetics*. 9:1003426.
- 55 Fowler DM, Fields S. 2014. Deep mutational scanning: a new  
56 style of protein science. *Nature Methods*. 11:801–807. Pub-  
57 lisher: Nature Publishing Group.
- 58 Fragata I, Blanckaert A, Dias Louro MA, Liberles DA, Bank C.  
59 2019. Evolution in the light of fitness landscape theory. *Trends  
60 in Ecology and Evolution*. 34:69–82. Publisher: Elsevier Ltd.
- 61 Gao Y, Lin KT, Jiang T, Yang Y, Rahman M, Gong S, Bai J, Wang L,  
62 Sun J, Sheng L *et al.* 2022. Systematic characterization of short  
intronic splicing-regulatory elements in *SMN2* pre-mRNA. *63  
Nucleic Acids Research*. 50:731–749. 64
- Gardner J, Pleiss G, Weinberger KQ, Bindel D, Wilson AG. 2018. *65  
Gpytorch: Blackbox matrix-matrix gaussian process inference  
66 with gpu acceleration. Advances in neural information pro-  
67 cessing systems*. 31. 68
- Gavrilets S. 2004. *69  
Fitness Landscapes and the Origin of Species*.  
70 Princeton University Press.
- Haldane A, Flynn WF, He P, Levy RM. 2018. Coevolutionary  
71 Landscape of Kinase Family Proteins: Sequence Probabilities  
72 and Functional Motifs. *Biophysical Journal*. 114:21–31. 73
- Haldane A, Flynn WF, He P, Vijayan R, Levy RM. 2016. Structural  
74 propensities of kinase family proteins from a Potts model of  
75 residue co-variation. *Protein Science : A Publication of the  
76 Protein Society*. 25:1378–1384. 77
- Haseman JK, Elston RC. 1972. The investigation of linkage be-  
78 tween a quantitative trait and a marker locus. *Behavior Genet-  
79 ics*. 2:3–19. 80
- Herrera-Álvarez S, Patton JEJ, Thornton JW. 2025. Ancient biases  
81 in phenotype production drove the functional evolution of a  
82 protein family. 83
- Hwang S, Schmiegel B, Ferretti L, Krug J. 2018. Universality  
84 Classes of Interaction Structures for NK Fitness Landscapes.  
85 *Journal of Statistical Physics*. 172:226–278. arXiv: 1708.06556  
86 Publisher: Springer US. 87
- Ishigami Y, Wong MS, Martí-Gómez C, Ayaz A, Kooshkbaghi M,  
88 Hanson SM, McCandlish DM, Krainer AR, Kinney JB. 2024.  
89 Specificity, synergy, and mechanisms of splice-modifying  
90 drugs. *Nature Communications*. 15:1880. 91
- Jalal AS, Tran NT, Stevenson CE, Chan EW, Lo R, Tan X, Noy  
92 A, Lawson DM, Le TB. 2020. Diversification of DNA-Binding  
93 Specificity by Permissive and Specificity-Switching Mutations  
94 in the ParB/Noc Protein Family. *Cell Reports*. 32:107928. 95
- Johnson MS, Reddy G, Desai MM. 2023. Epistasis and evolution:  
96 recent advances and an outlook for prediction. *BMC Biology*.  
97 21:120. 98
- Johnston KE, Almhjell PJ, Watkins-Dulaney EJ, Liu G, Porter NJ,  
99 Yang J, Arnold FH. 2024. A combinatorially complete epistatic  
100 fitness landscape in an enzyme active site. *Proceedings of the  
101 National Academy of Sciences of the United States of America*.  
102 121:e2400439121. 103
- Kauffman S, Levin S. 1987. Towards a general theory of adaptive  
104 walks on rugged landscapes. *Journal of Theoretical Biology*.  
105 128:11–45. 106
- Kauffman SA, Weinberger ED. 1989. The NK model of rugged  
107 fitness landscapes and its application to maturation of the  
108 immune response. *Journal of Theoretical Biology*. 141:211–245. 109
- Khan AI, Dinh DM, Schneider D, Lenski RE, Cooper TF. 2011.  
110 Negative epistasis between beneficial mutations in an evolu-  
111 ting bacterial population. *Science*. 332:1193–1196. 112
- Kingman JFC. 1978. A simple model for the balance between  
113 selection and mutation. *Journal of Applied Probability*. 15:1–  
114 12. 115
- Kinney JB, McCandlish DM. 2019. Massively Parallel Assays and  
116 Quantitative Sequence–Function Relationships. *Annual Re-  
117 view of Genomics and Human Genetics*. 20:annurev–genom-  
118 083118–014845. 119
- Kinney JB, Murugan A, Callan CG, Cox EC. 2010. Using deep  
120 sequencing to characterize the biophysical mechanism of a  
121 transcriptional regulatory sequence. *Proceedings of the Na-  
122 tional Academy of Sciences of the United States of America*.  
123 107:9158–9163. 124

- 1 Komarova ES, Chervontseva ZS, Osterman IA, Evfratov SA,  
2 Rubtsova MP, Zatsepin TS, Semashko TA, Kostryukova ES,  
3 Bogdanov AA, Gelfand MS *et al.* 2020. Influence of the spacer  
4 region between the Shine–Dalgarno box and the start codon  
5 for fine-tuning of the translation efficiency in *Escherichia coli*.  
6 *Microbial Biotechnology*. 13:1254–1261.
- 7 Kondrashov AS, Sunyaev S, Kondrashov FA. 2002. Dobzhansky-  
8 Muller incompatibilities in protein evolution. *Proceedings of the*  
9 *National Academy of Sciences of the United States of*  
10 *America*. 99:14878–14883.
- 11 Kuo ST, Chang JK, Chang C, Shen WY, Hsu C, Lai SW, Chou  
12 HHD. 2025. Unraveling the start element and regulatory di-  
13 vergence of core promoters across the domain Bacteria.
- 14 Kuo ST, Jahn RL, Cheng YJ, Chen YL, Lee YJ, Hollfelder F, Wen  
15 JD, Chou HHD. 2020. Global fitness landscapes of the Shine-  
16 Dalgarno sequence. *Genome Research*. 30:711–723.
- 17 Kvittek DJ, Sherlock G. 2011. Reciprocal sign epistasis between  
18 frequently experimentally evolved adaptive mutations causes  
19 a rugged fitness landscape. *PLoS genetics*. 7:e1002056.
- 20 Liao SE, Sudarshan M, Regev O. 2023. Deciphering RNA splicing  
21 logic with interpretable machine learning. *Proceedings of the*  
22 *National Academy of Sciences*. 120:e2221165120.
- 23 Lite TLV, Grant RA, Nocedal I, Littlehale ML, Guo MS, Laub  
24 MT. 2020. Uncovering the basis of protein-protein interaction  
25 specificity with a combinatorially complete library. *eLife*. 9:1–  
26 57.
- 27 Manrubia S, Cuesta JA, Aguirre J, Ahnert SE, Altenberg L, Cano  
28 AV, Catalán P, Diaz-Uriarte R, Elena SF, García-Martín JA *et al.*  
29 2021. From genotypes to organisms: State-of-the-art and per-  
30 spectives of a cornerstone in evolutionary dynamics. *Physics of*  
31 *Life Reviews*. 38:55–106. arXiv: 2002.00363 Publisher: Else-  
32 vier B.V.
- 33 Marks DS, Hopf TA, Sander C. 2012. Protein structure prediction  
34 from sequence variation. *Nature Biotechnology*. 30:1072–1080.
- 35 Martí-Gómez C, McCandlish DM, Kinney JB. 2026a. Gauge-  
36 Fixer: overcoming parameter non-identifiability in models  
37 of sequence-function relationships. ISSN: 2692-8205 Pages:  
38 2025.12.08.693054 Section: New Results.
- 39 Martí-Gómez C, Zhou J, Chen WC, Stoltzfus A, Kinney JB, Mc-  
40 Candlish DM. 2026b. Inference and Visualization of Complex  
41 Genotype–Phenotype Maps. *Molecular Biology and Evolution*.  
42 43:msag023.
- 43 Matsui T, Mullis MN, Roy KR, Hale JJ, Schell R, Levy SF, Ehren-  
44 reich IM. 2022. The interplay of additivity, dominance, and  
45 epistasis on fitness in a diploid yeast cross. *Nature Commu-*  
46 *nications*. 13:1463.
- 47 McCandlish DM. 2011. Visualizing fitness landscapes. *Evolution*.  
48 65:1544–1558. Publisher: John Wiley & Sons, Ltd.
- 49 Moore JH, Williams SM. 2009. Epistasis and Its Implications for  
50 Personal Genetics. *The American Journal of Human Genetics*.  
51 85:309–320. Publisher: Elsevier.
- 52 Moulana A, Dupic T, Phillips AM, Chang J, Roffler AA, Gre-  
53 aney AJ, Starr TN, Bloom JD, Desai MM. 2023. The landscape  
54 of antibody binding affinity in SARS-CoV-2 Omicron BA.1  
55 evolution. *eLife*. 12:e83442.
- 56 Muzellec B, Teleńczuk M, Cabeli V, Andreux M. 2023. PyDESeq2:  
57 a python package for bulk RNA-seq differential expression  
58 analysis. *Bioinformatics*. 39:btad547.
- 59 Neidhart J, Szendro IG, Krug J. 2013. Exact results for amplitude  
60 spectra of fitness landscapes. *Journal of Theoretical Biology*.  
61 332:218–227. arXiv: 1301.1923 Publisher: Elsevier.
- 62 Neidhart J, Szendro IG, Krug J. 2014. Adaptation in Tunably  
Rugged Fitness Landscapes: The Rough Mount Fuji Model. *63*  
*Genetics*. 198:699–721. *64*
- Nguyen Ba AN, Lawrence KR, Rego-Costa A, Gopalakrishnan  
S, Temko D, Michor F, Desai MM. 2022. Barcoded bulk QTL  
mapping reveals highly polygenic and epistatic architecture  
of complex traits in yeast. *eLife*. 11:e73983. *65*  
*66*  
*67*  
*68*
- Noderer WL, Flockhart RJ, Bhaduri A, Diaz De Arce AJ, Zhang  
J, Khavari PA, Wang CL. 2014. Quantitative analysis of mam-  
malian translation initiation sites by facs-seq. *Molecular Sys-*  
*tems Biology*. 10:748. *69*  
*70*  
*71*  
*72*
- N’Guessan A, Tong WY, Heydari H, Ba ANN. 2025. Refining  
the resolution of the yeast genotype-phenotype map using  
single-cell RNA-sequencing. *eLife*. 13. Publisher: eLife Sci-  
ences Publications Limited. *73*  
*74*  
*75*  
*76*
- Ogbunugafor CB, Wylie CS, Diakite I, Weinreich DM, Hartl  
DL. 2016. Adaptive Landscape by Environment Interactions  
Dictate Evolutionary Dynamics in Models of Drug Resistance.  
*PLOS Computational Biology*. 12:e1004710. Publisher: Public  
Library of Science. *77*  
*78*  
*79*  
*80*
- O’Maille PE, Malone A, Dellas N, Andes Hess B, Smentek L,  
Sheehan I, Greenhagen BT, Chappell J, Manning G, Noel JP.  
2008. Quantitative exploration of the catalytic landscape sepa-  
rating divergent plant sesquiterpene synthases. *Nature Chem-*  
*ical Biology*. 4:617–623. *81*  
*82*  
*83*  
*84*  
*85*  
*86*
- Otwinowski J, McCandlish DM, Plotkin JB. 2018. Inferring  
the shape of global epistasis. *Proceedings of the National*  
*Academy of Sciences*. 115:E7550–E7558. *87*  
*88*  
*89*
- Papkou A, Garcia-Pastor L, Escudero JA, Wagner A. 2023.  
A rugged yet easily navigable fitness landscape. *Science*.  
382:eadh3860. Publisher: American Association for the Ad-  
vancement of Science. *90*  
*91*  
*92*  
*93*
- Petti S, Martí-Gómez C, Kinney JB, Zhou J, McCandlish DM.  
2025. On learning functions over biological sequence space:  
relating gaussian process priors, regularization, and gauge  
fixing. arXiv preprint arXiv:2504.19034. . *94*  
*95*  
*96*  
*97*
- Poelwijk FJ, Socolich M, Ranganathan R. 2019. Learning the pat-  
tern of epistasis linking genotype and phenotype in a protein.  
*Nature Communications*. 10:1–11. Publisher: Springer US. *98*  
*99*  
*100*
- Pokusaeva VO, Usmanova DR, Putintseva EV, Espinar L, Sark-  
isyan KS, Mishin AS, Bogatyreva NS, Ivankov DN, Akopyan  
AV, Arvakumov SY *et al.* 2019. An experimental assay of the in-  
teractions of amino acids from orthologous sequences shaping  
a complex fitness landscape. *PLoS Genet.*. 15:e1008079. *101*  
*102*  
*103*  
*104*  
*105*
- Rabani M, Pieper L, Chew GL, Schier AF. 2017. A Massively  
Parallel Reporter Assay of 3 UTR Sequences Identifies In  
Vivo Rules for mRNA Degradation. *Molecular Cell*. 68:1083–  
1094.e5. Publisher: Elsevier. *106*  
*107*  
*108*  
*109*
- Rasmussen CE, Williams CKI. 2008. *Gaussian processes for machine*  
*learning*. Adaptive computation and machine learning. MIT  
Press. Cambridge, Mass.. third edition. *110*  
*111*  
*112*
- Reddy G, Desai MM. 2021. Global epistasis emerges from a  
generic model of a complex trait. *eLife*. 10:1–36. *113*  
*114*
- Rojas Echenique JI, Kryazhinskiy S, Nguyen Ba AN, Desai MM.  
2019. Modular epistasis and the compensatory evolution of  
gene deletion mutants. *PLoS genetics*. 15:e1007958. *115*  
*116*  
*117*
- Rosenberg AB, Patwardhan RP, Shendure J, Seelig G. 2015. Learn-  
ing the sequence determinants of alternative splicing from  
millions of random sequences. *Cell*. 163:698–711. *118*  
*119*  
*120*
- Rotrattanadumrong R, Yokobayashi Y. 2022. Experimental ex-  
ploration of a ribozyme neutral network using evolutionary  
algorithm and deep learning. *Nature Communications*. pp.  
1–14. Publisher: Springer US ISBN: 4146702232. *121*  
*122*  
*123*  
*124*

