

1 **Title**

2 MAVE-NN: learning genotype-phenotype maps from multiplex assays of variant effect

3 **Authors**

4 Ammar Tareen¹, Mahdi Kooshkbaghi¹, Anna Posfai¹, William T. Ireland^{2,3}, David M.

5 McCandlish¹, Justin B. Kinney^{1,*}

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

7 ² Department of Physics, California Institute of Technology, Pasadena, CA, 91125

8 ³ Present address: Department of Applied Physics, Harvard University, Cambridge, MA, 02134

9 * To whom correspondence should be addressed: jkinney@cshl.edu

10

11 **Abstract**

12 Multiplex assays of variant effect (MAVEs) are a family of methods that includes deep
13 mutational scanning (DMS) experiments on proteins and massively parallel reporter assays
14 (MPRAs) on gene regulatory sequences. However, a general strategy for inferring quantitative
15 models of genotype-phenotype (G-P) maps from MAVE data is lacking. Here we introduce
16 MAVE-NN, a neural-network-based Python package that implements a broadly applicable
17 information-theoretic framework for learning G-P maps—including biophysically interpretable
18 models—from MAVE datasets. We demonstrate MAVE-NN in multiple biological contexts, and
19 highlight the ability of our approach to deconvolve mutational effects from otherwise
20 confounding experimental nonlinearities and noise.

21

22 **Main Text**

23 **Introduction**

24 Over the last decade, the ability to quantitatively study genotype-phenotype (G-P) maps
25 has been revolutionized by the development of multiplex assays of variant effect (MAVEs),
26 which can measure molecular phenotypes for thousands to millions of genotypic variants in
27 parallel.^{1,2} MAVE is an umbrella term that describes a diverse set of experimental methods,
28 three examples of which are illustrated in **Fig. 1**. Deep mutational scanning (DMS) experiments³
29 are a type of MAVE commonly used to study protein sequence-function relationships. These
30 assays work by linking variant proteins to their coding sequences, either directly or indirectly,
31 then using deep sequencing to assay which variants survive a process of activity-dependent
32 selection (e.g., **Fig. 1a**). Massively parallel reporter assays (MPRAs) are another major class of
33 MAVE, and are commonly used to study DNA or RNA sequences that regulate gene expression
34 at a variety of steps, including transcription, mRNA splicing, cleavage and polyadenylation,
35 translation, and mRNA decay.⁴⁻⁷ MPRAs typically rely on either an RNA-seq readout of barcode
36 abundances (**Fig. 1c**) or the sorting of cells expressing a fluorescent reporter gene (**Fig. 1e**).

37 Most computational methods for analyzing MAVE data have focused on accurately
38 quantifying the activity of individual assayed sequences.⁸⁻¹⁴ However, MAVE measurements like
39 enrichment ratios or cellular fluorescence levels usually cannot be interpreted as providing
40 direct quantification of biologically meaningful activities, due to the presence of experiment-
41 specific nonlinearities and noise. Moreover, MAVE data is usually incomplete, as one often
42 wishes to understand G-P maps over vastly larger regions of sequence space than can be
43 exhaustively assayed. The explicit quantitative modeling of G-P maps can address both the
44 indirectness and incompleteness of MAVE measurements.^{1,15} The goal here is to determine a
45 mathematical function that, given a sequence as input, will return a quantitative value for that

46 sequence's molecular phenotype. Such quantitative modeling has been of great interest since
47 the earliest MAVE methods were developed,¹⁶⁻¹⁸ but no general-use software has yet been
48 described for inferring G-P maps of arbitrary functional form from MAVE data.

49 Here we introduce a unified conceptual framework for the quantitative modeling of
50 MAVE data. This framework is based on the use of latent phenotype models, which assume that
51 each assayed sequence has a well-defined latent phenotype (specified by the G-P map), of
52 which the MAVE experiment provides an indirect readout (described by the measurement
53 process). The quantitative forms of both the G-P map and the measurement process are then
54 inferred from MAVE data simultaneously. We further introduce an information-theoretic
55 approach for separately assessing the performance of the G-P map and the measurement
56 process components of latent phenotype models. This strategy is implemented in an easy-to-
57 use open-source Python package called MAVE-NN, which represents latent phenotype models
58 as neural networks and infers the parameters of these models from MAVE data using a
59 TensorFlow 2 backend.¹⁹

60 In what follows, we expand on this unified MAVE modeling strategy and apply it to a
61 diverse array of DMS and MPRA datasets. Doing so, we find that MAVE-NN provides
62 substantial advantages over other MAVE modeling approaches. Our results also highlight the
63 substantial benefits of including sequence variants with multiple mutations in assayed sequence
64 libraries, as doing so allows MAVE-NN to deconvolve the features of the G-P map from
65 potentially confounding effects of experimental nonlinearities and noise. Importantly, we find that
66 including just a modest number of multiple-mutation variants in a MAVE experiment can be
67 beneficial even when one is primarily interested in the effects of single mutations. Finally, we
68 illustrate how the ability of MAVE-NN to train custom G-P maps can shed light on the
69 biophysical mechanisms of gene regulation.

70 **Results**

71 ***Latent phenotype modeling strategy***

72 MAVE-NN supports the analysis of MAVE data on DNA, RNA, and protein sequences,
73 and can accommodate either continuous or discrete measurement values. Given a set of
74 sequence-measurement pairs, MAVE-NN aims to infer a probabilistic mapping from sequence
75 to measurement. Our primary enabling assumption, which is encoded in the structure of the
76 latent phenotype model (**Fig. 2a**), is that this mapping occurs in two stages. Each sequence is
77 first mapped to a latent phenotype by a deterministic G-P map, then this latent phenotype is
78 mapped to possible measurement values via a stochastic measurement process. During
79 training, the G-P map and measurement process are simultaneously learned by maximizing a
80 regularized form of likelihood.

81 MAVE-NN includes four types of built-in G-P maps: additive, neighbor, pairwise, and
82 black box. Additive G-P maps assume that each character at each position within a sequence
83 contributes independently to the latent phenotype. Neighbor G-P maps incorporate interactions
84 between adjacent (i.e., nearest-neighbor) characters in a sequence, while pairwise G-P maps
85 include interactions between all pairs of characters in a sequence regardless of the distance
86 separating the characters in each pair. Black box G-P maps have the form of a densely
87 connected multilayer perceptron, the specific architecture of which can be controlled by the
88 user. MAVE-NN also supports custom G-P maps that can be used, e.g., to represent specific
89 biophysical hypotheses about the mechanisms of sequence function.

90 To handle both discrete and continuous measurement values, two different strategies for
91 modeling measurement processes are provided. Measurement process agnostic (MPA)
92 regression uses techniques from the biophysics literature^{15,16,20,21} to analyze MAVE datasets
93 that report discrete measurements. Here the measurement process is represented by an

94 overparameterized neural network that takes the latent phenotype value as input and outputs
95 the probability of each possible measurement value (**Fig. 2b**). Global epistasis (GE) regression
96 (**Fig. 2c**), by contrast, leverages ideas previously developed in the evolution literature for
97 analyzing datasets that contain continuous measurements,²²⁻²⁵ and is becoming an increasingly
98 popular strategy for modeling DMS data.²⁶⁻²⁸ Here, the latent phenotype is nonlinearly mapped
99 to a prediction that represents the most probable measurement value. A noise model is then
100 used to describe the distribution of likely deviations from this prediction. MAVE-NN supports
101 both homoscedastic and heteroscedastic noise models based on three different classes of
102 probability distribution: Gaussian, Cauchy, and skewed-t. We note that the skewed-t distribution,
103 introduced by Jones and Faddy,²⁹ reduces to Gaussian and Cauchy distributions in certain limits
104 while also accommodating asymmetric experimental noise. **Fig. 2d** shows an example of a GE
105 measurement process with a heteroscedastic skewed-t noise model.

106 Readers should note that the current implementation of MAVE-NN places certain
107 constraints on input data and model architecture. Input sequences must be the same length,
108 and when analyzing continuous data, only scalar measurements (as opposed to vectors of
109 multiple measurements) can be used to train models. In addition, because our method for
110 learning the form of experimental nonlinearities depends on observing how multiple mutations
111 combine, MAVE-NN's functionality is more limited when analyzing MAVE libraries that comprise
112 only single-mutation variants. More information on these constraints and the reasons behind
113 them can be found below in the section "Constraints on datasets and models".

114 ***Information-theoretic measures of model performance***

115 We further propose three distinct quantities for assessing the performance of latent
116 phenotype models: intrinsic information, predictive information, and variational information (**Fig.**
117 **2e**). These quantities come from information theory and are motivated by thinking of G-P maps

118 in terms of information compression. In information theory, a quantity called mutual information
119 quantifies the amount of information that the value of one variable communicates about the
120 value of another.^{30,31} Mutual information is symmetric, nonnegative, is measured in units of
121 “bits”, and is equal to 0 bits only if the two variables are independent. Alternatively, if knowing
122 the value of one variable allows you to narrow down the value of the other variable to one of two
123 possibilities that would otherwise be equally likely, the mutual information between these two
124 variables will be 1.0 bits. If the value is narrowed down to one of four otherwise equally likely
125 possibilities, the mutual information will be 2.0 bits. Narrowing down to one of with eight
126 possibilities will yield 3.0 bits, and so on. But importantly, mutual information does not require
127 that the relationship between two variables in question be so clean cut, and mutual information
128 can in fact be computed between any two types of variables—discrete, continuous, multi-
129 dimensional, etc.. This property makes the information-based quantities we propose applicable
130 to all MAVE datasets, regardless of the specific type of experimental readout used. By contrast,
131 many of the standard model performance metrics have restricted domains of applicability:
132 accuracy can only be applied to data with discrete labels, R^2 can only be applied to data with
133 univariate continuous labels, and so on. We note, however, that estimating mutual information
134 and related quantities from finite data is nontrivial and that MAVE-NN uses a variety of
135 approaches to do this.

136 Intrinsic information, I_{int} , is the mutual information between the sequences and
137 measurements contained within a MAVE dataset. This quantity provides a benchmark against
138 which to compare the performance of inferred G-P maps. Predictive information, I_{pre} , is the
139 mutual information between MAVE measurements and the latent phenotype values predicted by
140 a G-P map of interest. This quantifies how well the G-P map preserves sequence-encoded
141 information that is determinative of experimental measurements. When evaluated on test data,
142 I_{pre} is bounded above by I_{int} , and equality is realized only when the latent phenotype losslessly

143 encodes relevant sequence-encoded information. Variational information, I_{var} , is a linear
144 transformation of log likelihood that provides a variational lower bound on I_{pre} .³²⁻³⁴ The
145 difference between I_{pre} and I_{var} quantifies how accurately the inferred measurement process
146 matches the observed distribution of measurements and latent phenotypes (see **Supplemental**
147 **Information**).

148 MAVE-NN infers model parameters by maximizing a (lightly) regularized form of
149 likelihood. These computations are performed using the standard backpropagation-based
150 training algorithms provided within the TensorFlow 2 backend. With certain caveats noted (see
151 **Methods**), this optimization procedure maximizes I_{pre} while avoiding the costly estimates of
152 mutual information at each iteration that have hindered the adoption of previous mutual-
153 information-based modeling strategies.¹⁶

154 ***Application: deep mutational scanning assays***

155 We now demonstrate the capabilities of MAVE-NN on three DMS datasets, starting with
156 the study of Olson et al.³⁵ on pairwise epistasis in protein G. Here the authors measured the
157 effects of all single and nearly all double mutations to residues 2-56 of the IgG binding domain.
158 This domain, called GB1, has long served as a model system for studying protein sequence-
159 function relationships. To assay the binding of GB1 variants to IgG, the authors combined
160 mRNA display with ultra-high-throughput DNA sequencing (**Fig. 1a**). The resulting dataset
161 reports log enrichment values for all 1,045 single- and 530,737 double-mutant GB1 variants
162 (**Fig. 1b**).