- 1 Rubin AF, Gelman H, Lucas N, Bajjalieh SM, Papenfuss AT,  
2 Speed TP, Fowler DM. 2017. A statistical framework for ana-  
3 lyzing deep mutational scanning data. *Genome Biology*. 18:1–  
4 15. Publisher: Genome Biology.
- 5 Sackton TB, Hartl DL. 2016. Genotypic Context and Epistasis in  
6 Individuals and Populations. *Cell*. 166:279–287.
- 7 Schmiedel JM, Lehner B. 2019. Determining protein structures  
8 using deep mutagenesis. *Nature Genetics*. 51:1177–1186.
- 9 Schmiegel B, Krug J. 2014. Evolutionary Accessibility of Modular  
10 Fitness Landscapes. *Journal of Statistical Physics*. 154:334–  
11 355.
- 12 Shine J, Dalgarno L. 1975. Determinant of cistron specificity in  
13 bacterial ribosomes. *Nature*. 254:34–38.
- 14 Somermeyer LG, Fleiss A, Mishin AS, Bozhanova NG, Igolkina  
15 AA, Meiler J, Pujol MEA, Putintseva EV, Sarkisyan KS, Kon-  
16 drashov FA. 2022. Heterogeneity of the *gfp* fitness landscape  
17 and data-driven protein design. *Elife*. 11:e75842.
- 18 Soo VW, Swadling JB, Faure AJ, Warnecke T. 2021. Fitness  
19 landscape of a dynamic RNA structure. *PLoS Genetics*.  
20 17:e1009353. ISBN: 1111111111.
- 21 Soyk S, Benoit M, Lippman ZB. 2020. New Horizons for Dis-  
22 secting Epistasis in Crop Quantitative Trait Variation. *Annual*  
23 *Review of Genetics*. 54:287–307. Publisher: Annual Reviews.
- 24 Stadler PF. 1996. Landscapes and their correlation functions.  
25 *Journal of Mathematical Chemistry*. 20:1–45.
- 26 Stadler PF. 2002. Fitness landscapes, In: Lässig M, Valleriani  
27 A, editors, *Biological Evolution and Statistical Physics*, Springer.  
28 Berlin, Heidelberg. pp. 183–204.
- 29 Stadler PF, Happel R. 1999. Random field models for fitness  
30 landscapes. *Journal of Mathematical Biology*. 38:435–478.
- 31 Starr TN, Picton LK, Thornton JW. 2017. Alternative evolution-  
32 ary histories in the sequence space of an ancient protein. *Nature*.  
33 549:409–413.
- 34 Sundar V, Tu B, Guan L, Esvelt K. 2024. A NEW ULTRA-HIGH-  
35 THROUGHPUT ASSAY FOR MEASURING PROTEIN FIT-  
36 NESS. *Proceedings of the Generalization and Epistemic Mea-*  
37 *sures (GEM) Workshop at the International Conference on*  
38 *Learning Representations (ICLR)*. .
- 39 Szendro IG, Schenk MF, Franke J, Krug J, de Visser JAGM. 2013.  
40 Quantitative analyses of empirical fitness landscapes. *Journal*  
41 *of Statistical Mechanics: Theory and Experiment*. 2013:P01005.  
42 Publisher: IOP Publishing and SISSA.
- 43 Vaishnav ED, de Boer CG, Molinet J, Yassour M, Fan L, Adiconis  
44 X, Thompson DA, Levin JZ, Cubillos FA, Regev A. 2022. The  
45 evolution, evolvability and engineering of gene regulatory  
46 DNA. *Nature*. Publisher: Springer US.
- 47 Wang K, Pleiss G, Gardner J, Tyree S, Weinberger KQ, Wilson  
48 AG. 2019. Exact Gaussian Processes on a Million Data Points.  
49 *NeurIPS*. .
- 50 Wang T, Zhao D, Tian S. 2015. An overview of kernel alignment  
51 and its applications. *Artificial Intelligence Review*. 43:179–192.
- 52 Weinberger E. 1990. Correlated and uncorrelated fitness land-  
53 scapes and how to tell the difference. *Biological Cybernetics*.  
54 63:325–336.
- 55 Weinreich DM, Delaney NF, DePristo MA, Hartl DL. 2006. Dar-  
56 winian evolution can follow only very few mutational paths  
57 to fitter proteins. *Science*. 312:111–114.
- 58 Weinreich DM, Lan Y, Jaffe J, Heckendorn RB. 2018. The Infl-  
59 uence of Higher-Order Epistasis on Biological Fitness Land-  
60 scape Topography. *Journal of Statistical Physics*. 172:208–225.  
61 Publisher: Springer US.
- 62 Weinreich DM, Watson RA, Chao L. 2005. Perspective: sign  
epistasis and genetic constraint on evolutionary trajectories. 63  
*Evolution*. 59:1165–1174. 64
- Westmann CA, Goldbach L, Wagner A. 2024a. Entangled adap- 65  
tive landscapes facilitate the evolution of gene regulation by 66  
exaptation. Pages: 2024.11.10.620926 Section: New Results. 67
- Westmann CA, Goldbach L, Wagner A. 2024b. The highly rugged 68  
yet navigable regulatory landscape of the bacterial transcrip- 69  
tion factor TetR. *Nature Communications*. 15:10745. Publisher: 70  
Nature Publishing Group. 71
- Wong MS, Kinney JB, Krainer AR. 2018. Quantitative Activity 72  
Profile and Context Dependence of All Human 5' Splice Sites. 73  
*Molecular Cell*. 71:1012–1026.e3. Publisher: Elsevier Inc. 74
- Wright S. 1932. The roles of mutation, inbreeding, crossbreeding 75  
and selection in evolution. *Proceedings of the Sixth Interna-*  
76 *tional Congress of Genetics*. pp. 356–366. 77
- Wu NC, Dai L, Olson CA, Lloyd-Smith JO, Sun R. 2016. Adap- 78  
tation in protein fitness landscapes is facilitated by indirect 79  
paths. *eLife*. 5:1–21. 80
- Zarin T, Lehner B. 2024. A complete map of specificity encoding 81  
for a partially fuzzy protein interaction. 82
- Zebell SG, Martí-Gómez C, Fitzgerald B, Cunha CP, Lach M,  
83 Seman BM, Hendelman A, Sretenovic S, Qi Y, Bartlett M *et al.*  
84 2025. Cryptic variation fuels plant phenotypic change through  
85 hierarchical epistasis. *Nature*. pp. 1–9. Publisher: Nature Pub-  
86 lishing Group. 87
- Zhou J, Martí-Gómez C, Petti S, McCandlish DM. 2025. Learning 88  
sequence-function relationships with scalable, interpretable 89  
Gaussian processes. *eLife*. 14. Publisher: eLife Sciences Publi-  
90 cations Limited. 91
- Zhou J, McCandlish DM. 2020. Minimum epistasis interpolation 92  
for sequence-function relationships. *Nature Communications*.  
93 11. 94
- Zhou J, Wong MS, Chen Wc, Krainer AR, Justin B, Mccandlish  
95 DM. 2022. Higher-order epistasis and phenotypic prediction.  
96 *Proc. Natl. Acad. Sci. USA*. 119. 97

## 1 Supplementary Information

### 2 Epistasis on fitness landscapes

3 Let  $f$  be an  $\alpha^\ell$ -dimensional vector encoding the fitness associated with each genotype in the space of possible haploid sequences with  
4  $\alpha$  alleles and  $\ell$  sites  $S = \{1, 2, \dots, \ell\}$ . In this section, we review two common ways to quantify epistasis in a fitness landscape: one  
5 based on local epistatic coefficients (Zhou and McCandlish 2020; Chen et al. 2021) and a second one based on the variance explained by  
6 genetic interactions of different orders (Stadler and Happel 1999; Zhou et al. 2022) or across subsets of sites (Martí-Gómez et al. 2026b).

7 **Variance components** Any fitness landscape  $f$  can be decomposed into orthogonal components  $f_k$  corresponding to genetic interactions  
8 of order  $k$

$$f = \sum_{k=0}^{\ell} f_k. \quad (S1)$$

9 These components can be obtained by projecting the  $f$  into the  $k$ th-subspace using the orthogonal projection matrix  $P_k$  given by the  
10 Krawtchouk polynomials (Stadler and Happel 1999; Zhou et al. 2022):

$$P_k(x, x') = \alpha^{-\ell} \sum_{q=0}^k (-1)^q (\alpha - 1)^{k-q} \binom{d(x, x')}{q} \binom{\ell - d(x, x')}{k - q}. \quad (S2)$$

11 Moreover, each  $k$ -th order subspace can be further decomposed into orthogonal components corresponding to the contributions of  
12 interactions between specific subsets of sites  $U \subseteq S$  of size  $k$ :

$$f = \sum_{U \in \mathcal{P}(S)} f_U, \quad (S3)$$

13 where  $\mathcal{P}(S)$  is the set of all subsets of  $S$  (i.e. the power set of  $S$ ) and  $f_U = P_U f$ , where  $P_U$  is an orthogonal projection matrix onto the  
14 corresponding subspace. Each such  $P_U$  can be expressed as a Kronecker product of a series of site-specific projection matrices onto the  
15 constant subspace, defined by  $P_{\text{con}} = \frac{1}{\alpha} 11^T$ , and the orthogonal or additive subspace, given by  $P_{\text{add}} = I - P_{\text{con}}$ . In particular

$$P_U = \bigotimes_{i=1}^{\ell} P_{i \in U}, \quad (S4)$$

16 where  $P_{i \in U} = P_{\text{add}}$  when site  $i \in U$  and  $P_{i \in U} = P_{\text{con}}$  otherwise (Martí-Gómez et al. 2026b):

$$P_U(x, x') = \alpha^{-\ell} \prod_{i \in U: x_i = x'_i} (\alpha - 1) \prod_{i \in U: x_i \neq x'_i} (-1). \quad (S5)$$

17 Because  $U$ -components are orthogonal to each other, the total variance in a fitness landscape can be expressed as a sum of vari-  
18 ances explained by genetic interactions of order  $k$  (i.e. summing over all  $U$  with  $|U| = k$ ) or by interactions between the sites  $U$ ,  
19 respectively (Martí-Gómez et al. 2026b):

$$\text{Var}[f] = \sum_{k=1}^{\ell} \text{Var}[f_k] = \sum_{U \in \mathcal{P}(S)} \text{Var}[f_U]. \quad (S6)$$

20 **Local epistatic coefficients** The classical way to quantify epistasis is through the definition of the epistatic coefficient  $\epsilon$  for a pair of  
21 mutations at different loci  $A \rightarrow a$  and  $B \rightarrow b$ , which quantifies the difference in the effect of mutation  $A \rightarrow a$  in the presence of allele  $B$   
22 compared to that in the presence of allele  $b$  at the other locus:

$$\epsilon = (f_{AB} - f_{aB}) - (f_{Ab} - f_{ab}). \quad (S7)$$

23 This epistatic coefficient is defined locally, as it depends only on four specific genotypes sharing the same genetic background. Still, we  
24 can compute the average magnitude of local epistatic interactions  $\bar{\epsilon}^2$  in a given fitness landscape by averaging the its squared values  
25 across every possible pair of mutations at every possible genetic background:

$$\bar{\epsilon}^2 = \frac{1}{s} f^T \Delta^{(2)} f, \quad (S8)$$

26 where  $s = \binom{\ell}{2} \binom{\alpha}{2} \alpha^{\ell-2}$  is the number of epistatic coefficients and  $\Delta^{(2)}$  is a positive semi-definite sparse matrix (Zhou and McCandlish  
27 2020; Chen et al. 2021). These results can be generalized to quantify local epistatic coefficients of any order  $P$  using

$$\Delta^{(P)}(x, x') = \begin{cases} 0 & \text{if } d(x, x') > P \\ (-1)^{d(x, x')} (\alpha - 1)^{P - d(x, x')} \binom{\ell - P}{P - d(x, x')} & \text{if } d(x, x') \leq P, \end{cases} \quad (S9)$$

28 where  $d(x, x')$  represents the Hamming distance i.e. number of single point mutations, separating sequences  $x$  and  $x'$ , such that the  
29 sum of the squared  $P$ -th order local epistatic coefficients is given by  $f^T \Delta^{(2)} f$  and the mean squared  $P$ -th order local epistatic coefficient

is obtained by dividing by the number of such coefficients  $\binom{\ell}{P} \binom{\alpha}{2}^P \alpha^{\ell-P}$  (Chen *et al.* 2021). Interestingly,  $\Delta^{(P)}$  can also be expressed as a weighted sum of the projection operators into the  $k$ -th order subspaces given by  $P_k$ :

$$\Delta^{(P)} = \sum_{k=0}^{\ell} \lambda_k P_k = \sum_{k=P}^{\ell} \lambda_k P_k, \quad (\text{S10})$$

where  $\lambda_k = \alpha^P \binom{k}{P}$  (Chen *et al.* 2021) correspond to the eigenvalues of the  $\Delta^{(P)}$  operator. Note that the second equality arises because  $\binom{k}{P} = 0$  for  $k < P$ , corresponding to the fact that all local  $P$ -th order epistatic coefficients are zero for any function  $f$  with maximal order less than  $P$  (so that all such functions are in the null space of  $\Delta^{(P)}$ ) and that moreover alteration of any component of order less than  $P$  for an arbitrary  $f$  leaves all  $P$ -th order local epistatic coefficients unchanged.

### Epistatic coefficients among subsets of sites

In this section, we describe how the  $\Delta^{(P)}$  operator that extracts the sum of squared epistatic interactions in a fitness landscape can be decomposed into the sum of simpler operators that we call  $\Delta^{(U)}$ , corresponding to local epistatic interactions only among a subset of sites  $U \subseteq S$ . These results generalize the main text results focused on local epistatic interactions between pairs of sites  $i$  and  $j$  i.e.  $U = \{i, j\}$ .

Let  $E_U$  be a  $s_U \times \alpha^{\ell}$  matrix such that the entries in the  $E_U f$  vector encode all the local  $|U|$ -th order epistatic coefficient for a subset of positions  $U$ .  $s_U$  is the number of different epistatic coefficients corresponding to local  $|U|$ -th order mutant epistatic coefficients within sites  $U$ , and is given by the product of the number of genetic backgrounds at sites not in  $U$   $\alpha^{\ell-|U|}$  and the number of possible combinations of mutations across the sites  $U$   $\binom{\alpha}{2}^{|U|}$ :

$$s_U = \alpha^{\ell-|U|} \binom{\alpha}{2}^{|U|} = \alpha^{\ell-|U|} \left( \frac{\alpha(\alpha-1)}{2} \right)^{|U|} = \frac{\alpha^{\ell}}{2^{|U|}} (\alpha-1)^{|U|}. \quad (\text{S11})$$

Thus,

$$\overline{e_U^2} = \frac{1}{s_U} (E_U f)^T E_U f = \frac{1}{s_U} f^T \Delta^{(U)} f, \quad (\text{S12})$$

where the entries of  $\Delta^{(U)}$  for a pair of sequences  $x$  and  $x'$  can be obtained by summing over all possible local  $|U|$ -mutant epistatic coefficients between sites in  $U$

$$\Delta^{(U)}(x, x') = \sum_{m=1}^{s_U} E_U(m, x) E_U(m, x'). \quad (\text{S13})$$

$E_U(m, x') = 0$  if sequence  $x'$  is not involved in epistatic coefficient  $m$ , and takes values  $-1$  or  $1$  otherwise. Thus, we only need to sum over local epistatic coefficients involving both  $x$  and  $x'$ . If  $x_i \neq x'_i$  for any position  $i \notin U$ , then sequences  $x$  and  $x'$  cannot be involved in any epistatic coefficient and thus  $\Delta^{(U)}(x, x') = 0$ . Otherwise,  $E_U(m, x) E_U(m, x') = (-1)^{d(x, x')}$  depending on the Hamming distance between them  $d(x, x')$ . Moreover, the number of local  $|U|$ -mutant epistatic coefficients that involve both  $x$  and  $x'$  is given by  $(\alpha-1)^{|U|-d(x, x')}$ . Thus

$$\Delta^{(U)}(x, x') = \begin{cases} (-1)^{d(x, x')} (\alpha-1)^{|U|-d(x, x')} & \text{if } x_i = x'_i \forall i \notin U \\ 0 & \text{otherwise.} \end{cases} \quad (\text{S14})$$

We can verify that we can recover the  $\Delta^{(P)}$  operator by summing  $\Delta^{(U)}$  over all possible subsets of sites  $U$  of size  $P$ . If  $d(x, x') > P$ , there is no single set of sites  $U$  for which the context at sequences  $x$  and  $x'$  can match and thus  $\sum_{U:|U|=P} \Delta^{(U)}(x, x') = 0$ . For  $d(x, x') \leq P$ , the entries of  $\Delta^{(U)}(x, x')$  are the same but they are summed over multiple  $U$ . As we only sum over  $U$ 's such that  $x$  and  $x'$  share the context, the number of times we are summing them corresponds to the number of ways we can choose  $P - d(x, x')$  sites that match out of the  $\ell - P$  sites in the shared context between them. Thus,

$$\sum_{U:|U|=P} \Delta^{(U)}(x, x') = \begin{cases} (-1)^{d(x, x')} (\alpha-1)^{P-d(x, x')} \binom{\ell-P}{P-d(x, x')} & \text{if } d(x, x') \leq P \\ 0 & \text{if } d(x, x') > P, \end{cases} \quad (\text{S15})$$

which exactly matches the  $\Delta^{(P)}$  operator (Zhou and McCandlish 2020; Chen *et al.* 2021).

### Relationship between local epistatic coefficients and variance components

In this section, we describe some properties of the new  $\Delta^{(U)}$  operator and the relationships with the projection operators into the subspaces corresponding to interactions of different orders and subsets of sites  $U$ .

One useful property of the  $\Delta^{(U)}$  operator is that it can be expressed as a Kronecker product of site-specific matrices

$$\Delta^{(U)} = \bigotimes_{i=1}^{\ell} \Delta_{i \in U} = \bigotimes_{i=1}^{\ell} \begin{cases} \alpha P_{\text{add}} & \text{if } i \in U \\ I & \text{if } i \notin U, \end{cases} \quad (\text{S16})$$

1 with entry-wise formula given by:

$$\Delta^{(U)}(x, x') = \prod_{\substack{i \in U \\ x_i = x'_i}} (\alpha - 1) \prod_{\substack{i \in U \\ x_i \neq x'_i}} (-1) \prod_{\substack{i \notin U \\ x_i = x'_i}} (1) \prod_{\substack{i \notin U \\ x_i \neq x'_i}} 0. \quad (\text{S17})$$

2 Thus, we can use the mixed-product property to show that the columns of  $P_{U'}$  are eigenvectors of  $\Delta^{(U)}$  with eigenvalue  $\alpha^{|U|}$  if  $U \subseteq U'$   
 3 or are in the null space of  $\Delta^{(U)}$  whenever  $U \not\subseteq U'$ .

$$\Delta^{(U)} P_{U'} = \bigotimes_{i=1}^{\ell} \Delta_{i \in U} \bigotimes_{i=1}^{\ell} P_{i \in U'} = \bigotimes_{i=1}^{\ell} \Delta_{i \in U} P_{i \in U'}, \quad (\text{S18})$$

4 where

$$\Delta_{i \in U} P_{i \in U'} = \begin{cases} P_{\text{con}} & \text{if } i \notin U \wedge i \notin U' \\ P_{\text{add}} & \text{if } i \notin U \wedge i \in U' \\ 0 & \text{if } i \in U \wedge i \notin U' \\ \alpha P_{\text{add}} & \text{if } i \in U \wedge i \in U'. \end{cases} \quad (\text{S19})$$

5 Therefore

$$\Delta^{(U)} P_{U'} = \begin{cases} \alpha^{|U|} P_{U'} & \text{if } U \subseteq U' \\ 0 & \text{if } U \not\subseteq U'. \end{cases} \quad (\text{S20})$$

6 Because the projection operators  $P_{U'}$  are orthogonal to each other, we can express  $\Delta^{(U)}$  as a linear combination of  $P_{U'}$  weighted by  
 7 their eigenvalues (either 0 or  $\alpha^{|U|}$ ):

$$\Delta^{(U)} = \alpha^{|U|} \sum_{U': U \subseteq U'} P_{U'}, \quad (\text{S21})$$

8 so that, using the fact that sums of projection matrices into orthogonal subspaces are themselves projection matrices, we see that  $\Delta^{(U)}$   
 9 is itself just a  $|U|$ -dependent constant times a specific projection matrix. Moreover, this relation between  $\Delta^{(U)}$  and the  $P_{U'}$  with  $U \subseteq U'$   
 10 can be used to derive the relationship between the  $\Delta^{(P)}$  operator and the projection operator into the  $k$ th order subspace  $P_k$ :

$$\Delta^{(P)} = \sum_{U: |U|=P} \Delta^{(U)} = \sum_{U: |U|=P} \alpha^{|U|} \sum_{U': U \subseteq U'} P_{U'} = \alpha^P \sum_{U: |U|=P} \sum_{U': U \subseteq U'} P_{U'}. \quad (\text{S22})$$