163 Inspired by the work of Otwinowski et al.,²⁵ we used MAVE-NN to infer a latent
164 phenotype model comprising an additive G-P map and a GE measurement process. This
165 inference procedure required only about 5 minutes on a single node of a computer cluster

166 **(Supplemental Fig. S1). Fig. 3a** illustrates the inferred additive G-P map via the effects that
167 every possible single-residue mutation has on the latent phenotype. From this heatmap of
168 additive effects, we can immediately identify all of the critical GB1 residues, including the IgG
169 interacting residues at 27, 31, and 43.³⁵ We also observe that missense mutations to proline
170 throughout the GB1 domain tend to negatively impact IgG binding, as expected due to this
171 amino acid's exceptional conformational rigidity. **Fig. 3b** illustrates the corresponding GE
172 measurement process, revealing a sigmoidal relationship between log enrichment
173 measurements and the latent phenotype values predicted by the G-P map. Nonlinearities like
174 this are ubiquitous in DMS data due to the presence of background and saturation effects.
175 Unless they are explicitly accounted for in one's quantitative modeling efforts, as they are here,
176 these nonlinearities can greatly distort the parameters of inferred G-P maps. **Fig. 3c** shows that
177 accounting for this nonlinearity yields predictions that correlate quite well with measurement
178 values. Moreover, every latent phenotype model inferred by MAVE-NN can be used as a MAVE
179 dataset simulator (see **Methods**). By analyzing simulated data generated by our inferred model
180 for this GB1 experiment, we further observed that MAVE-NN can accurately and robustly
181 recover the GE nonlinearity and ground-truth G-P map parameters (**Supplementary Fig. S1**).

182 **Fig. 3d** summarizes the values of our information-theoretic metrics for model
183 performance. On held-out test data, we find that $I_{\text{var}} = 2.178 \pm 0.027$ bits and $I_{\text{pre}} = 2.225 \pm$
184 0.017 bits. The similarity of these two values suggests that the inferred GE measurement
185 process, which includes a heteroscedastic skewed-t noise model, has nearly sufficient accuracy
186 to fully describe the distribution of residuals. We further find that 2.741 ± 0.013 bits $\leq I_{\text{int}} \leq$
187 3.215 ± 0.007 bits (see **Methods**), meaning that the inferred G-P map accounts for 69%-81% of
188 the total sequence-dependent information in the dataset. While this performance is impressive,
189 the additive G-P map evidently misses some relevant sequence features. This observation
190 motivates the more complex biophysical model for GB1 discussed later in **Results**.

191 The ability of MAVE-NN to deconvolve experimental nonlinearities from additive G-P
192 maps requires that some of the assayed sequences contain multiple mutations. This is because
193 such nonlinearities are inferred by reconciling the effects of single mutations with the effects
194 observed for combinations of two or more mutations. To investigate how many multiple-mutation
195 variants are required, we performed GE inference on subsets of the GB1 dataset containing all
196 1,045 single-mutation sequences and either 50,000, 5,000, or 500 double-mutation sequences
197 (see **Methods**). The shapes of the resulting GE nonlinearities are illustrated in **Figs. 3e-g**.
198 Remarkably, MAVE-NN is able to recover the underlying nonlinearity using only about 500
199 randomly selected double mutants, which represent only ~0.1% of all possible double mutants.
200 The analysis of simulated data also supports the ability to accurately recover ground-truth model
201 predictions using highly reduced datasets (**Supplemental Fig. S1**). These findings have
202 important implications for the design of DMS experiments: even if one only wants to determine
203 an additive G-P map, including a modest number of multiple-mutation sequences in the assayed
204 library is often advisable because it may allow the removal of artifactual nonlinearities.

205 To test the capabilities of MAVE-NN on less complete DMS datasets, we analyzed
206 recent experiments on amyloid beta ($A\beta$)³⁶ and TDP-43,³⁷ both of which exhibit aggregation
207 behavior in the context of neurodegenerative diseases. In these experiments, protein
208 functionality was assayed using selective growth in genetically modified *Saccharomyces*
209 *cerevisiae*: Seuma et al.³⁶ performed a selection against $A\beta$ toxicity, whereas Bolognesi et al.³⁷
210 positively selected for TDP-43 aggregation. Like with GB1, the variant libraries used in these
211 two experiments included a substantial number of multiple-mutation sequences: 499 single- and
212 15,567 double-mutation sequences for $A\beta$; 1,266 single- and 56,730 double-mutation
213 sequences for TDP-43. But unlike with GB1, these datasets are highly incomplete due to the
214 use of mutagenic PCR (for $A\beta$) or doped oligo synthesis (TDP-43) for variant library
215 construction.

216 We used MAVE-NN to infer additive G-P maps from these two datasets, adopting the
217 same type of latent phenotype model used for GB1. **Fig. 4a** illustrates the additive G-P map
218 inferred from aggregation measurements of A β variants. In agreement with the original study,
219 we see that most amino acid mutations between positions 30-40 have a negative effect on
220 nucleation, suggesting that this region plays a major role in nucleation behavior. **Fig. 4b** shows
221 the corresponding measurement process (see also **Supplemental Information Fig. S2**). Even
222 though these data are much sparser than the GB1 data, the inferred model performs well on
223 held-out test data ($I_{\text{var}} = 1.142 \pm 0.065$ bits, $I_{\text{pre}} = 1.187 \pm 0.050$ bits, $R^2 = 0.763 \pm 0.024$).
224 Similarly, **Figs. 4c-d** show the G-P map parameters and GE measurement process inferred
225 from toxicity measurements of TDP-43 variants, revealing among other things the toxicity-
226 determining hot-spot observed by Bolognesi et al.³⁷ at positions 310-340. The resulting latent
227 phenotype model performs well on held-out test data ($I_{\text{var}} = 1.834 \pm 0.035$ bits, $I_{\text{pre}} = 1.994 \pm$
228 0.023 bits, $R^2 = 0.914 \pm 0.007$).

229 ***Application: a massively parallel splicing assay***

230 Exon/intron boundaries are defined by 5' splice sites (5'ss), which bind the U1 snRNP
231 during the initial stages of spliceosome assembly. To investigate how 5'ss sequence
232 quantitatively controls alternative mRNA splicing, Wong et al.³⁸ used a massively parallel
233 splicing assay (MPSA) to measure percent-spliced-in (PSI) values for nearly all 32,768 possible
234 5'ss of the form NNN/GYNNNN in three different genetic contexts (**Fig. 1c,d**). Applying MAVE-
235 NN to data from the BRCA2 exon 17 context, we inferred four different types of G-P maps:
236 additive, neighbor, pairwise, and black box. As with GB1, these G-P maps were each inferred
237 using GE regression with a heteroscedastic skewed-t noise model. For comparison, we also
238 inferred an additive G-P map using the epistasis package of Sailer and Harms.²⁴

239 **Fig. 5a** compares the performance of these G-P map models on held-out test data, while
240 **Figs. 5b-d** illustrate the corresponding inferred measurement processes. We observe that the
241 additive G-P map inferred using the epistasis package²⁴ exhibits less predictive information
242 ($I_{\text{pre}} = 0.180 \pm 0.011$ bits) than the additive G-P map found using MAVE-NN ($P = 3.8 \times 10^{-6}$,
243 two-sided z-test). This is likely because the epistasis package estimates the parameters of the
244 additive G-P map prior to estimating the GE nonlinearity. We also note that, while the epistasis
245 package provides a variety of options for modeling the GE nonlinearity, none of these options
246 appear to work as well as our mixture-of-sigmoids approach (compare **Figs. 5b,c**). This finding
247 again demonstrates that the accurate inference of G-P maps requires the explicit and
248 simultaneous modeling of experimental nonlinearities.

249 We also observe that increasingly complex G-P maps exhibit increased accuracy. For
250 example, the additive G-P map gives $I_{\text{pre}} = 0.257 \pm 0.013$ bits, whereas the pairwise G-P map
251 (**Figs. 5e,f**) attains $I_{\text{pre}} = 0.374 \pm 0.014$ bits. Using MAVE-NN's built-in parametric bootstrap
252 approach for quantifying parameter uncertainty, we find that both the additive and pairwise G-P
253 map parameters are very precisely determined (see **Supplemental Information Fig. S3**). The
254 black box G-P map, which is comprised of 5 densely connected hidden layers of 10 nodes each,
255 performed the best of all four G-P maps, achieving $I_{\text{pre}} = 0.458 \pm 0.015$ bits. Remarkably, this
256 last predictive information value exceeds the lower bound of $I_{\text{int}} \geq 0.462 \pm 0.009$ bits, which was
257 estimated from replicate experiments (see **Methods**). We thus conclude that pairwise
258 interaction models are not flexible enough to fully account for how 5' splice sequences control
259 splicing. More generally, these results underscore the need for software that is capable of
260 inferring and assessing a variety of different G-P maps through a uniform interface.

261 **Application: biophysically interpretable G-P maps**

262 Biophysical models, unlike the phenomenological models considered thus far, have
263 mathematical structures that reflect specific hypotheses about how sequence-dependent
264 interactions between macromolecules mechanistically define G-P maps. Thermodynamic
265 models, which rely on a quasi-equilibrium assumption, are the most commonly used type of
266 biophysical model.^{39–41} Previous studies have shown that precise thermodynamic models can
267 be inferred from MAVE datasets,¹⁶ but no software intended for use by the broader MAVE
268 community has yet been developed for doing this. MAVE-NN meets this need by enabling the
269 inference of custom G-P maps. We now demonstrate this biophysical modeling capability in the
270 contexts of protein-ligand binding (using DMS data; **Fig. 1a**) and bacterial transcriptional
271 regulation (using sort-seq MPRA data; **Fig. 1e**). An expanded discussion of how these models
272 are mathematically formulated and specified within MAVE-NN is provided in the “Biophysical
273 modeling” section of **Supplemental Information**.

274 Otwinowski⁴² showed that a three-state thermodynamic G-P map (**Fig. 6a**), one that
275 accounts for GB1 folding energy in addition to GB1-IgG binding energy,⁴³ can explain the DMS
276 data of Olson et al.³⁵ better than a simple additive G-P map does. This biophysical model
277 subsequently received impressive confirmation in the work of Nisthal et al.,⁴⁴ who measured the
278 thermostability of 812 single-mutation GB1 variants. We tested the ability of MAVE-NN to
279 recover the same type of thermodynamic model that Otwinowski had inferred using custom
280 analysis scripts. Our analysis yielded a G-P map with significantly improved performance on the
281 data of Olson et al. ($I_{\text{var}} = 2.303 \pm 0.013$ bits, $I_{\text{pre}} = 2.357 \pm 0.007$ bits, $R^2 = 0.947 \pm 0.001$)
282 relative to the additive G-P map of **Fig. 3**. **Fig. 6b** shows the two inferred energy matrices that
283 respectively describe the effects of every possible single-residue mutation on the Gibbs free
284 energies of protein folding and protein-ligand binding. The folding energy predictions of our
285 model also correlate as well with the data of Nisthal et al. ($R^2 = 0.570 \pm 0.049$) as the

286 predictions of Otwinowski's model does ($R^2 = 0.515 \pm 0.056$). This demonstrates that MAVE-NN
287 can infer accurate and interpretable quantitative models of protein biophysics.