11 It is now easy to see that the number of times we are summing each  $P_{U'}$  depends on the size of  $U'$ . In particular, it corresponds to the  
 12 number of  $U: |U|=P$  subsets in  $U'$ , which can be calculated as the number of ways of choosing  $|U|=P$  sites out of  $|U'|$ . Then, we  
 13 can use the fact that  $P_k = \sum_{U: |U|=k} P_U$  (Martí-Gómez *et al.* 2026b) to recover the eigendecomposition of the  $\Delta^{(P)}$  operator (Chen *et al.*  
 14 2021):

$$\Delta^{(P)} = \alpha^P \sum_{U': |U'| \geq P} \binom{|U'|}{P} P_{U'} = \sum_{k \geq P} \alpha^P \binom{k}{P} \sum_{U': |U'|=k} P_{U'} = \sum_{k \geq P} \alpha^P \binom{k}{P} P_k. \quad (\text{S23})$$

#### 15 Relationship between local $|U|$ -way epistatic coefficients and the epistatic variance explained by the subset $U$

16 A useful low dimensional summary statistic to describe the patterns of genetic interactions across subsets of sites in a fitness landscape  
 17 is the epistatic variance explained by a set of sites  $U$  (Crawford *et al.* 2017; Reddy and Desai 2021; Martí-Gómez *et al.* 2026b). We define  
 18 this variance  $\text{Var}^{(U)}[f]$  to include not only the variance explained by  $|U|$ -way interactions between mutations at the sites  $U$ , but also  
 19 all other interactions of order higher than  $|U|$  where all sites in  $U$  are involved, so that  $\text{Var}^{(U)}[f]$  provides an overall measure of the  
 20 amount of  $|U|$ -way and higher epistasis that the subset of sites  $U$  is involved in. This quantity can be computed by projecting the  
 21 fitness landscape  $f$  onto the  $2^\ell$  subspaces defined by every possible  $U'$  given by  $f_{U'} = P_{U'} f$ , summing over the  $f_{U'}$ 's of interest and  
 22 computing the inner product.

$$\text{Var}^{(U)}[f] = \left( \sum_{U': U \subseteq U'} f_{U'} \right)^T \left( \sum_{U': U \subseteq U'} f_{U'} \right) = \left( \sum_{U': U \subseteq U'} P_{U'} f \right)^T \left( \sum_{U': U \subseteq U'} P_{U'} f \right) = f^T \left( \sum_{U': U \subseteq U'} P_{U'} \right) f, \quad (\text{S24})$$

23 where

$$\sum_{U': U \subseteq U'} P_{U'} = \sum_{U': U \subseteq U'} \bigotimes_{i=1}^{\ell} P_{i \in U'}. \quad (\text{S25})$$

24 As the Kronecker product is commutative up to row and column permutations given by  $Q_1$  and  $Q_2$ , we can write the factors in order  
 25 depending on whether the sites are in  $U$  or not:

$$\begin{aligned} \sum_{U': U \subseteq U'} P_{U'} &= Q_1 \left( \sum_{U': U \subseteq U'} \bigotimes_{i \in U} P_{i \in U'} \bigotimes_{i \notin U} P_{i \in U'} \right) Q_2 = Q_1 \left( \bigotimes_{i \in U} P_{\text{add}} \sum_{U': U \subseteq U'} \bigotimes_{i \notin U} P_{i \in U'} \right) Q_2 \\ &= Q_1 \left( \bigotimes_{i \in U} P_{\text{add}} \bigotimes_{i \notin U} I \right) Q_2 = \bigotimes_{i=1}^{\ell} \begin{cases} P_{\text{add}} & \text{if } i \in U \\ I & \text{if } i \notin U \end{cases} = \alpha^{-|U|} \Delta^{(U)}. \end{aligned} \quad (\text{S26})$$

Thus, this shows that the sum of squared epistatic coefficients among a subset of sites  $U$  given by  $f^T \Delta^{(U)} f$  is equal to the epistatic variance explained by those sites  $\text{Var}^{(U)}[f]$  multiplied by a factor of  $\alpha^{|U|}$ . Noting that  $f^T \Delta^{(U)} f = \overline{\epsilon_U^2} s_U$  and simplifying  $s_U / \alpha^{|U|}$  yields Equation 2 in the main text. Similarly, the sum of squared local  $P$ -epistatic coefficients across the complete fitness landscape given by  $f^T \Delta^{(P)} f$  can also be interpreted as the sum the epistatic variances of all possible combinations of  $P$  sites multiplied by a factor of  $\alpha^P$ .

### Relationship between mean and variance of local $|U|$ -way epistatic coefficients and the variance components

In the previous sections, we have shown the relationships between the average squared epistatic coefficients between mutations at sites  $U$  and the total variance explained by interactions of order  $|U|$  or larger involving all sites in  $U$ . In this section, we describe how these quantities relate to the mean and variance of the distribution of  $|U|$ -way epistatic coefficients for a set of mutations within sites  $U$ . To do so, we first decompose the epistatic coefficients between mutations at sites  $U$  into the epistatic coefficients for each possible combination of mutations at sites  $U$ . Using this more granular statistic, we derive the mean of the epistatic coefficients between a given set of mutations, and use its square, together with the average squared epistatic coefficient to compute the variance. Finally, we take the average of these variances over all possible combinations of mutations at sites  $U$  and derive its relationship to the variance components defined over the subsets of sites  $U$ .

Let  $C = \{\{c_{11}, c_{12}\}, \dots, \{c_{k1}, c_{k2}\}\}$  be the set of  $k$  pairs of characters at a set of  $k$  positions  $U$  and  $\epsilon_{UC}$  be the  $\alpha^{\ell-k}$ -dimensional vector of  $k$ -way epistatic coefficients between the mutations specified by  $C$  across every possible genetic background. Let  $E_{UC}$  be an  $\alpha^{\ell-k} \times \alpha^\ell$  such that  $\epsilon_{UC} = E_{UC} f$  and note that this matrix that can be expressed as a Kronecker product of site specific factors:

$$E_{UC} = \bigotimes_i^\ell \begin{cases} E_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ I & \text{if } i \notin U, \end{cases} \quad (\text{S27})$$

where  $E_i^{c_1, c_2}$  is a row vector with entries given by

$$E_i^{c_1, c_2}(x) = \begin{cases} 1 & \text{if } x_i = c_1 \\ -1 & \text{if } x_i = c_2 \\ 0 & \text{otherwise.} \end{cases} \quad (\text{S28})$$

Using these expressions, we can compute the average squared epistatic coefficients restricted to the sets of mutations defined by  $C$  as

$$\overline{\epsilon_{UC}^2} = \frac{1}{\alpha^{\ell-|U|}} (E_{UC} f)^T (E_{UC} f) = \frac{1}{\alpha^{\ell-|U|}} f^T E_{UC}^T E_{UC} f = \frac{1}{\alpha^{\ell-|U|}} f^T \Delta^{(UC)} f, \quad (\text{S29})$$

where  $\Delta^{(UC)}$  can be easily computed by using the Kronecker factorization of  $E_{UC}$

$$\Delta^{(UC)} = \bigotimes_i^\ell \begin{cases} \Delta_i^{c_1, c_2} & \text{if } i \in U \\ I & \text{if } i \notin U. \end{cases} \quad (\text{S30})$$

The factor  $\Delta_i^{c_1, c_2} = (E_i^{c_1, c_2})^T E_i^{c_1, c_2}$  takes the form

$$\Delta_i^{c_1, c_2}(x, x') = \begin{cases} 1 & \text{if } (x_i = c_1 \wedge x'_i = c_1) \vee (x_i = c_2 \wedge x'_i = c_2) \\ -1 & \text{if } (x_i = c_1 \wedge x'_i = c_2) \vee (x_i = c_2 \wedge x'_i = c_1) \\ 0 & \text{otherwise.} \end{cases} \quad (\text{S31})$$

It is easy to see that  $\sum_{c_1, c_2} \Delta_i^{c_1, c_2} = \alpha P_{\text{add}}$  and that  $\Delta^{(U)}$  can be expressed as a sum of  $\Delta^{(UC)}$  for all possible combinations of pairs of alleles  $C$  at sites  $U$ .

$$\sum_C \Delta^{(UC)} = \sum_C \bigotimes_i^\ell \begin{cases} \Delta_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ I & \text{if } i \notin U \end{cases} = \bigotimes_i^\ell \begin{cases} \sum_C \Delta_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ I & \text{if } i \notin U \end{cases} = \Delta^{(U)}. \quad (\text{S32})$$

Next, we use the  $E_{UC}$  to compute the average epistatic coefficients  $\overline{\epsilon_{UC}}$  between the set of mutations defined by  $C$  at sites  $U$  across all possible genetic backgrounds:

$$\overline{\epsilon_{UC}} = \frac{1}{\alpha^{\ell-|U|}} 1^T E_{UC} f. \quad (\text{S33})$$

This quantity corresponds to the previously proposed background-averaged epistatic coefficient for parametrizing and describing sequence-function relationships (Faure et al. 2024a; Petti et al. 2025). As the sign of  $\overline{\epsilon_{UC}}$  depends on the ordering of alleles  $c_{i1}$  and  $c_{i2}$  for each position  $i$ , we can use its squared value  $\overline{\epsilon_{UC}^2}$  to describe the magnitude of the average epistatic coefficients independently of the choice of reference allele at each position:

$$\overline{\epsilon_{UC}^2} = \left( \frac{1}{\alpha^{\ell-|U|}} 1^T E_{UC} f \right)^2 = \frac{1}{\alpha^{2\ell-2|U|}} f^T E_{UC}^T 11^T E_{UC} f. \quad (\text{S34})$$

- 1 Let  $A_{UC} = \frac{1}{\alpha^{\ell-|U|}} E_{UC}^T 11^T E_{UC}$ , so that  $\overline{\epsilon_{UC}^2}$  can be expressed as a function of a quadratic form with this matrix given by  $\frac{1}{\alpha^{\ell-|U|}} f^T A_{UC} f$ .  
 2 To derive a simple expression for  $A_{UC}$  and understand its relationship with other objects, we use the fact that all of the matrices  $E_{UC}$   
 3 and  $11^T$  can be expressed as Kronecker products of site specific factors:

$$A_{UC} = \frac{1}{\alpha^{\ell-|U|}} E_{UC}^T 11^T E_{UC} = \frac{1}{\alpha^{\ell-|U|}} \bigotimes_i \begin{cases} (E_i^{c_{i1}, c_{i2}})^T E_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ 11^T & \text{if } i \notin U \end{cases} = \bigotimes_i \begin{cases} \Delta_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ \frac{1}{\alpha} 11^T & \text{if } i \notin U. \end{cases} \quad (\text{S35})$$

- 4 It is easy to see that  $A_{UC}$  is in the subspace defined by genetic interactions of order  $|U|$  between sites  $U$  since  $P_U A_{UC} = A_{UC}$ . In fact,  
 5 if we sum over all possible combinations of mutations  $C$ , we obtain a matrix that is proportional to the projection operator into the  
 6  $U$ -subspace  $P_U$ :

$$\sum_C A_{UC} = \sum_C \bigotimes_i \begin{cases} \Delta_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ \frac{1}{\alpha} 11^T & \text{if } i \notin U. \end{cases} = \bigotimes_i \begin{cases} \sum_{c_{i1}, c_{i2}} \Delta_i^{c_{i1}, c_{i2}} & \text{if } i \in U \\ \frac{1}{\alpha} 11^T & \text{if } i \notin U. \end{cases} = \bigotimes_i \begin{cases} \alpha P_{\text{add}} & \text{if } i \in U \\ P_{\text{con}} & \text{if } i \notin U. \end{cases} = \alpha^{|U|} P_U. \quad (\text{S36})$$