288 To test MAVE-NN's ability to infer thermodynamic models of transcriptional regulation,
289 we re-analyzed the MPRA data of Kinney et al.,¹⁶ in which random mutations to a 75 bp region
290 of the *Escherichia coli lac* promoter were assayed. This promoter region binds two regulatory
291 proteins, σ^{70} RNA polymerase (RNAP) and the transcription factor CRP. As in Kinney et al.,¹⁶
292 we proposed a four-state thermodynamic model that quantitatively explains how promoter
293 sequences control transcription rate (**Fig. 6c**). The parameters of this G-P map include the
294 Gibbs free energy of interaction between CRP and RNAP (ΔG_I), as well as energy matrices that
295 describe the CRP-DNA and RNAP-DNA interaction energies. Because the sort-seq MPRA of
296 Kinney et al. yielded discrete measurement values (**Figs. 1e,f**), we used an MPA measurement
297 process in our latent phenotype model (**Fig. 6d**). The biophysical parameter values we thus
298 inferred (**Fig. 6e**), including a CRP-RNAP interaction energy of $\Delta G_I = -2.598 \pm 0.018$ kcal/mol,
299 largely match those of Kinney et al., but were obtained far more rapidly (in ~10 min versus
300 multiple days) thanks to the use of stochastic gradient descent rather than Metropolis Monte
301 Carlo.

302 **Constraints on datasets and models**

303 As stated above, MAVE-NN places certain limitations on both input datasets and latent
304 phenotype models. Some of these constraints have been adopted to simplify the initial release
305 of MAVE-NN and can potentially be relaxed in future updates. Others reflect fundamental
306 mathematical properties of latent phenotype models. Here we summarize the primary
307 constraints users should be aware of.

308 MAVE-NN currently requires that all input sequences be the same length. This constraint
309 has been adopted because a large fraction of MAVE datasets have this form, and all of the built-

310 in G-P maps operate only on fixed-length sequences. Users who wish to analyze variable length
311 sequences can still do so by padding the ends of sequences with dummy characters.
312 Alternatively, users can provide a multiple-sequence alignment as input and include the gap
313 character as one of the characters to consider when training models.

314 As stated above, MAVE-NN can analyze MAVE datasets that have either continuous or
315 discrete measurements. At present, both types of measurements must be one-dimensional, i.e.,
316 users cannot fit a single model to vectors of multiple measurements (e.g., joint measurements of
317 protein binding affinity and protein stability, as in Ref. ²⁷). This constraint has been adopted only
318 to simplify the user interface of the initial release. It is not a fundamental limitation of latent
319 phenotype models and is scheduled to be relaxed in upcoming versions of MAVE-NN.

320 The current implementation of MAVE-NN also supports only one-dimensional latent
321 phenotypes (though the latent phenotype of custom G-P maps can depend on multiple
322 precursor phenotypes, e.g., binding energy or folding energy). This restriction was made
323 because accurately interpreting multi-dimensional latent phenotypes is substantially more
324 fraught than interpreting one-dimensional latent phenotypes, and we believe that additional
325 computational tools need to be developed to facilitate such interpretation. That being said, the
326 mathematical form of latent phenotype models is fully compatible with multi-dimensional latent
327 phenotypes. Indeed, this modeling strategy has been used in other work,^{20,27,28,45} and we plan to
328 enable this functionality in future updates to MAVE-NN.

329 More fundamental constraints come into play when analyzing MAVE data that contains
330 only single-mutation variants. In such experiments, the underlying effects of individual mutations
331 are hopelessly confounded by the biophysical, physiological, and experimental nonlinearities
332 that may be present. By contrast, when the same mutation is observed in multiple genetic
333 backgrounds, MAVE-NN can use systematic differences in the mutational effects observed

334 between stronger and weaker backgrounds to remove these confounding influences. Thus, for
335 datasets that comprise only single-mutant effects, we limit MAVE-NN to inferring only additive
336 G-P maps using GE regression, and while the noise model in the GE measurement process is
337 allowed to be heteroscedastic, the nonlinearity is constrained to be linear.

338 We emphasize that, in practice, only a modest number of multiple-mutant variants are
339 required for MAVE-NN to learn the form of a non-linear measurement process (see Fig. 3e-g).
340 In this way, including a small fraction of the possible double-mutation variants in MAVE libraries
341 can be beneficial even just for determining the effects of single mutations. Adding such non-
342 comprehensive sets of double mutants to MAVE libraries is experimentally straight-forward, and
343 our numerical experiments suggest that assaying roughly the same number of double-mutation
344 variants as single-mutation variants should often suffice. We therefore recommend that
345 experimentalists—even those primarily interested in the effects of single mutations—consider
346 augmenting their MAVE libraries with a small subset of double-mutation variants.

347 **Discussion**

348 In this work we have presented a unified strategy for inferring quantitative models of G-P
349 maps from diverse MAVE datasets. At the core of our approach is the conceptualization of G-P
350 maps as a form of information compression, i.e., the G-P map first compresses an input
351 sequence into a latent phenotype value, which the MAVE then reads out indirectly via a noisy
352 nonlinear measurement process. By explicitly modeling this measurement process, one can
353 remove potentially confounding effects from the G-P map, as well as accommodate diverse
354 experimental designs. We have also introduced three information-theoretic metrics for
355 assessing the performance of the resulting models. These capabilities have been implemented
356 within an easy-to-use Python package called MAVE-NN.

357 We have demonstrated the capabilities of MAVE-NN in diverse biological contexts,
358 including in the analysis of both DMS and MPRA data. We have also demonstrated the superior
359 performance of MAVE-NN relative to the epistasis package of Sailer and Harms.²⁴ Along the
360 way, we observed that MAVE-NN can deconvolve experimental nonlinearities from additive G-P
361 maps when a relatively small number of sequences containing multiple mutations are included
362 in the assayed libraries. This capability provides a compelling reason for experimentalists to
363 include such sequences in their MAVE libraries, even if they are primarily interested in the
364 effects of single mutations. Finally, we showed how MAVE-NN can learn biophysically
365 interpretable G-P maps from both DMS and MPRA data.

366 MAVE-NN thus fills a critical need in the MAVE community, providing user-friendly
367 software capable of learning quantitative models of G-P maps from diverse MAVE datasets.
368 MAVE-NN has a streamlined user interface and is readily installed from PyPI by executing “pip
369 install mavenn” at the command line. Comprehensive documentation and step-by-step tutorials
370 are available at <http://mavenn.readthedocs.io>.

371 **Acknowledgements.** This work was supported by NIH grant 1R35GM133777 (awarded to
372 JBK), NIH Grant 1R35GM133613 (awarded to DMM), an Alfred P. Sloan Research Fellowship
373 (awarded to DMM), a grant from the CSHL/Northwell Health partnership, and funding from the
374 Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory.

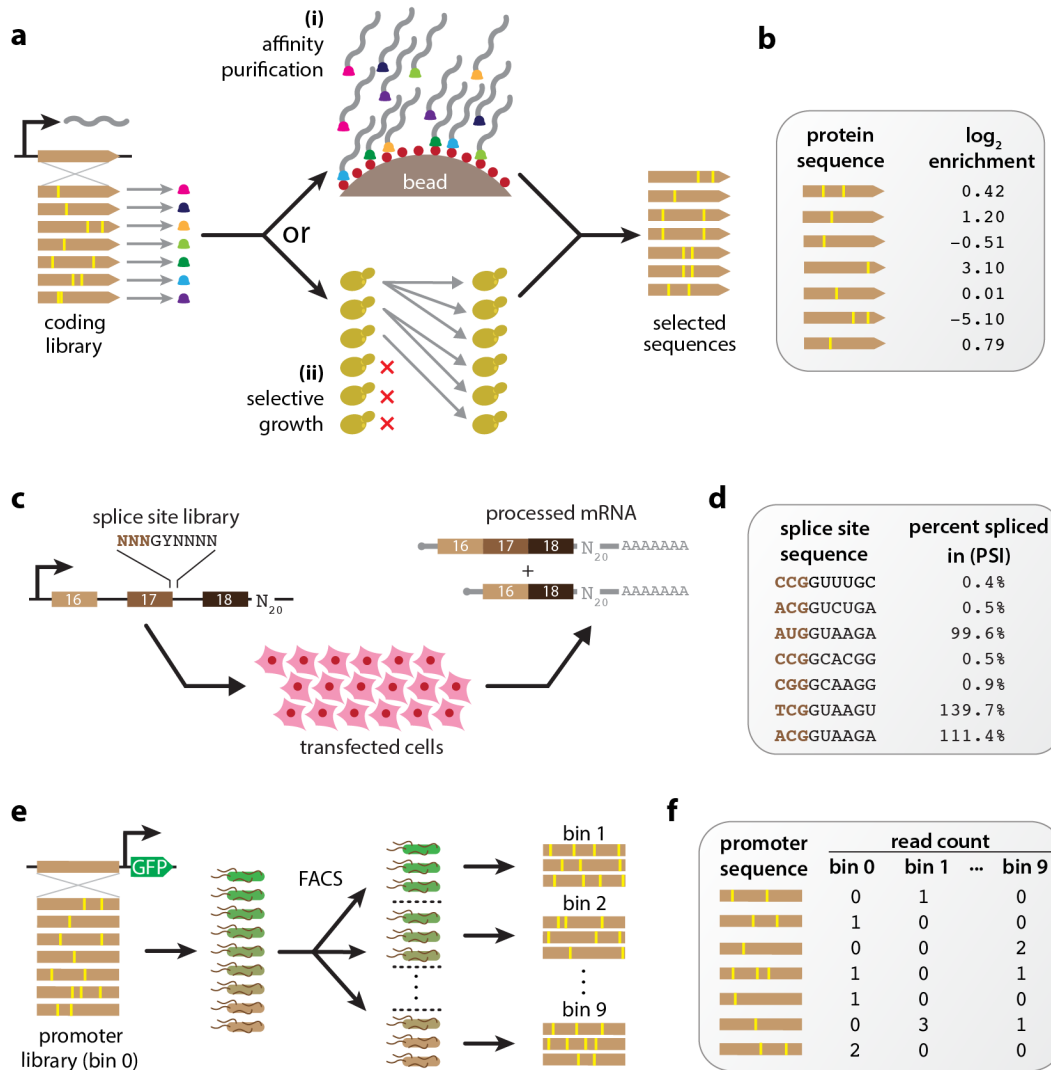
375 **Author Contributions.** AT, WTI, DMM, and JBK conceived the project. AT and JBK wrote the
376 software with assistance from AP and MK. WTI and JBK wrote a preliminary version of the
377 software. AT, MK, and JBK performed the data analysis. AT, DMM, and JBK wrote the
378 manuscript with contributions from MK and AP.

379 **Availability of Data and Materials.** MAVE-NN can be installed from PyPI by executing “pip
380 install mavenn” at the POSIX command line. Comprehensive documentation, including step-by-

381 step tutorials, is available at <http://mavenn.readthedocs.io>. Source code, the data sets analyzed
382 in this paper, and the scripts used for training the models and making the figures presented
383 herein, are available under an MIT open-source license at <https://github.com/jbkinney/mavenn>.
384 MAVE-NN version 1.0.0 was used for all of the analysis described in this manuscript and is
385 archived on Zenodo at <https://doi.org/10.5281/zenodo.5812439>.

386 **Conflicts of Interest.** The authors declare that they have no known conflicts of interest.

387

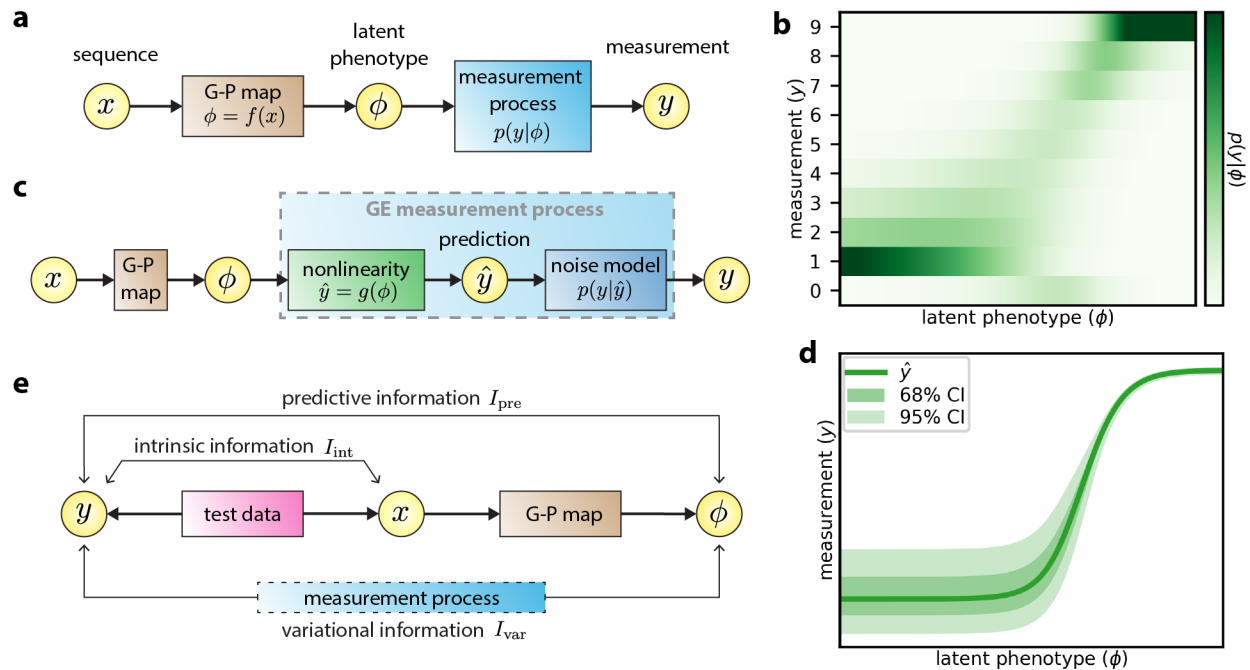


388

389 **Figure 1.** Examples illustrating the diversity of MAVEs. (a) DMS assays using either affinity purification or selective
 390 growth. (i) The DMS assay of Olson et al.³⁵ used a library of variant GB1 proteins covalently linked to their coding
 391 mRNAs via mRNA display. Functional GB1 proteins were then enriched using IgG beads. (ii) The DMS studies of
 392 Seuma et al.³⁶ and Bolognesi et al.³⁷ used selective growth in genetically modified *Saccharomyces cerevisiae* cells to
 393 respectively assay the functionality of variant A β and TDP-43 proteins. In all three experiments, deep sequencing was
 394 used to determine an enrichment ratio for each protein variant. (b) The resulting DMS dataset consists of variant
 395 protein sequences and their corresponding log enrichment values. (c) The MPSA of Wong et al.³⁸ A library of 3-exon
 396 minigenes was constructed from exons 16, 17, and 18 of *BRCA2*, with each minigene having a variant 5' ss at exon
 397 17 and a random 20 nt barcode in the 3' UTR. This library was transfected into HeLa cells, and deep sequencing was
 398 used to quantify mRNA isoform abundance. (d) The resulting MPSA dataset comprises variant 5' splice sites with (noisy) PSI

399 values. (e) The sort-seq MPRA of Kinney et al.¹⁶ A plasmid library was generated in which randomly mutagenized
400 versions of the *Escherichia coli lac* promoter drove the expression of GFP. Cells carrying these plasmids were sorted
401 using FACS, and the variant promoters in each bin of sorted cells as well as the initial library were sequenced. (f) The
402 resulting dataset comprises a list of variant promoter sequences, as well as a matrix of counts for each variant in
403 each FACS bin. MAVE: multiplex assay of variant effect; DMS: deep mutational scanning; MPSA: massively parallel
404 splicing assay; 5'ss: 5' splice site(s); PSI: percent spliced in; GFP: green fluorescent protein; FACS: fluorescence-
405 activated cell sorting.

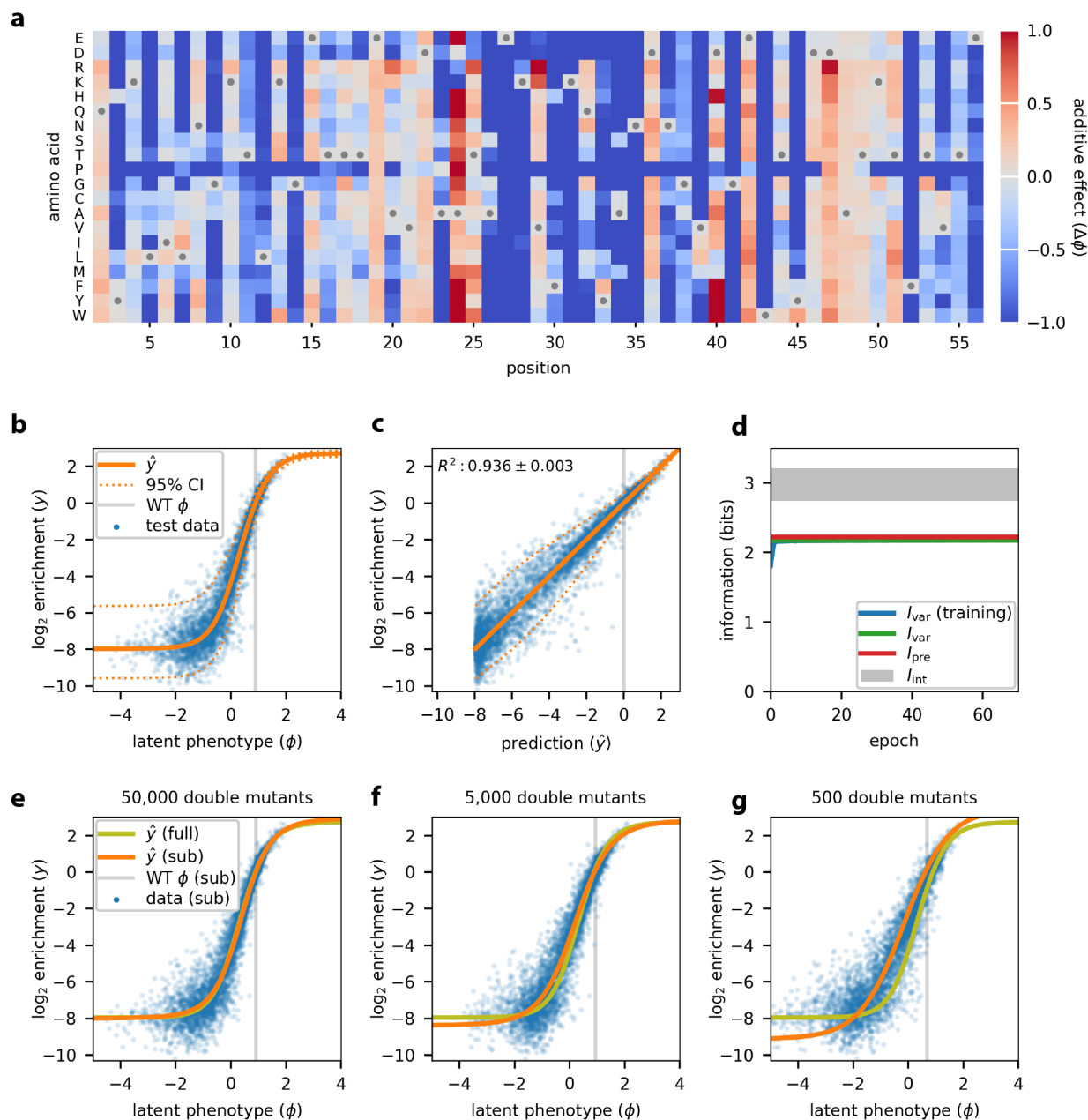
406



407

408 **Figure 2.** MAVE-NN quantitative modeling strategy. (a) Structure of latent phenotype models. A G-P map $f(x)$ maps
 409 each sequence x to a latent phenotype ϕ , after which a measurement process $p(y|\phi)$ determines the measurement
 410 y . (b) Example of an MPA measurement process inferred from the sort-seq MPRA data of Kinney et al.¹⁶ MPA
 411 measurement processes are used when y values are discrete. (c) Structure of a GE regression model, which is used
 412 when y is continuous. A GE measurement process assumes that the mode of $p(y|\phi)$, called the prediction \hat{y} , is given
 413 by a nonlinear function $g(\phi)$, and the scatter about this mode is described by a noise model $p(y|\hat{y})$. (d) Example of a
 414 GE measurement process inferred from the DMS data of Olson et al.³⁵ Shown are the nonlinearity, the 68% CI, and
 415 the 95% CI. (e) Information-theoretic quantities used to assess model performance. Intrinsic information, I_{int} , is the
 416 mutual information between sequences x and measurements y . Predictive information, I_{pre} , is the mutual information
 417 between measurements y and the latent phenotype values ϕ assigned by a model. Variational information, I_{var} , is a
 418 linear transformation of log likelihood. The inequality $I_{\text{int}} \geq I_{\text{pre}} \geq I_{\text{var}}$ always holds on test data (modulo finite data
 419 uncertainties), with $I_{\text{int}} = I_{\text{pre}}$ when the G-P map is correct, and $I_{\text{pre}} = I_{\text{var}}$ when the measurement process correctly
 420 describes the distribution of y conditioned on ϕ . G-P: genotype-phenotype; MPA: measurement process agnostic;
 421 GE: global epistasis; CI: confidence interval.