- 7 Using this equivalence, we can show that the average squared mean epistatic coefficient for mutations at sites  $U$  is proportional to the  
 8 variance explained by genetic interactions of order  $|U|$  between those sites:

$$\begin{aligned} \frac{1}{\binom{\alpha}{2}^{|U|}} \sum_C (\overline{\epsilon_{UC}})^2 &= \frac{1}{\binom{\alpha}{2}^{|U|}} \sum_C \frac{1}{\alpha^{\ell-|U|}} f^T A_{UC} f = \frac{1}{\binom{\alpha}{2}^{|U|} \alpha^{\ell-|U|}} f^T \left( \sum_C A_{UC} \right) f = \frac{1}{\binom{\alpha}{2}^{|U|} \alpha^{\ell-|U|}} f^T \left( \sum_C A_{UC} \right) f \\ &= \frac{1}{\binom{\alpha}{2}^{|U|} \alpha^{\ell-|U|}} \alpha^{|U|} f^T P_U f = \frac{2^{|U|}}{\alpha^\ell} \left( \frac{\alpha}{\alpha-1} \right)^{|U|} \text{Var}[f_U] \end{aligned} \quad (\text{S37})$$

- 9 Finally, we can derive the variance over local epistatic coefficients between mutations at sites  $U$  by using the relationship between  
 10 the variance and the raw second moments  $\text{Var}[\epsilon_{UC}] = \overline{\epsilon_{UC}^2} - \overline{\epsilon_{UC}}^2$  as well as the average variance across all possible sets of mutations  
 11  $C$  at sites  $U$  and show that it is proportional to the variance explained by genetic interactions of order higher than  $|U|$  involving sites  $U$ :

$$\begin{aligned} \text{Var}[\epsilon_U] &= \frac{1}{\binom{\alpha}{2}^{|U|}} \sum_C \text{Var}[\epsilon_{UC}] = \frac{1}{\binom{\alpha}{2}^{|U|}} \sum_C (\overline{\epsilon_{UC}^2} - \overline{\epsilon_{UC}}^2) = \frac{1}{\binom{\alpha}{2}^{|U|}} \sum_C \overline{\epsilon_{UC}^2} - \frac{1}{\binom{\alpha}{2}^{|U|}} \sum_C \overline{\epsilon_{UC}}^2 \\ &= \frac{2^{|U|}}{\alpha^\ell} \left( \frac{\alpha}{\alpha-1} \right)^{|U|} (\text{Var}^{(U)}[f] - \text{Var}[f_U]) = \frac{2^{|U|}}{\alpha^\ell} \left( \frac{\alpha}{\alpha-1} \right)^{|U|} \text{Var}_{k>|U|}^{(U)}[f]. \end{aligned} \quad (\text{S38})$$

## 12 Covariance function between sequences differing at subsets of sites

- 13 In addition to characterizing interaction structure through squared epistatic coefficients or epistatic variances of specific subsets of  
 14 sites  $U$ , it is often useful to summarize an *arbitrary* landscape  $f$  by how similar fitness values are for sequences that differ at particular  
 15 sets of sites. In particular, let  $H(x, x')$  be the subset of sites at which two sequences  $x$  and  $x'$  differ

$$H(x, x') = \{i \in S : x_i \neq x'_i\}. \quad (\text{S39})$$

- 16 We define the covariance function for a mismatch set  $D$  as

$$C_f(D) = \frac{1}{N_D} \sum_{(x, x') : H(x, x') = D} (f(x) - \bar{f})(f(x') - \bar{f}), \quad (\text{S40})$$

- 17 where  $\bar{f} = \alpha^{-\ell} \sum_x f(x)$  and  $N_D = \alpha^\ell (\alpha - 1)^{|D|}$  is the number of sequence pairs that differ at sites  $D$ .

- 18  $C_f(D)$  can be expressed as a quadratic form  $C_f(D) = \frac{1}{N_D} (f - \bar{f})^T A_D (f - \bar{f})$ , where  $A_D$  is an  $\alpha^\ell \times \alpha^\ell$  matrix given by

$$A_D(x, x') = \begin{cases} 1 & \text{if } H(x, x') = D \\ 0 & \text{if } H(x, x') \neq D. \end{cases} \quad (\text{S41})$$

- 19 This matrix can be written as a Kronecker product of single-site matrices, enabling efficient computation of this quantity without  
 20 explicit evaluation of the matrices  $A_D$  (Martí-Gómez *et al.* 2026b):

$$A_D = \bigotimes_{i=1}^{\ell} \begin{cases} 11^T - I & \text{if } i \in D \\ I & \text{if } i \notin D. \end{cases} \quad (\text{S42})$$

- 21 Moreover, if we consider the decomposition of  $f$  into its orthogonal components  $f_U = P_U f$  such that  $f = \sum_{U \in \mathcal{P}(S)} f_U$ , we can  
 22 express the covariance function as a sum of covariance over all the possible  $U$ -components:

$$C_f(D) = \frac{1}{N_D} \left( \sum_{U \in \mathcal{P}(S)} P_U f \right)^T A_D \left( \sum_{U \in \mathcal{P}(S)} P_U f \right) = \frac{1}{N_D} \sum_{U \in \mathcal{P}(S)} f_U^T P_U A_D P_U f_U = \sum_{U \in \mathcal{P}(S)} C_{f_U}(D). \quad (\text{S43})$$

Next, we note that  $P_U A_D P_U$  can be easily calculated through multiplication of their Kronecker factors individually:

$$P_U^T A_D P_U = \bigotimes_{i=1}^{\ell} \begin{cases} -P_{\text{add}} & \text{if } i \in D \wedge i \in U \\ (\alpha - 1)P_{\text{con}} & \text{if } i \in D \wedge i \notin U \\ P_{\text{add}} & \text{if } i \notin D \wedge i \in U \\ P_{\text{con}} & \text{if } i \notin D \wedge i \notin U. \end{cases} \quad (\text{S44})$$

In fact, this can be summarized as  $P_U^T A_D P_U = (-1)^{|D \cap U|} (\alpha - 1)^{|D| - |D \cap U|} P_U$ . Since  $f_U$  is in the  $U$ -subspace,  $P_U f_U = f_U$  and thus

$$\begin{aligned} C_{f_U}(D) &= \frac{1}{N_D} f_U^T P_U A_D P_U f_U = \left( \frac{1}{\alpha^{\ell} (\alpha - 1)^{|D|}} \right) \left( f_U^T (-1)^{|D \cap U|} (\alpha - 1)^{|D| - |D \cap U|} P_U f_U \right) \\ &= \alpha^{-\ell} (-1)^{|D \cap U|} (\alpha - 1)^{-|D \cap U|} f_U^T f_U = \frac{\|f_U\|^2}{(\alpha - 1)^{|U|}} \alpha^{-\ell} (-1)^{|D \cap U|} (\alpha - 1)^{|U| - |D \cap U|} \\ &= \frac{\|f_U\|^2}{(\alpha - 1)^{|U|}} w_U(D). \end{aligned} \quad (\text{S45})$$

Here, we denote by  $w_U(D)$  the contribution of genetic interactions among sites  $U$  to the covariance between sequences that differ at sites  $D$ . We refer to these quantities as subset-resolved covariance weights and they are related to the classical Krawtchouk polynomials  $w_k(d)$  arising in the Fourier analysis of fitness landscapes (Stadler 1996; Zhou et al. 2022). In particular, the functions  $w_k(d)$ , which depend only on the Hamming distance  $d = |D|$ , are recovered by summing  $w_U(D)$  over all subsets  $U$  of size  $k$ ,

$$w_k(d) = \sum_{U:|U|=k} w_U(D), \quad \text{for any } D \text{ such that } |D| = d.$$

### Gaussian random field landscapes

In the previous sections, we have explained different ways to characterize the patterns of epistatic interactions in a given fitness landscape  $f$ . Here we aim to define a probabilistic ensemble of fitness landscapes via a Gaussian distribution  $p(f)$ . This Gaussian distribution can be used as a theoretical random fitness landscape model, similar to the classical random fitness landscape models e.g. Rough Mount Fuji, NK or House of Cards landscapes (Kingman 1978; Kauffman and Weinberger 1989; Aita and Husimi 1998), or as a prior distribution for Gaussian process inference

$$p(f) = \mathcal{N}(\mu, K), \quad (\text{S46})$$

as previously proposed (Zhou et al. 2022, 2025). Here, we generally assume that the prior Gaussian distribution has zero mean i.e.  $\mu = 0$ . To define this prior, let  $Q_U$  be an  $\alpha^{\ell} \times (\alpha - 1)^{|U|}$ -dimensional matrix containing an orthonormal basis for the subspace defined by  $|U|$ -way interactions among the sites  $U$  (i.e. the column space of  $P_U$ ). We assume that the coefficients for these basis vectors are drawn independently from a zero-mean Gaussian with variance  $\lambda_U$ . In particular, let  $b_U$  be the  $(\alpha - 1)^{|U|}$ -dimensional vector of coefficients:

$$b_U \sim \mathcal{N}(0, \lambda_U I). \quad (\text{S47})$$

If we then define  $f = \sum_{U \in \mathcal{P}(S)} Q_U b_U$ , it is easy to see that  $f$  is also a zero-mean Gaussian distribution with covariance matrix  $K$  given by

$$K = \mathbb{E}_f [f f^T] = \mathbb{E}_b \left[ \left( \sum_{U \in \mathcal{P}(S)} Q_U b_U \right) \left( \sum_{U' \in \mathcal{P}(S)} Q_{U'} b_{U'} \right)^T \right] = \sum_{U, U'} Q_U \mathbb{E}_b [b_U b_{U'}^T] Q_{U'}^T = \sum_{U \in \mathcal{P}(S)} \lambda_U Q_U Q_U^T = \sum_{U \in \mathcal{P}(S)} \lambda_U P_U. \quad (\text{S48})$$

Since the columns of  $Q_U$  are orthonormal,  $P_U$  corresponds to the projection matrix into the function subspace spanned by interactions between sites in  $U$ , as given by Eq. S4. Here, we note that the entries of  $P_U(x, x')$  only depend on the Hamming distance between sequences  $x, x'$  at sites in  $U$ , but not on the alleles at which they match or differ. Thus, the covariance matrix  $K(x, x')$  only depends on the subset of sites at which  $x$  and  $x'$  differ.