422

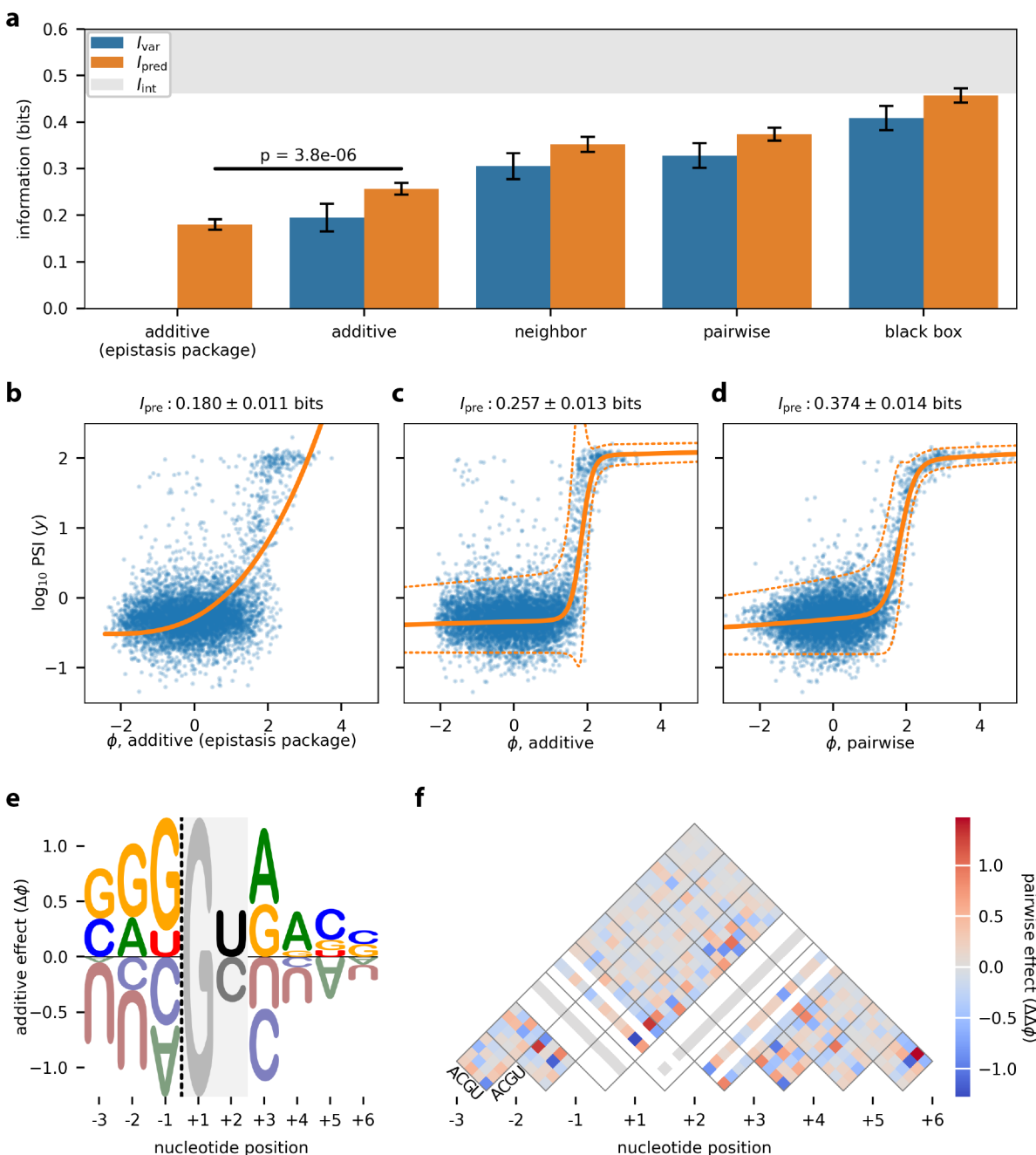


423

424 **Figure 3.** Analysis of DMS data for protein GB1. MAVE-NN was used to infer a latent phenotype model, consisting of
 425 an additive G-P map and a GE measurement process having a heteroskedastic skewed-t noise model, from the DMS
 426 data of Olson et al.³⁵ All 530,737 pairwise variants reported for positions 2 to 56 of the GB1 domain were analyzed.
 427 Data were split 90:5:5 into training, validation, and test sets. (a) The G-P map parameters inferred from all pairwise
 428 variants. Gray dots indicate wildtype residues. Amino acids are ordered as in Olson et al.³⁵ (b) GE plot showing
 429 measurements versus predicted latent phenotype values for 5,000 randomly selected test-set sequences (blue dots),

430 alongside the inferred nonlinearity (solid orange line) and the 95% CI (dashed lines) of the noise model. Gray line
431 indicates the latent phenotype value of the wildtype sequence. (c) Measurements plotted against \hat{y} predictions for
432 these same sequences. Dashed lines indicate the 95% CI of the noise model. Gray line indicates the wildtype
433 sequence \hat{y} . (d) Corresponding information metrics computed during model training (using training data) or for the
434 final model (using test data); uncertainties in these estimates are roughly the width of the plotted lines. Gray shaded
435 area indicates allowed values for intrinsic information based on upper and lower bounds estimated as described in
436 **Methods**. (e-g) Test set predictions (blue dots) and GE nonlinearities (orange lines) for models trained using subsets
437 of the GB1 data containing all single mutants and 50,000 (e), 5,000 (f), or 500 (g) double mutants. The GE
438 nonlinearity from panel b is shown for reference (yellow-green lines). Uncertainties reflect standard errors. GE: global
439 epistasis; G-P: genotype-phenotype; CI: confidence interval.

440

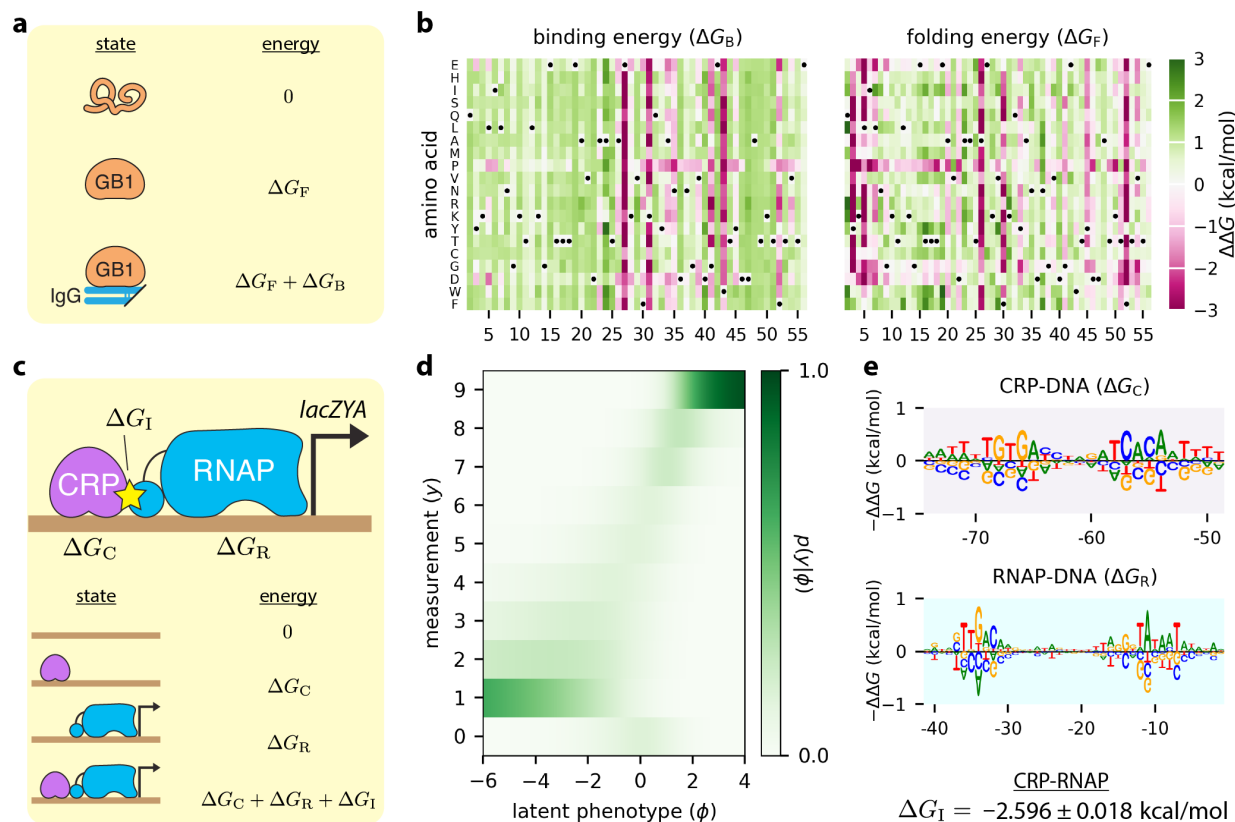


457

458 **Figure 5.** Analysis of MPSA data from Wong et al.³⁸ This dataset reports PSI values, measured in the *BRCA2* exon
 459 17 context, for nearly all 32,768 variant 5' ss of the form NNN/GYNNNN. Data were split 60:20:20 into training,
 460 validation, and test sets. Latent phenotype models with one of four types of G-P map (additive, neighbor, pairwise, or
 461 black box), as well as a GE measurement process with a heteroscedastic skewed-t noise model, were inferred. The
 462 epistasis package of Sailer and Harms²⁴ was also used to infer an additive G-P map and GE nonlinearity. **(a)**

463 Performance of trained models as quantified by I_{var} and I_{pre} , computed on test data. The lower bound on I_{int} was
464 estimated from experimental replicates (see **Methods**). p-value reflects a two-sided z-test. I_{var} was not computed for
465 the additive (epistasis package) model because that package does not infer an explicit noise model. **(b-d)**
466 Measurement values versus latent phenotype values, computed on test data, using the additive (epistasis package)
467 model **(b)**, the additive model **(c)**, and the pairwise model **(d)**. The corresponding GE measurement processes are
468 also shown. **(e)** Sequence logo⁴⁶ illustrating the additive effects component of the pairwise G-P map. Dashed line
469 indicates the exon/intron boundary. G at +1 serves as a placeholder because no other bases were assayed at this
470 position. Only values for U and C at +2 were inferred. **(f)** Heatmap showing the pairwise effects component of the
471 pairwise G-P map. White diagonals correspond to unobserved bases. Error bars indicate standard errors.
472 **Supplemental Information Fig. S3** shows the uncertainties in these parameters. MPSA: massively parallel splicing
473 assay; PSI: percent spliced in; G-P: genotype-phenotype; GE: global epistasis.

474



475

476 **Figure 6.** Biophysical models inferred from DMS and MPRA data. (a) Thermodynamic model for IgG binding by GB1.
 477 This model comprises three GB1 microstates (unfolded, folded-unbound, and folded-bound). The Gibbs free energies
 478 of folding (ΔG_F) and binding (ΔG_B) are computed from sequence using additive models called energy matrices. The
 479 latent phenotype is given by the fraction of time GB1 is in the folded-bound state. (b) The $\Delta\Delta G$ parameters of the
 480 energy matrices for folding and binding, inferred from the data of Olson et al.³⁵ using GE regression. **Supplemental**
 481 **Fig. S5** plots folding energy predictions against the measurements of Nisthal et al.⁴⁴ (c) A four-state thermodynamic
 482 model for transcriptional activation at the *E. coli lac* promoter. The Gibbs free energies of RNAP-DNA binding (ΔG_R)
 483 and CRP-DNA binding (ΔG_C) are computed using energy matrices, whereas the CRP-RNAP interaction energy ΔG_I is
 484 a scalar. The latent phenotype is the fraction of time a promoter is bound by RNAP. (d,e) The latent phenotype model
 485 inferred from the sort-seq MPRA of Kinney et al.,¹⁶ including both the MPA measurement process (d) and the
 486 parameters of the thermodynamic G-P map (e). Supplemental **Fig. S4** provides detailed definitions of the
 487 thermodynamic models in panels a and c. Sequence logos in panel e were generated using Logomaker,⁴⁶ and
 488 standard errors for the protein-protein interaction energy were determined using the built-in parametric bootstrapping

489 approach described in Methods. GE: global epistasis. RNAP: RNA polymerase. MPA: measurement-process

490 agnostic. G-P: genotype-phenotype.

491

492 **Online Methods**

493 **Notation**

494 We represent each MAVE dataset as a set of N observations, $\{(x_n, y_n)\}_{n=0}^{N-1}$, where each
495 observation consists of a sequence x_n and a measurement y_n . Here, y_n can be either a
496 continuous real-valued number, or a nonnegative integer representing the “bin” in which the n th
497 sequence was found. Note that, in this representation the same sequence x can be observed
498 multiple times, potentially with different values for y due to experimental noise.

499 **G-P maps**

500 We assume that all sequences have the same length L , and that at each of the L
501 positions in each sequence there is one of C possible characters. MAVE-NN represents
502 sequences using a vector of one-hot encoded features of the form

$$503 \quad x_{l:c} = \begin{cases} 1 & \text{if character } c \text{ occurs at position } l \\ 0 & \text{otherwise} \end{cases}, \quad (1)$$

504 where $l = 0, 1, \dots, L - 1$ indexes positions within the sequence, and c indexes the C distinct
505 characters. MAVE-NN supports built-in alphabets for DNA, RNA and protein (with or without
506 stop codons), as well as user-defined sequence alphabets.

507 We assume that the latent phenotype is given by a linear function $\phi(x; \theta)$ that depends
508 on a set of G-P map parameters θ . As mentioned in the main text, MAVE-NN supports four
509 types of G-P map models, all of which can be inferred using either GE regression or MPA
510 regression. The additive model is given by,

$$511 \quad \phi_{\text{additive}}(x; \theta) = \theta_0 + \sum_{l=0}^{L-1} \sum_c \theta_{l:c} x_{l:c}, \quad (2)$$

512 and thus each position in x contributes independently to the latent phenotype. The neighbor
513 model is given by,

$$514 \quad \phi_{\text{neighbor}}(x; \theta) = \theta_0 + \sum_{l=0}^{L-1} \sum_c \theta_{l:c} x_{l:c} + \sum_{l=0}^{L-2} \sum_{c,c'} \theta_{l:c,l+1:c'} x_{l:c} x_{l+1:c'}, \quad (3)$$

515 and further accounts for potential epistatic interactions between neighboring positions. The
516 pairwise model is given by,

$$517 \quad \phi_{\text{pairwise}}(x; \theta) = \theta_0 + \sum_{l=0}^{L-1} \sum_c \theta_{l:c} x_{l:c} + \sum_{l=0}^{L-2} \sum_{l'=l+1}^{L-1} \sum_{c,c'} \theta_{l:c,l':c'} x_{l:c} x_{l':c'}, \quad (4)$$

518 and includes interactions between all pairs of positions. Note our convention of requiring $l' > l$ in
519 the pairwise parameters $\theta_{l:c,l':c'}$.

520 Unlike these three parametric models, the black box G-P map does not have a fixed
521 functional form. Rather, it is given by a multilayer perceptron that takes a vector of sequence
522 features (additive, neighbor, or pairwise) as input, contains multiple fully-connected hidden
523 layers with nonlinear activations, and has a single node output with a linear activation. Users are
524 able to specify the number of hidden layers, the number of nodes in each hidden layer, and the
525 activation function used by these nodes.

526 MAVE-NN further supports custom G-P maps that users can define by subclassing the G-
527 P map base class. These G-P maps can have arbitrary functional form, e.g., representing specific
528 biophysical hypotheses of sequence function. This feature of MAVE-NN is showcased in the
529 analyses of **Fig. 6**.

530 **Gauge modes and diffeomorphic modes**

531 G-P maps typically have non-identifiable degrees of freedom that must be fixed, i.e.,
532 pinned down, before the values of individual parameters can be meaningfully interpreted or
533 compared between models. These degrees of freedom come in two flavors: gauge modes and
534 diffeomorphic modes. Gauge modes are changes to θ that do not alter the values of the latent
535 phenotype ϕ . Diffeomorphic modes^{15,20} are changes to θ that do alter ϕ , but do so in ways that
536 can be undone by transformations of the measurement process $p(y|\phi)$. As shown by Kinney
537 and Atwal,^{15,20} the diffeomorphic modes of linear G-P maps like those considered here will in
538 general correspond to affine transformations of ϕ , although additional unconstrained modes can
539 occur in special situations.

540 MAVE-NN fixes both gauge modes and diffeomorphic modes of inferred models (except
541 when using custom G-P maps). The diffeomorphic modes of G-P maps are fixed by
542 transforming θ via

$$543 \quad \theta_0 \rightarrow \theta_0 - a, \quad (5)$$

544 and then

$$545 \quad \theta \rightarrow \frac{\theta}{b}, \quad (6)$$

546 where $a = \text{mean}(\{\phi_n\})$ and $b = \text{std}(\{\phi_n\})$ are the mean and standard deviation of ϕ values
547 computed on the training data. This produces a corresponding change in latent phenotype
548 values $\phi \rightarrow (\phi - a)/b$. To avoid altering likelihood values, MAVE-NN makes a corresponding
549 transformation to the measurement process $p(y|\phi)$. In GE regression this is done by adjusting
550 the GE nonlinearity via

$$551 \quad g(\phi) \rightarrow g(a + b\phi), \quad (7)$$

552 while keeping the noise model $p(y|\hat{y})$ fixed. In MPA regression MAVE-NN transforms the full
553 measurement process via

554
$$p(y|\phi) \rightarrow p(y|a + b\phi). \quad (8)$$

555 For the three parametric G-P maps, gauge modes are fixed using what we call the
 556 “hierarchical gauge.” Here, the parameters θ are adjusted so that the lower-order terms in
 557 $\phi(x; \theta)$ account for the highest possible fraction of variance in ϕ . This procedure requires a
 558 probability distribution on sequence space with respect to which these variances are computed.
 559 MAVE-NN assumes that such distributions factorize by position, and can thus be represented
 560 by a probability matrix with elements $p_{l:c}$, denoting the probability of character c at position l .
 561 MAVE-NN provides three built-in choices for this distribution: uniform, empirical, or wildtype.
 562 The corresponding values of $p_{l:c}$ are given by

563
$$p_{l:c} = \begin{cases} 1/C & \text{for uniform} \\ n_{l:c}/N & \text{for empirical} \\ x_{l:c}^{\text{wt}} & \text{for wildtype} \end{cases}, \quad (9)$$

564 where $n_{l:c}$ denotes the number of sequences (out of N total) that have character c at position l ,
 565 and $x_{l:c}^{\text{wt}}$ is the one-hot encoding of a user-specified wildtype sequence. In particular, the
 566 wildtype gauge was used for illustrating the additive G-P maps in **Fig. 3** and **Fig. 4**, while the
 567 uniform gauge was used for illustrating the pairwise G-P map in **Fig. 5** and the energy matrices
 568 in **Fig. 6**. After a sequence distribution is chosen, MAVE-NN fixes the gauge of the pairwise G-P
 569 map by transforming

570
$$\begin{aligned} \theta_0 \rightarrow & \theta_0 \\ & + \sum_l \sum_{c'} \theta_{l:c'} p_{l:c'} \\ & + \sum_l \sum_{l' > l} \sum_{c,c'} \theta_{l,c,l':c'} p_{l:c} p_{l':c'}, \end{aligned} \quad (10)$$

571

$$\begin{aligned}
 \theta_{l:c} &\rightarrow \theta_{l:c} \\
 &- \sum_{c'} \theta_{l:c'} p_{l:c'} \\
 &+ \sum_{l'>l} \sum_{c'} \theta_{l:c,l':c'} p_{l':c'} \\
 &+ \sum_{l'<l} \sum_{c'} \theta_{l':c',l:c} p_{l':c'} \\
 &- \sum_{l'>l} \sum_{c',c''} \theta_{l:c',l':c''} p_{l:c'} p_{l':c''} \\
 &- \sum_{l'<l} \sum_{c',c''} \theta_{l':c'',l':c'} p_{l:c'} p_{l':c''},
 \end{aligned} \tag{11}$$

572 and

573

$$\begin{aligned}
 \theta_{l:c,l':c'} &\rightarrow \theta_{l:c,l':c'} \\
 &- \sum_{c''} \theta_{l:c'',l':c'} p_{l:c''} \\
 &- \sum_{c''} \theta_{l:c,l':c''} p_{l':c''} \\
 &+ \sum_{c'',c'''} \theta_{l:c'',l':c'''} p_{l:c''} p_{l':c'''}.
 \end{aligned} \tag{12}$$

574 This transformation is also used for the additive and neighbor G-P maps, but with $\theta_{l:c,l':c'} = 0$ for
 575 all l, l' (additive) or whenever $l' \neq l + 1$ (neighbor).

576 **GE nonlinearities**

577 GE models assume that each measurement y is a nonlinear function of the latent
 578 phenotype $g(\phi)$ plus some noise. In MAVE-NN, this nonlinearity is represented as a sum of
 579 tanh sigmoids:

580

$$g(\phi; \alpha) = a + \sum_{k=0}^{K-1} b_k \tanh(c_k \phi + d_k). \tag{13}$$

581 Here, K specifies the number of hidden nodes contributing to the sum, and $\alpha = \{a, b_k, c_k, d_k\}$ are
582 trainable parameters. We note that this mathematical form is an example of the bottleneck
583 architecture previously used by^{23,45} for modeling GE nonlinearities. By default, MAVE-NN
584 constrains $g(\phi; \alpha)$ to be monotonic in ϕ by requiring all $b_k \geq 0$ and $c_k \geq 0$, but this constraint
585 can be relaxed.

586 **GE noise models**

587 MAVE-NN supports three types of GE noise model: Gaussian, Cauchy, and skew-t.
588 These all support the analytic computation of quantiles and confidence intervals, as well as the
589 rapid sampling of simulated measurement values. The Gaussian noise model is given by

$$590 \quad p_{\text{gauss}}(y|\hat{y}; s) = \frac{1}{\sqrt{2\pi s^2}} \exp\left[-\frac{(y - \hat{y})^2}{2s^2}\right], \quad (14)$$

591 where s denotes the standard deviation. Importantly, MAVE-NN allows this noise model to be
592 heteroskedastic by representing s as an exponentiated polynomial in \hat{y} , i.e.,

$$593 \quad s(\hat{y}) = \exp\left[\sum_{k=0}^K a_k \hat{y}^k\right], \quad (15)$$

594 where K is the order of the polynomial and $\{a_k\}$ are trainable parameters. The user has the
595 option to set K , and setting $K = 0$ renders this noise model homoscedastic. Quantiles are
596 computed using $y_q = \hat{y} + s\sqrt{2} \operatorname{erf}^{-1}(2q - 1)$ for user-specified values of $q \in [0,1]$. Similarly, the
597 Cauchy noise model is given by

$$598 \quad p_{\text{cauchy}}(y|\hat{y}; s) = \left[\pi s \left(1 + \frac{(y - \hat{y})^2}{s^2}\right)\right]^{-1}, \quad (16)$$

599 where the scale parameter s is an exponentiated K 'th order polynomial in \hat{y} , and quantiles are
600 computed using $y_q = \hat{y} + s \tan\left[\pi\left(q - \frac{1}{2}\right)\right]$.

601 The skew- t noise model is of the form described by Jones and Faddy,²⁹ and is given by

602
$$p_{\text{skewt}}(y|\hat{y}; s, a, b) = s^{-1}f(t; a, b), \quad (17)$$

603 where

604
$$t = t^* + \frac{y - \hat{y}}{s}, \quad t^* = \frac{(a - b)\sqrt{a + b}}{\sqrt{2a + 1}\sqrt{2b + 1}}, \quad (18)$$

605 and

606
$$f(t; a, b) = \frac{2^{1-a-b} \Gamma(a + b)}{\sqrt{a + b} \Gamma(a)\Gamma(b)} \left[1 + \frac{t}{\sqrt{a + b + t^2}} \right]^{a+\frac{1}{2}} \times \left[1 - \frac{t}{\sqrt{a + b + t^2}} \right]^{b+\frac{1}{2}}. \quad (19)$$

607 Note that the t statistic here is an affine function of y chosen so that the distribution's mode
 608 (corresponding to t^*) is positioned at \hat{y} . The three parameters of this noise model, $\{s, a, b\}$, are
 609 each represented using K -th order exponentiated polynomials with trainable coefficients.