Moreover, the variance of the regression coefficients  $\lambda_U$  for each subset of sites  $U$  determines the expected variance explained by interactions between exactly sites in  $U$  and is given by  $\mathbb{E}_{f \sim \mathcal{N}(0, K)} [\|f_U\|^2] = (\alpha - 1)^{|U|} \lambda_U$ .

$$\mathbb{E}_f [\|f_U\|^2] = \mathbb{E}_f [f_U^T f_U] = \mathbb{E}_f [(P_U f)^T (P_U f)] = \mathbb{E}_f [f^T P_U f]. \quad (\text{S49})$$

We note that samples from the random field model  $f \sim \mathcal{N}(0, K)$  can be drawn by first drawing  $z \sim \mathcal{N}(0, I)$  and then computing  $f = K^{\frac{1}{2}} z$ , such that

$$\mathbb{E}_f [f^T P_U f] = \mathbb{E}_z [(K^{\frac{1}{2}} z)^T P_U (K^{\frac{1}{2}} z)]. \quad (\text{S50})$$

Moreover, as the columns of  $P_U$  are eigenvectors of  $K$  with eigenvalue  $\lambda_U$  and  $P_U$  matrices are orthogonal to each other,  $K^{\frac{1}{2}} = \sum_{U \in \mathcal{P}(S)} \sqrt{\lambda_U} P_U$  and thus

$$\mathbb{E}_f [\|f_U\|^2] = \mathbb{E}_z [\lambda_U z^T P_U z] = \lambda_U \text{tr}(P_U) = \lambda_U \sum_x P_U(x, x) = \lambda_U \alpha^{\ell} (\alpha^{-\ell} (\alpha - 1)^{|U|}) = (\alpha - 1)^{|U|} \lambda_U. \quad (\text{S51})$$

## 1 A prior distribution for fitness landscapes

2 Because for sequences of length  $\ell$  these priors are defined by there  $2^\ell$  values  $\lambda_U$ , defining and interpreting this large number of the  
3 parameters can be challenging. Here, we use the relationship between the  $\Delta^{(U)}$  and the projection operators into the  $U$ -subspace  
4 shown in Eq. S21 to define a simplified prior where the  $\lambda_U$  corresponding to interactions of order higher than  $P$  are parametrized via  
5 the different  $\binom{\ell}{p}$  values of  $a_U$ :

$$\lambda_U = \begin{cases} \tilde{\lambda}_U & \text{if } |U| < P \\ \frac{1}{\alpha^{|U|} \sum_{U' \subseteq U} a_{U'}} & \text{if } |U| \geq P. \end{cases} \quad (\text{S52})$$

6 To understand the role of the parameters  $a_U$  in the prior, we can write it as

$$p(f) \propto e^{-f^T K^{-1} f} = e^{-\sum_{U:|U|<P} \frac{1}{\lambda_U} f^T P_U f - \sum_{U:|U|=P} \frac{a_U}{s_U} f^T \Delta^{(U)} f}, \quad (\text{S53})$$

7 where  $s_U$  is the number of local epistatic coefficients between sites in  $U$  and the value of  $a_U$  modulates how much the prior penalizes  
8 these coefficients for a given  $f$ . Interestingly, one can see that as  $a_U \rightarrow \infty$ ,  $\lambda_U$  approaches 0 for all function subspaces explained by  
9 interactions involving the sites in  $U$ , essentially removing interactions of order equal or higher than  $|U|$  involving the set of sites  
10  $U$ . For instance, for  $P = 2$ , setting  $a_{ij} = \infty$  enforces the assumption that there are no genetic interactions between sites  $i$  and  $j$  or,  
11 equivalently, that the effects of mutations at site  $i$  never change when introducing an additional mutation at site  $j$ . This property allows  
12 us to define prior distributions where genetic interactions of order higher than  $|U|$  are allowed, but are constrained to specific sets of  
13 sites. Similarly, if we set all the  $\tilde{\lambda}_U = c$  for  $|U| < P$  we see that the first sum in the exponent of Equation S53 becomes proportional to  
14 the squared norm of the projection of  $f$  into the null space of  $\Delta^{(P)}$ . Taking the limit  $c \rightarrow \infty$  is then equivalent to imposing a uniform  
15 prior on these directions, and for the special case  $P = 2$  we recover Equation 4 in the main text.

16 Next, we consider the expected average squared epistatic coefficients involving sites  $U$  under the prior distribution  $f \sim \mathcal{N}(0, K)$ ,  
17 parameterized by the lower-order variances  $\tilde{\lambda}_U$  and  $a_U$  for  $|U| \geq P$ . To draw samples from this prior, we can first sample  $z \sim \mathcal{N}(0, I)$   
18 and then set  $f = K^{1/2} z$ , so that  $\text{Cov}[f] = K^{1/2} (K^{1/2})^T = K$  (for any matrix square root  $K^{1/2}$ ). One way to define the matrix square  
19 root is via its known eigendecomposition  $K^{1/2} = \sum_{U \in \mathcal{P}(S)} \sqrt{\lambda_U} P_U$ . Then, we consider the average squared epistatic coefficient for any  
20  $f$  and compute its expectation when  $f$  are drawn from the prior distribution as follows

$$\mathbb{E}_f \left[ \frac{1}{s_U} f^T \Delta^{(U)} f \right] = \mathbb{E}_z \left[ \frac{1}{s_U} z^T (K^{1/2})^T \Delta^{(U)} K^{1/2} z \right] = \frac{1}{s_U} \text{tr} \left( (K^{1/2})^T \Delta^{(U)} K^{1/2} \right). \quad (\text{S54})$$

21 We compute the  $K^{1/2} \Delta^{(U)} K^{1/2}$  product by using the known eigendecompositions of  $K$  and  $\Delta^{(U)}$ :

$$\begin{aligned} \mathbb{E}_f \left[ \frac{1}{s_U} f^T \Delta^{(U)} f \right] &= \frac{1}{s_U} \text{tr} \left( \sum_{U': U \subseteq U'} \lambda_{U'} \alpha^{|U|} P_U \right) = \frac{1}{s_U} \sum_{U': U \subseteq U'} \lambda_{U'} \alpha^{|U|} \text{tr}(P_{U'}) \\ &= \frac{1}{s_U} \sum_{U': U \subseteq U'} \frac{1}{\alpha^{|U'|} \sum_{U'' \subseteq U'} a_{U''}} \alpha^{|U|} (\alpha - 1)^{|U'|} = \frac{\alpha^{|U|}}{s_U} \sum_{U': U \subseteq U'} \left( \frac{\alpha - 1}{\alpha} \right)^{|U'|} \frac{1}{\sum_{U'' \subseteq U'} a_{U''}} \\ &= \frac{\alpha^{|U|} 2^{|U|}}{\alpha^\ell (\alpha - 1)^{|U|}} \sum_{U': U \subseteq U'} \left( \frac{\alpha - 1}{\alpha} \right)^{|U'|} \frac{1}{\sum_{U'' \subseteq U'} a_{U''}} \\ &= \left( \frac{\alpha^\ell}{2^{|U|}} \left( \frac{\alpha - 1}{\alpha} \right)^{|U|} \right)^{-1} \sum_{U': U \subseteq U'} \left( \frac{\alpha - 1}{\alpha} \right)^{|U'|} \frac{1}{\sum_{U'' \subseteq U'} a_{U''}}. \end{aligned} \quad (\text{S55})$$

22 Thus, we can see that the expected average squared local epistatic coefficient for sites  $U$  does not depend only on the value  $a_U$ , but also  
23 on all other  $a_{U'}$  such that both sets have at least one site in common ( $|U \cap U'| > 0$ ). For instance, if we let  $a_U \rightarrow \infty$  for a particular  $U$ ,  
24  $\lambda_{U'}$  for all  $U'$  that include the whole set of sites in  $U$  will be set to zero, decreasing the expected average squared epistatic coefficients  
25 for sites involving sites in  $U'$ .

## 26 Relationship with the Connectedness Model

27 In this section, we review the Connectedness Model and its relationship to Local Epistasis Regression. The Connectedness Model was  
28 first introduced as a Gaussian random field model by [Reddy and Desai \(2021\)](#) to allow different loci to have different probability of  
29 being involved in epistatic interactions with mutations at other sites. Then, it was used as a prior distribution in a Gaussian process  
30 model, uncovering the set of sites that are more strongly involved in epistatic interactions and using that information for inference of  
31 high-dimensional empirical fitness landscapes ([Zhou et al. 2025](#)). Here, we show that the Connectedness Model can be derived as a  
32 particular case of Eq. S48 as a function of the variance associated to the constant component  $\tilde{\lambda}_0$  and the variance associated to the  
33 additive contribution of each site  $\tilde{\lambda}_i$

$$\lambda_U = \prod_{i \in U} \tilde{\lambda}_i \prod_{i \notin U} \tilde{\lambda}_0. \quad (\text{S56})$$

If we plug this into Eq. S48, we can see that the resulting kernel can be expressed as a product of site-specific kernels

$$\begin{aligned}
 K &= \sum_{U \in \mathcal{P}(S)} \lambda_U P_U = \sum_{U \in \mathcal{P}(S)} \prod_{i \in U} \tilde{\lambda}_i \prod_{i \notin U} \tilde{\lambda}_0 \bigotimes_{i=1}^{\ell} P_{i \in U} = \sum_{U \in \mathcal{P}(S)} \bigotimes_{i=1}^{\ell} \begin{cases} \tilde{\lambda}_0 P_{\text{con}} & i \notin U \\ \tilde{\lambda}_i P_{\text{add}} & i \in U \end{cases} \\
 &= \bigotimes_{i=1}^{\ell} (\tilde{\lambda}_0 P_{\text{con}} + \tilde{\lambda}_i P_{\text{add}}) = \bigotimes_{i=1}^{\ell} K_i,
 \end{aligned} \tag{S57}$$

with entry-wise formula given by

$$k(x, x') = \alpha^{-\ell} \prod_{i: x_i = x'_i} (\tilde{\lambda}_0 + \tilde{\lambda}_i(\alpha - 1)) \prod_{i: x_i \neq x'_i} (\tilde{\lambda}_0 - \tilde{\lambda}_i). \tag{S58}$$

This kernel can be reparametrized as a function of the prior variance  $\sigma^2$  and the correlation under this kernel parametrized by  $\mu_i = \tilde{\lambda}_i / \tilde{\lambda}_0$  assuming  $\tilde{\lambda}_0 > 0$  as in Zhou *et al.* (2025):

$$k(x, x') = \sigma^2 \prod_{i: x_i \neq x'_i} \frac{1 - \mu_i}{1 + \mu_i(\alpha - 1)}. \tag{S59}$$

This construction shows that the variance explained by epistatic interactions under the Connectedness Model is completely specified by the variance explained by the additive contribution of each individual site. In contrast, Local Epistasis Regression allows arbitrary relationships between the additive and pairwise contribution of individual sites and pairs of sites, but the variance for higher-order interactions is fully specified by the variance associated to pairwise interactions between specific pairs of sites (Eq. S52). These models also make different assumptions on how the variances associated to lower-order interactions combine to specify the variances associated to higher-order interactions. In particular, in the Connectedness Model these variances combine multiplicatively (Eq. S56), whereas in Local Epistasis Regression they are proportional to the harmonic mean of the variances associated to pairwise interactions between each possible pair of sites in  $U$

$$\lambda_U = \frac{1}{\sum_{i < j \in U} \frac{1}{\tilde{\lambda}_{ij}}}, \tag{S60}$$

where  $\tilde{\lambda}_{ij} = \frac{1}{\alpha^2 a_{ij}}$ .

### Inference of fitness landscapes under the prior

In this section, we review how to do Gaussian process inference of a complete fitness landscape  $f$  from high throughput experimental data (Martí-Gómez *et al.* 2026b). We start by defining a Gaussian prior distribution over the space of possible fitness landscapes  $p(f) \sim \mathcal{N}(0, K)$  that assigns higher probability to fitness landscapes that we believe are more plausible *a priori*.

Let  $y$  be an  $n$ -dimensional vector of measurements with known experimental Gaussian error given by the variance  $n$ -dimensional vector  $y_{\text{var}}$  for a subset of  $n \leq \alpha^\ell$  sequences  $X$ . As both the prior distribution and the likelihood function are Gaussian, the posterior distribution is also multivariate Gaussian with closed form analytical solution (Rasmussen and Williams 2008) for the mean

$$\hat{f} = K_{*X}(K_{XX} + D_{\text{var}})^{-1}y \tag{S61}$$

and covariance matrix

$$K - K_{*X}(K_{XX} + D_{\text{var}})^{-1}K_{X*}, \tag{S62}$$

where  $K_{XX}$ ,  $K_{*X}$ ,  $K_{X*}$  are submatrices of  $K$  indexed by sequences  $X$  and  $*$ , where  $*$  represents all possible sequences, and  $D_{\text{var}}$  is a diagonal matrix with the known experimental variances  $y_{\text{var}}$  along the diagonal.

Despite the simplicity of the solution, practical evaluation of these expressions becomes challenging as the number of observations increases. Traditional approaches rely to computation of the Cholesky decomposition of the  $K_{XX} + D_{\text{var}}$  matrix to then use efficient triangular solves to compute the solutions to the linear systems rather than using direct matrix inversion for higher numerical stability. However, the algorithm for computing this decomposition is  $O(n^3)$  and cannot be parallelized, which has traditionally limited the applicability of Gaussian process models to datasets with at most few thousand data points (Rasmussen and Williams 2008). In previous work, we have circumvented this limitation by leveraging the mathematical properties of the specific precision or kernel matrices, which could be expressed as polynomials in the Laplacian of the Hamming graph representing sequence space (Zhou and McCandlish 2020; Zhou *et al.* 2022; Chen *et al.* 2021; Martí-Gómez *et al.* 2026b). This property allowed us to encode these matrices as linear operators that allow computing matrix-vector products efficiently without explicitly storing them in memory, and use these operators together with iterative methods for solving systems of linear equations to scale these methods up to a few million data points.

Here we use a similar strategy by finding a representation of the kernel matrices that enables efficient computation of matrix-vector products without explicitly constructing them in memory, even if the kernel matrices presented here cannot be represented as polynomials in the Laplacian of the Hamming graph.

In the case of the Connectedness Model (Zhou *et al.* 2025), the kernel can in fact be expressed as a Kronecker product of  $\ell$  site-specific  $\alpha \times \alpha$  matrices (Eq. S57), which enables fast computation of matrix-vector products by leveraging the mixed-product property, which reduces the complexity to that of computing  $\ell$  products of an  $\alpha \times \alpha$  matrix with an  $\alpha \times \alpha^{\ell-1}$  matrix (Martí-Gómez *et al.* 2026b,a).

1 In the case of Local Epistasis Regression, the kernel matrix can be expressed as the sum of  $2^\ell$  matrices, each of which is Kronecker  
2 factorizable

$$K = \sum_{U \in \mathcal{P}(S)} \lambda_U \bigotimes_{i=1}^{\ell} P_{i \in U}. \quad (\text{S63})$$

3 While this enables computation of matrix-vector products with a total of  $2^\ell \ell \alpha \times \alpha$  by  $\alpha \times \alpha^{\ell-1}$  matrix products, the computational  
4 burden of these calculations is much larger compared with previous approaches e.g. requires  $2^\ell$  times the computation needed in the  
5 Connectedness model. However, here we note that all of the  $2^\ell$  kernel matrices represent different combinations of only two different  
6 Kronecker factors  $P_{\text{con}}$  and  $P_{\text{add}}$ , which allow us to re-use parts of the computation. In particular, we can decompose any  $P_U$  as follows

$$P_U = \bigotimes_{i=1}^{\ell} P_{i \in U} = (P_{1 \in U} \otimes I \otimes \dots \otimes I)(I \otimes P_{2 \in U} \otimes \dots \otimes I) \dots (I \otimes I \otimes \dots \otimes P_{\ell \in U}) = M_{1 \in U} M_{2 \in U} \dots M_{\ell \in U}. \quad (\text{S64})$$

7 While the  $M$  matrices are still Kronecker products of  $\ell$  matrices,  $\ell - 1$  of the factors correspond to the identity matrix and thus leave  
8 the matrices they act on unchanged, reducing the computation to a single  $\alpha \times \alpha$  by  $\alpha \times \alpha^{\ell-1}$  matrix products. Importantly,  $P_U$  matrices  
9 differing only at the factor at the first site can be computed with a single additional operation from the same intermediate result

$$P_U v = M_{1 \in U} (M_{2 \in U} \dots M_{\ell \in U} v). \quad (\text{S65})$$

10 The intermediate results can also be computed in the same fashion

$$M_{2 \in U} M_{3 \in U} \dots M_{\ell \in U} v = M_{2 \in U} (M_{3 \in U} \dots M_{\ell \in U} v) \quad (\text{S66})$$

11 so that the computation can be shared with products of the same vector with other  $P_{U'}$  for a different subset of sites  $U'$ . These  
12 computational dependencies can be represented by a bifurcating tree where nodes represent  $\alpha^\ell$ -dimensional vectors and edges  
13 represent matrix-vector products with specific  $M_{i \in U}$  matrices. The vector  $v$  is located at the root of the tree, allowing computation of  
14 the intermediate vectors at each of the nodes of the tree up to the tips containing all the  $P_U v$  for every possible  $U$ . This algorithm  
15 reduces the total number of operations from  $2^\ell \ell$  to  $2^\ell \sum_{i=0}^{\ell} 2^{-i}$ . This series converges relatively fast to  $2^{\ell+1}$ , resulting in an approximate  
16  $\ell/2$ -fold speedup even when  $\ell$  is small. For instance, for  $\ell = 8$ , the number of matrix-matrix products goes from  $8 \times 2^8 = 2048$  under  
17 the naive approach to 510, nearly achieving the expected 4-fold increase in computational efficiency.

## 18 Hyperparameter optimization

19 In order to infer a fitness landscape under a given prior distribution, we must first choose the parameters that define the properties of  
20 the prior, also known as hyperparameters. Here, we generally consider prior distributions defined over the space of possible fitness  
21 landscapes  $f$  defined by a kernel function  $k(x, x')$  that returns the covariance for any pair of sequences  $x$  and  $x'$  given by

$$k(x, x') = \sum_{U \in \mathcal{P}(S)} \lambda_U P_U(x, x'), \quad (\text{S67})$$

22 where  $P_U(x, x')$  is the covariance due to interactions between exactly  $U$  sites between sequences  $x$  and  $x'$ , which depends only on the  
23 combination of sites at which they differ  $P_U(x, x') = w_U(H(x, x'))$ ; and where the parameters  $\lambda_U$  can be free or a function of a smaller  
24 set of parameters generally called  $\theta$  ( $\lambda_U = g(\theta)$ ) e.g. Eq. S52 in Local Epistasis Regression and Eq. S56 in the Connectedness Model.

25 In this work, we use a strategy known as kernel alignment (Wang et al. 2015; Zhou et al. 2022) or Haseman-Elston regression (Hase-  
26 man and Elston 1972), in which the prior covariance approximates as closely as possible the patterns of covariance in the empirical  
27 data. Specifically, this is done by finding the parameter values  $\hat{\theta}$  that minimize the Frobenius norm of the difference between the prior  
28 predictive covariance  $K_{XX} + D_{var}$  and the empirical second moment matrix  $yy^T$ :

$$\hat{\theta} = \arg \min_{\theta} \left\| yy^T - (K_{XX}(\theta) + D_{var}) \right\|_F^2. \quad (\text{S68})$$

29 Naively solving this minimization problem is challenging, as we need to work with  $n \times n$  matrices, where the number of measured  
30 sequences  $n$  can be in the order of hundreds of thousands to millions. However, as the prior covariance between two sequences only  
31 depends on the set of sites at which they differ, the dimensionality of the problem can be reduced to a more manageable  $2^\ell$ -dimensional  
32 weighted least squares problem as follows:

$$\begin{aligned} \left\| yy^T - (K_{XX}(\theta) + D_{var}) \right\|_F^2 &= \sum_{x \in X} \left[ (y_x^2 - y_{x,var}) - k_{\theta}(\emptyset) \right]^2 + \sum_{D \neq \emptyset} \sum_{x, x': H(x, x')=D} [y_x y_{x'} - k_{\theta}(D)]^2 \\ &= N_{\emptyset} [t(\emptyset) - k_{\theta}(\emptyset)]^2 + \sum_{x \in X} \left[ (y_x^2 - y_{x,var}) - t(\emptyset) \right]^2 + \sum_{D \neq \emptyset} N_D [t(D) - k_{\theta}(D)]^2 + \sum_{D \neq \emptyset} \sum_{x, x': H(x, x')=D} [y_x y_{x'} - t(D)]^2 \\ &= \sum_D N_D [t(D) - k_{\theta}(D)]^2 + \sum_{x \in X} \left[ (y_x^2 - y_{x,var}) - t(\emptyset) \right]^2 + \sum_{D \neq \emptyset} \sum_{x, x': H(x, x')=D} [y_x y_{x'} - t(D)]^2. \end{aligned} \quad (\text{S69})$$

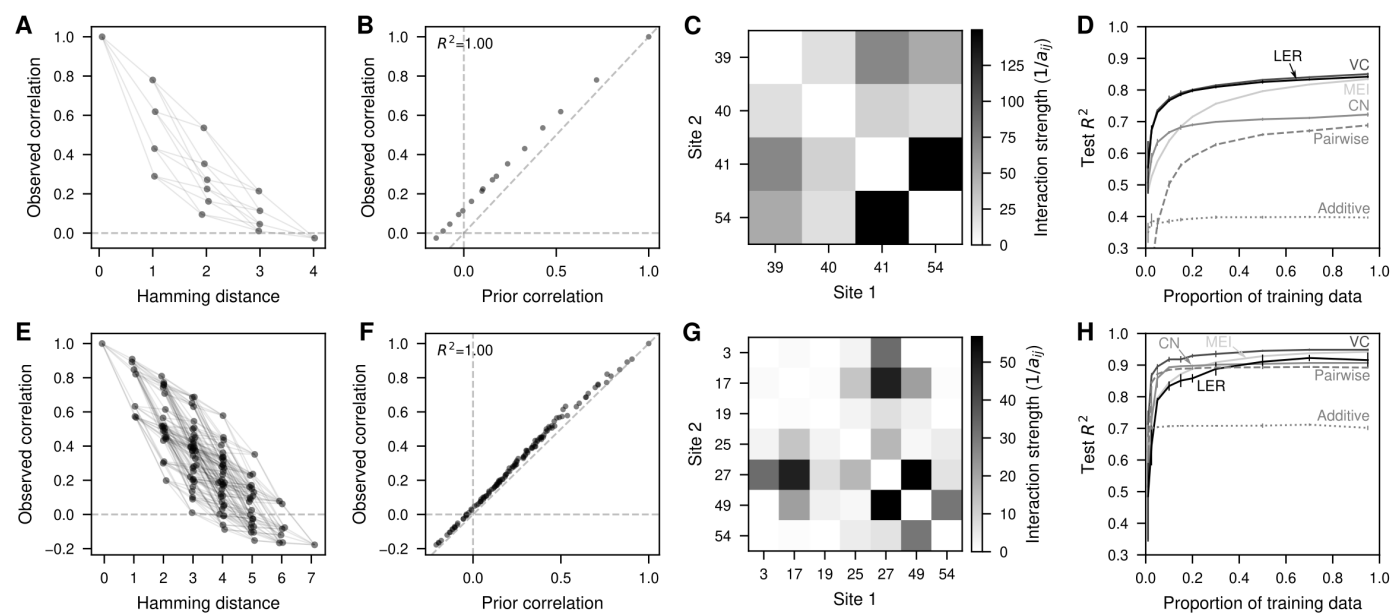
Since the last two terms are independent of  $\theta$ , we can find by optimal hyperparameter values  $\hat{\theta}$  simply as

$$\hat{\theta} = \arg \min_{\theta} \sum_{D \in \mathcal{P}(S)} N_D \left[ t(D) - \sum_{U \in \mathcal{P}(S)} \lambda_U(\theta) w_U(D) \right]^2, \quad (\text{S70})$$