610 Quantiles are computed using

611
$$y_q = \hat{y} + (t_q - t^*)s, \quad (20)$$

612 where

613
$$t_q = \frac{(2x_q - 1)\sqrt{a + b}}{\sqrt{1 - (2x_q - 1)^2}}, \quad x_q = I_q^{-1}(a, b), \quad (21)$$

614 and I^{-1} denotes the inverse of the regularized incomplete Beta function $I_x(a, b)$.

615 Empirical noise models

616 MAVE-NN further supports the inference of GE regression models that account for user-
 617 specified measurement noise. In such cases, the user provides a set of measurement-specific

618 uncertainties $\{s_n\}_{n=0}^{N-1}$ along with the corresponding sequences and measurements. These uncertainties
 619 can, for example, be estimated by using a software package like Enrich2¹¹ or DiMSum¹⁴. MAVE-NN then
 620 trains the parameters of latent phenotype models by assuming a Gaussian noise model of the form

$$621 \quad p_{\text{empirical}}(y_n | \hat{y}_n, s_n) = \frac{1}{\sqrt{2\pi s_n^2}} \exp\left[-\frac{(y_n - \hat{y}_n)^2}{2s_n^2}\right], \quad (14)$$

622 where $\hat{y}_n = g(f(x_n; \theta); \alpha)$ is the expected measurement for sequence x_n , θ denotes G-P map
 623 parameters, and α denotes the parameters of the GE nonlinearity. This noise model thus has
 624 the advantage of having no free parameters, but it may be problematically mis-specified if the
 625 true error distribution is heavy-tailed or skewed.

626 **MPA measurement process**

627 In MPA regression, MAVE-NN directly models the measurement process $p(y|\phi)$. At
 628 present, MAVE-NN only supports MPA regression for discrete values of y indexed using
 629 nonnegative integers. MAVE-NN supports two alternative forms of input for MPA regression.
 630 One is a set of sequence-measurement pairs $\{(x_n, y_n)\}_{n=0}^{N-1}$, where N is the total number of
 631 reads, $\{x_n\}$ is a set of (typically) non-unique sequences, each $y_n \in \{0, 1, \dots, Y-1\}$ is a bin
 632 number, and Y is the total number of bins. The other is a set of sequence-count-vector pairs
 633 $\{(x_m, c_m)\}_{m=0}^{M-1}$, where M is the total number of unique sequences and $c_m = (c_{m0}, c_{m1}, \dots, c_{m(Y-1)})$
 634 is a vector that lists the number of times c_{my} that the sequence x_m was observed in each bin y .
 635 MPA measurement processes are represented as multilayer perceptron with one hidden layer
 636 (having tanh activations) and a softmax output layer. Specifically,

$$637 \quad p(y|\phi) = \frac{w_y(\phi)}{\sum_{y'} w_{y'}(\phi)}, \quad (22)$$

638 where

$$639 \quad w_y(\phi) = \exp\left[a_y + \sum_{k=0}^{K-1} b_{yk} \tanh(c_{yk}\phi + d_{yk})\right] \quad (23)$$

640 and K is the number of hidden nodes per value of y . The trainable parameters of this
 641 measurement process are $\eta = \{a_y, b_{yk}, c_{yk}, d_{yk}\}$.

642 Loss function

643 Let θ denote the G-P map parameters, and η denote the parameters of the
 644 measurement process. MAVE-NN optimizes these parameters using stochastic gradient
 645 descent on a loss function given by

$$646 \quad \mathcal{L} = \mathcal{L}_{\text{like}} + \mathcal{L}_{\text{reg}}, \quad (24)$$

647 where $\mathcal{L}_{\text{like}}$ is the negative log likelihood of the model, given by

$$648 \quad \mathcal{L}_{\text{like}}[\theta, \eta] = - \sum_{n=0}^{N-1} \log [p(y_n | \phi_n; \eta)] \quad (25)$$

649 where $\phi_n = \phi(x_n; \theta)$, and \mathcal{L}_{reg} provides for regularization of the model parameters.

650 In the context of GE regression, we can write $\eta = (\alpha, \beta)$ where α represents the
 651 parameters of the GE nonlinearity $g(\phi; \alpha)$, and β denotes the parameters of the noise model
 652 $p(y | \hat{y}; \beta)$. The likelihood contribution from each observation n then becomes $p(y_n | \phi_n; \eta) =$
 653 $p(y_n | \hat{y}_n; \beta)$ where $\hat{y}_n = g(\phi_n; \alpha)$. In the context of MPA regression with a dataset of the form
 654 $\{(x_m, c_m)\}_{m=0}^{M-1}$, the loss function simplifies to

$$655 \quad \mathcal{L}_{\text{like}}[\theta, \eta] = - \sum_{m=0}^{M-1} \sum_{y=0}^{Y-1} c_{my} \log [p(y | \phi_m; \eta)] \quad (26)$$

656 where $\phi_m = \phi(x_m; \theta)$. For the regularization term, MAVE-NN uses an L_2 penalty of the form

$$657 \quad \mathcal{L}_{\text{reg}}[\theta, \eta] = \lambda_\theta |\theta|^2 + \lambda_\eta |\eta|^2, \quad (27)$$

658 where the user-adjusted parameters λ_θ and λ_η respectively control the strength of regularization
659 for the G-P map and measurement process parameters.

660 **Predictive information**

661 In what follows, we use $p_{\text{model}}(y|\phi)$ to denote a measurement process inferred by
662 MAVE-NN, whereas $p_{\text{true}}(y|\phi)$ denotes the empirical conditional distribution of y and ϕ values
663 that would be observed in the limit of infinite test data.

664 Predictive information $I_{\text{pre}} = I[y; \phi]$, where $I[\cdot; \cdot]$ represents mutual information computed
665 on data not used for training (i.e., a held-out test set or data from a different experiment), I_{pre}
666 provides a measure of how strongly a G-P map predicts experimental measurements.
667 Importantly, this quantity does not depend on the corresponding measurement process
668 $p_{\text{model}}(y|\phi)$. To estimate I_{pre} , we use k 'th nearest neighbor (kNN) estimators of entropy and
669 mutual information adapted from the NPEET Python package.⁴⁷ Here, the user has the option of
670 adjusting k , which controls a variance/bias tradeoff. When y is discrete (MPA regression), I_{pre} is
671 computed using the classic kNN entropy estimator^{48,49} via the decomposition $I[y; \phi] = H[\phi] -$
672 $\sum_y p(y) H_y[\phi]$, where $H_y[\phi]$ denotes the entropy of $p_{\text{true}}(\phi|y)$. When y is continuous (GE
673 regression), $I[y; \phi]$ is estimated using the kNN-based Kraskov Stögbauer Grassberger (KSG)
674 algorithm.⁴⁹ This approach optionally supports the local nonuniformity correction of Gao et al.,
675 which is important when y and ϕ exhibit strong dependencies, but which also requires
676 substantially more time to compute.

677 **Variational information**

678 We define variational information as an affine transformation of $\mathcal{L}_{\text{like}}$,

$$679 \quad I_{\text{var}} = H[y] - \frac{\log_2(e)}{N} \mathcal{L}_{\text{like}}. \quad (28)$$

680 Here, $H[y]$ is the entropy of the data $\{y_n\}$, which is estimated using the k 'th nearest neighbor
681 (kNN) estimator from the NPEET package.⁴⁷ Noting that this quantity can also be written as
682 $I_{\text{var}} = H[y] - \text{mean}(\{Q_n\})$, where $Q_n = -\log_2 p(y_n|\phi_n)$, we estimate the associated uncertainty
683 using

$$684 \quad \delta I_{\text{var}}[y; \phi] = \sqrt{\delta H[y]^2 + \frac{\text{var}(\{Q_n\})}{N}}. \quad (29)$$

685 The inference strategy used by MAVE-NN is based on the fact that I_{var} provides a tight
686 variational lower bound on I_{pre} . Indeed, in the large data limit,

$$687 \quad I_{\text{pre}} = I_{\text{var}} + D_{\text{KL}}(p_{\text{true}}||p_{\text{model}}), \quad (30)$$

688 where $D_{\text{KL}}(\cdot) \geq 0$ is the Kullback-Leibler divergence, and thus quantifies the accuracy of the
689 inferred measurement process. From **Eq. 30** one can see that, with appropriate caveats,
690 maximizing I_{var} (or equivalently, $\mathcal{L}_{\text{like}}$) will also maximize I_{pre} .²⁰ But unlike I_{pre} , I_{var} is readily
691 compatible with backpropagation and stochastic gradient descent. See Supplemental
692 Information for a derivation of **Eq. 30** and an expanded discussion of this key point. Note:
693 Sharpee et al.⁵⁰ cleverly showed that I_{pre} can, in fact, be optimized using stochastic gradient
694 descent. Computing gradients of I_{pre} , however, requires a time-consuming density estimation
695 step. Optimizing I_{var} , on the other hand, can be done using standard per-datum
696 backpropagation.

697 **Intrinsic information**

698 Intrinsic information, $I_{\text{int}} = I[x; y]$, is the mutual information between the sequences x
699 and measurements y in a dataset. This quantity is somewhat tricky to estimate due to the high-
700 dimensional nature of sequence space. We instead used three different methods to obtain the
701 upper and lower bounds on I_{int} shown in **Fig. 3d** and **Fig. 5a**. More generally, we believe the

702 development of both computational and experimental methods for estimating I_{int} is be an
703 important avenue for future research.

704 To compute the upper bound on I_{int} for GB1 data (in **Fig. 3d**), we used the fact that

$$705 \quad I[x; y] = H[y] - \langle H_x[y] \rangle_x, \quad (31)$$

706 where $H[y]$ is the entropy of all measurements y , $H_x[y]$ is the entropy of $p(y|x)$ for a specific
707 choice of sequence x , and $\langle \cdot \rangle_x$ indicates averaging over all sequences x . In this dataset, the
708 measurement values were computed using

$$709 \quad y = \log_2 \left[\frac{c_s + 1}{c_i + 1} \right], \quad (32)$$

710 where c_i is the input read count and c_s is the selected read count. $H[y]$ was estimated using the
711 KNN estimator.⁴⁸ We estimated the uncertainty in y by propagating errors expected due to
712 Poisson fluctuations in read counts, which gives

$$713 \quad \delta y = \log_2(e) \sqrt{\frac{1}{c_s+1} + \frac{1}{c_i+1}}. \quad (33)$$

714 Then, assuming $p(y|x)$ to be approximately Gaussian, we find the corresponding conditional
715 entropy to be

$$716 \quad H_x[y] = \frac{1}{2} \log_2(2\pi e \delta y^2). \quad (34)$$

717 These $H[y]$ and $H_x[y]$ values were then used in **Eq. 31** to estimate I_{int} . This should provide an
718 upper bound on the true value of I_{int} because uncertainty in y must be at least that expected
719 under Poisson sampling of reads. We note, however, that the use of linear error propagation
720 and the assumption that $p(y|x)$ is approximately Gaussian complicate this conclusion. Also,
721 when applied to MPSA data, this method yielded an upper bound of 0.96 bits. We believe this
722 value is likely to be far higher than the true value of I_{int} , and that this mismatch probably
723 resulted from read counts in the MPSA data being over-dispersed.