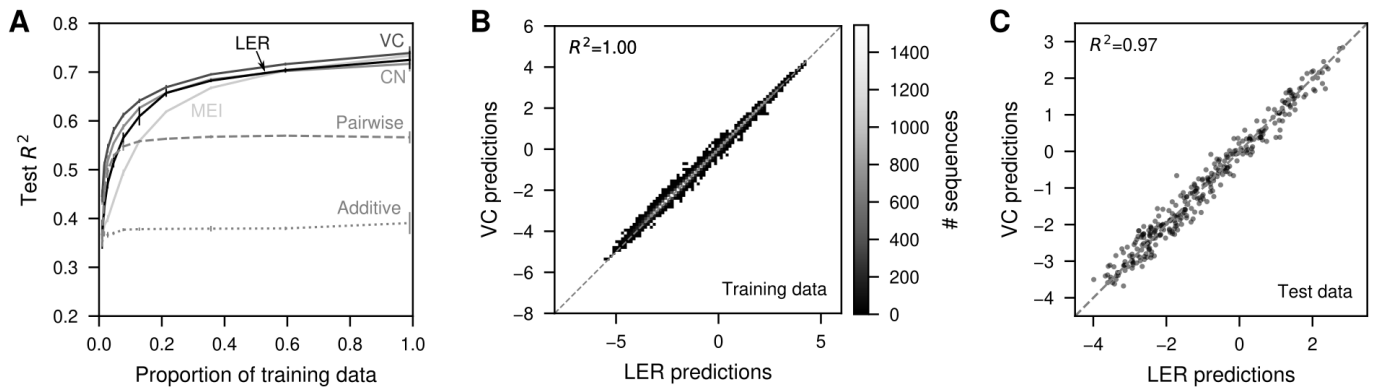
where  $t(D)$  corresponds to the second moment, related with the empirical autocovariance function  $c(D)$  ( $t(D) = c(D) + \bar{y}^2$ ) given by

$$c(D) = \begin{cases} \frac{1}{n} \sum_{x \in X} (y_x - \bar{y})^2 - \overline{yvar} & D = \emptyset \\ \frac{1}{N_D} \sum_{x, x' \in X: H(x, x') = D} (y_x - \bar{y})(y_{x'} - \bar{y}) & D \neq \emptyset. \end{cases} \quad (\text{S71})$$

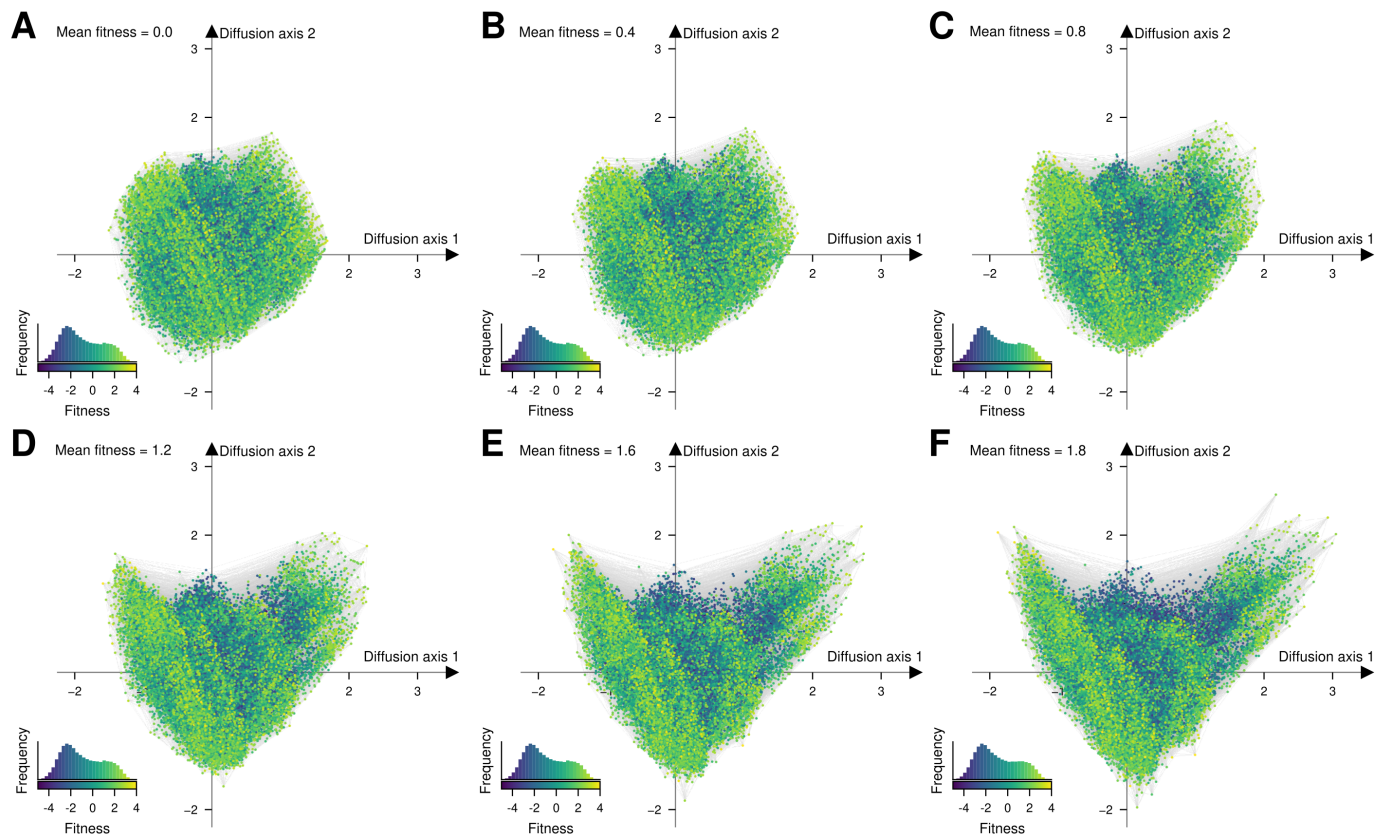
## Supplementary figures



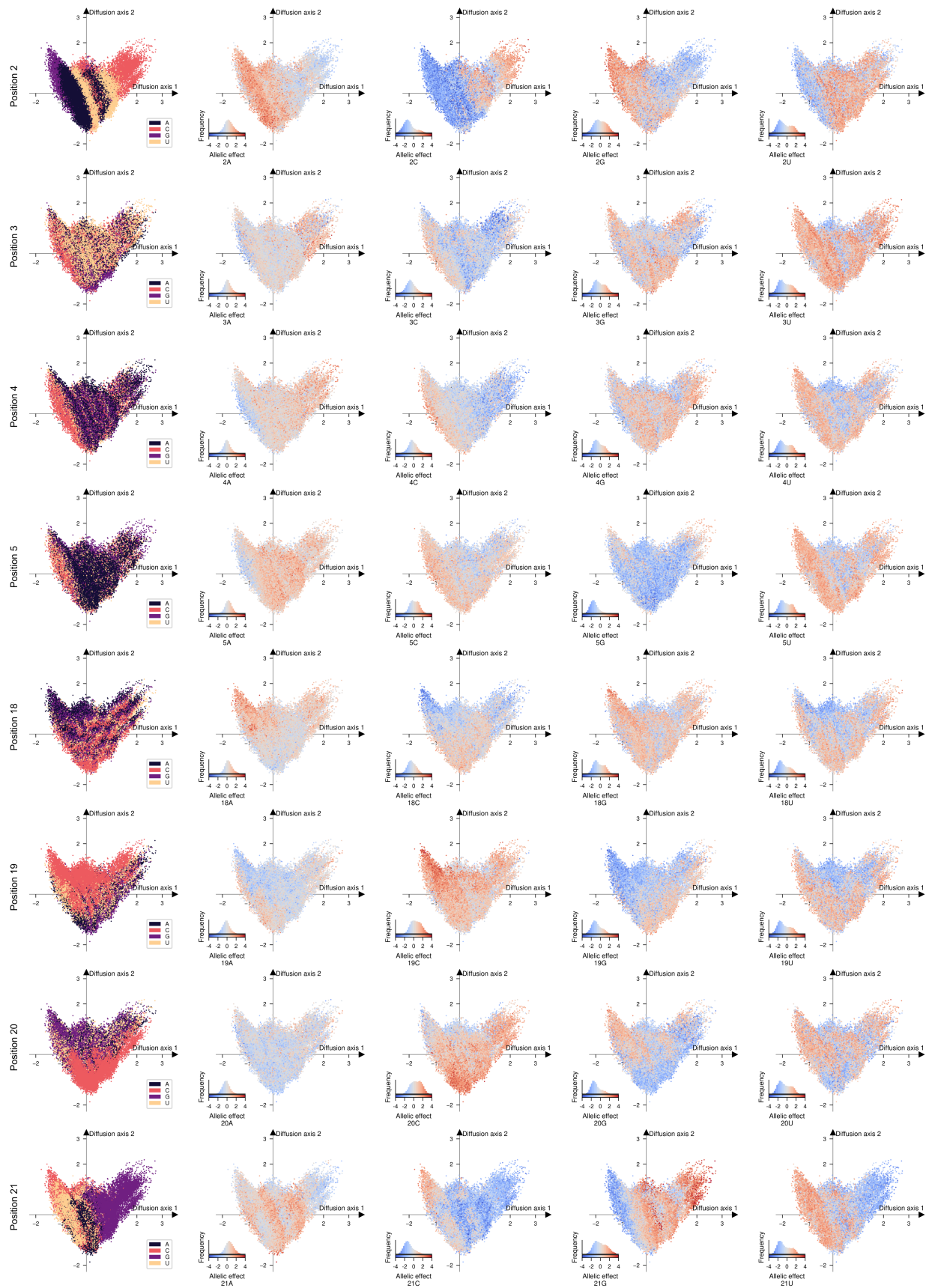
**Figure S1** Application of Local Epistasis Regression to protein datasets. (A,E) Correlation in the measured fitness values for pairs of sequences differing at each possible subset of sites  $D$  arranged according to the Hamming distance  $d = |D|$ . Each dot represents a single distance class  $D$  and are joined by lines whenever the distance classes differ by a single position from each other. (B,F) Comparison of the observed correlation values in the data and the values under the estimated prior ones using Local Epistasis Regression for every possible distance class  $D$  (each dot represents a different  $D$ ). Correlations were estimated using 80% of the data for training. (C,G) Heatmap representing the inferred model hyperparameters as  $1/a_{ij}$  for every pair of sites  $i, j$  highlighting the patterns of genetic interactions across sites under the prior. (D,H) Predictive performance evaluated by the  $R^2$  between the predicted and the measured fitness of held-out test sequences when using different amounts of training data for different models (MEI: Minimum Epistasis Interpolation, VC: Variance Component regression, CN: Connectedness Model regression, LER: Local Epistasis Regression). Predicted values are the maximum a posteriori estimate given by each method, which is equal to the posterior mean  $\hat{f}$ . Error bars represent the standard deviation across 3 different random samples for each fraction of training data. Each row represents a fitness landscape: GB1 (A,B,C,D); FYN-SH3 (E,F,G,H).



**Figure S2** Comparison of model predictions on the self-splicing intron dataset. (A) Predictive performance evaluated by the  $R^2$  between the predicted and the true fitness values of held-out test sequences when using different amounts of training data for different models (MEI: Minimum Epistasis Interpolation, VC: Variance Component regression, CN: Connectedness Model regression, LER: Local Epistasis Regression). Predicted values are the maximum a posteriori estimate given by each method, which is equal to the posterior mean  $\hat{f}$ . Error bars represent the standard deviation across 3 different random samples for each fraction of training data. (B) 2D histogram comparing the fitness landscape reconstructions of the self-splicing intron dataset under Local Epistasis Regression (LER) and Variance Component regression (VC). (C) Scatterplot comparing the predictions in held-out sequences of the self-splicing intron dataset under Local Epistasis Regression (LER) and Variance Component Regression (VC).



**Figure S3** Visualization of the inferred fitness landscape using Local Epistasis Regression under different strengths of selection, where here we quantify the strength of selection by the mean fitness achieved at stationarity, i.e. under long-term purifying selection. Every dot represents one of the possible  $4^8$  possible sequences and is colored according to the predicted fitness. The inset represents the phenotypic distribution along with their corresponding color in the map. Sequences are laid out according to the first two Diffusion axes and dots are plotted in order according to Diffusion axis 3.



**Figure S4** Visualizing alleles and allelic preferences across the fitness landscape visualization. Visualization of the inferred fitness landscape using Local Epistasis Regression. Every dot represents one of the possible  $4^8$  possible sequences and is colored according to the allele (first column) or the difference in the fitness of the sequence obtained when placing a specific allele at a specific position relative to the average fitness of the four possible alleles (four last columns). The inset represents the allelic effect distribution along with their corresponding color in the map. Sequences are laid out according to the first two Diffusion axes and dots are plotted in order according to Diffusion axis 3.