724 To compute the lower bound on I_{int} for GB1 data (**Fig. 3d**) we used the predictive
725 information I_{pre} (on test data) of a GE regression model having a blackbox G-P map. This
726 provides a lower bound because $I_{\text{int}} \geq I_{\text{pre}}$ for any model (when evaluated on test data) due to
727 the Data Processing Inequality and the Markov Chain nature of the dependencies $y \leftarrow x \rightarrow \phi$ in
728 **Fig. 2e**.^{20,31}

729 To compute a lower bound on I_{int} for MPSA data (**Fig. 5c**), we leveraged the availability
730 of replicate data in Wong et al..³⁸ Let y and y' represent the original and replicate
731 measurements obtained for a sequence x . Because $y \leftarrow x \rightarrow y'$ forms a Markov chain, $I[x; y] \geq$
732 $I[y; y']$.³¹ We therefore used an estimate of $I[y; y']$, computed using the KSG method,^{47,49} as the
733 lower bound for I_{int} .

734 **Uncertainties in kNN estimates of mutual information**

735 MAVE-NN quantifies uncertainties in $H[y]$ and $I[y; \phi]$ using multiple random samples of
736 half the data. Let $\mathcal{D}_{100\%}$ denote a full dataset, and let $\mathcal{D}_{50\%,r}$ denote a 50% subsample (indexed
737 by r) of this dataset. Given an estimator $E(\cdot)$ of either entropy or mutual information, as well as
738 the number of subsamples R to use, the uncertainty in $E(\mathcal{D}_{100\%})$ is estimated as

$$739 \quad \delta E(\mathcal{D}_{100\%}) = \frac{1}{\sqrt{2}} \text{std} \left[\left\{ E(\mathcal{D}_{50\%,r}) \right\}_{r=0}^{R-1} \right]. \quad (35)$$

740 MAVE-NN uses $R = 25$ by default. We note that computing such uncertainty estimates
741 substantially increases computation time, as $E(\cdot)$ needs to be evaluated $R + 1$ times instead of
742 just once. We also note that bootstrap resampling^{51,52} is often inadvisable in this context, as it
743 systematically underestimates $H[y]$ and overestimates $I[y; z]$.

744 **Uncertainties in G-P map parameters**

745 Given a trained latent phenotype model, having G-P map parameters θ^* and
746 measurement process parameters η^* , MAVE-NN can optionally assess model uncertainty using
747 the following parametric bootstrap approach. Using the trained model with parameters (θ^*, η^*)
748 as “ground truth”, MAVE-NN simulates R (chosen by the user) different MAVE datasets $\mathcal{D}_r =$
749 $\{(x_n, y_n^{(r)})\}_{n=0}^{N-1}$, where $r = 0, 1, \dots, R - 1$ indexes the different simulations. Note that the
750 sequences in these simulated datasets are the same as the original training sequence, but the
751 measurements are different. For each simulated dataset \mathcal{D}_r , MAVE-NN then trains a new
752 model, by default using the same hyperparameters as were used for the ground truth model.
753 This procedure yields a set $\{(\tilde{\theta}^{(r)}, \eta^{(r)})\}_{r=0}^{R-1}$ of simulation-inferred G-P map parameters and
754 corresponding measurement process parameters. Users can then use this sampling of G-P map
755 parameters to estimate uncertainties, e.g., by reporting $\delta\theta_k = \text{std} \left[\left\{ \theta_k^{(r)} \right\}_{r=0}^{R-1} \right]$.

756 An important detail when assessing parameter uncertainty is to ensure that both the
757 gauge modes and the diffeomorphic modes of each model are fixed. This is necessary so that
758 differences in the parameters that cannot result in changes to model predictions do not inflate
759 the uncertainty estimates. For additive, neighbor, and pairwise G-P maps, MAVE-NN
760 automatically implements the procedure described in the “Gauge modes and diffeomorphic
761 modes” section above, thereby removing these extra degrees of freedom. However, for more
762 complex models such as those implemented by MAVE-NN’s custom G-P map functionality (e.g.,
763 representing biophysical models) different gauge freedoms and diffeomorphic modes may arise
764 depending on the details of the model, and users must take care to determine and fix these
765 prior to assessing parameter uncertainty. We also note that no meaningful computation of

766 individual parameter uncertainties is likely to be possible for highly over-parameterized models,
767 such as the “black box” deep neural network models supported by MAVE-NN.

768 Datasets

769 For the GB1 DMS dataset of Olson et al.,³⁵ measurements were computed using

$$770 \quad y_n = \log_2 \frac{(c_n^{\text{out}}+1)/(c_{\text{WT}}^{\text{out}}+1)}{(c_n^{\text{in}}+1)/(c_{\text{WT}}^{\text{in}}+1)},$$

771 where c_n^{in} and c_n^{out} respectively represent the number of reads from the input and output
772 samples (i.e., pre-selection and post-selection libraries), and $n = \text{WT}$ represents the 55 aa
773 wildtype sequence, corresponding to positions 2-56 of the GB1 domain. To infer the model in
774 **Fig. 3b** and to compute the information metrics in **Fig. 3c**, only double-mutant sequences with
775 $c_n^{\text{in}} \geq 10$ were used; these represent 530,737 out of the 536,085 possible double mutants. For
776 the models in **Figs. 3d-f**, y_n values for the 1045 single-mutant were also used in the inference
777 procedure.

778 For the A β DMS data of Seuma et al.³⁶ and TDP-43 DMS data of Bolognesi et al.,³⁷ y_n
779 values respectively represent nucleation scores and toxicity scores reported by the authors.

780 For the MPSA data of Wong et al.,³⁸ we used the data of library 1 replicate 1 obtained
781 for the *BRCA2* minigene data. Measurements were computed as

$$782 \quad y_n = \log_{10} \left[100 \times \frac{(c_n^{\text{inc}}+1)/(c_{\text{CONS}}^{\text{inc}}+1)}{(c_n^{\text{tot}}+1)/(c_{\text{CONS}}^{\text{tot}}+1)} \right],$$

783 where c_n^{inc} and c_n^{tot} respectively represent the number of barcode reads obtained from exon
784 inclusion isoforms and from total mRNA, and $n = \text{CONS}$ corresponds to the consensus 5' ss

785 sequence CAG/GUAAGU. Corresponding PSI values were computed as $PSI_n = 10^{y_n}$. Only
786 sequences with $c_n^{\text{tot}} \geq 10$ were used, representing 30,483 of the 32,768 possible sequences of
787 the form NNN/GYNNNN.

788 For the *lac* promoter sort-seq MPRA data of Kinney et al.,¹⁶ we used data from the “full-
789 wt” experiment (available at https://github.com/jbkinney/09_sortseq).

790

791 **References**

792

793 1. Kinney, J. B. & McCandlish, D. M. Massively parallel assays and quantitative sequence-
794 function relationships. *Annu Rev Genom Hum G* **20**, 99-127 (2019).

795 2. Starita, L. M. *et al.* Variant interpretation: functional assays to the rescue. *Am J Hum Genetics*
796 **101**, 315-325 (2017).

797 3. Fowler, D. M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nat*
798 *Methods* **11**, 801-807 (2014).

799 4. Levo, M. & Segal, E. In pursuit of design principles of regulatory sequences. *Nat Rev Genet*
800 **15**, 453-468 (2014).

801 5. White, M. A. Understanding how cis-regulatory function is encoded in DNA sequence using
802 massively parallel reporter assays and designed sequences. *Genomics* **106**, 165-170 (2015).

803 6. Inoue, F. & Ahituv, N. Decoding enhancers using massively parallel reporter assays.
804 *Genomics* **106**, 159-164 (2015).

805 7. Peterman, N. & Levine, E. Sort-seq under the hood: implications of design choices on large-
806 scale characterization of sequence-function relations. *BMC Genomics* **17**, 206 (2016).

807 8. Fowler, D. M., Araya, C. L., Gerard, W. & Fields, S. Enrich: software for analysis of protein
808 function by enrichment and depletion of variants. *Bioinformatics* **27**, 3430-3431 (2011).

809 9. Alam, K. K., Chang, J. L. & Burke, D. H. FASTAptamer: A bioinformatic toolkit for high-
810 throughput sequence analysis of combinatorial selections. *Mol Ther-Nucleic Acids* **4**, e230
811 (2015).

812 10. Bloom, J. D. Software for the analysis and visualization of deep mutational scanning data.
813 *BMC Bioinformatics* **16**, 168 (2015).

814 11. Rubin, A. F. *et al.* A statistical framework for analyzing deep mutational scanning data.
815 *Genome Biol* **18**, 1-15 (2017).

816 12. Ashuach, T. *et al.* MPRAnalyze: statistical framework for massively parallel reporter assays.
817 *Genome Biol* **20**, 183 (2019).

818 13. Niroula, A., Ajore, R. & Nilsson, B. MPRAScore: robust and non-parametric analysis of
819 massively parallel reporter assays. *Bioinformatics* **35**, 5351-5353 (2019).

820 14. Faure, A. J., Schmiedel, J. M., Baeza-Centurion, P. & Lehner, B. DiMSum: an error model
821 and pipeline for analyzing deep mutational scanning data and diagnosing common experimental
822 pathologies. *Genome Biol* **21**, 207 (2020).

- 823 15. Atwal, G. S. & Kinney, J. B. Learning quantitative sequence-function relationships from
824 massively parallel experiments. *J Stat Phys* **162**, 1203-1243 (2016).
- 825 16. Kinney, J. B., Murugan, A., Callan, C. G. & Cox, E. C. Using deep sequencing to
826 characterize the biophysical mechanism of a transcriptional regulatory sequence. *Proc Natl*
827 *Acad Sci USA* **107**, 9158-9163 (2010).
- 828 17. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human
829 cells using a massively parallel reporter assay. *Nat Biotechnol* **30**, 271-277 (2012).
- 830 18. Mogno, I., Kwasnieski, J. C. & Cohen, B. A. Massively parallel synthetic promoter assays
831 reveal the in vivo effects of binding site variants. *Genome Res* **23**, 1908-1915 (2013).
- 832 19. Abadi, M. *et al.* TensorFlow: A Systems for Large-Scale Machine Learning. in *Proceedings*
833 *of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI '16)*
834 (2016).
- 835 20. Kinney, J. B. & Atwal, G. S. Parametric inference in the large data limit using maximally
836 informative models. *Neural Comput* **26**, 637-653 (2014).
- 837 21. Kinney, J. B., Tkačik, G. & Callan, C. G. Precise physical models of protein-DNA interaction
838 from high-throughput data. *Proc Natl Acad Sci USA* **104**, 501-506 (2007).
- 839 22. Otwinowski, J. & Nemenman, I. Genotype to phenotype mapping and the fitness landscape
840 of the E. coli lac promoter. *PLoS ONE* **8**, e61570 (2013).
- 841 23. Sarkisyan, K. S. *et al.* Local fitness landscape of the green fluorescent protein. *Nature* **533**,
842 397-401 (2016).
- 843 24. Sailer, Z. R. & Harms, M. J. Detecting high-order epistasis in nonlinear genotype-phenotype
844 maps. *Genetics* **205**, 1079-1088 (2017).
- 845 25. Otwinowski, J., McCandlish, D. M. & Plotkin, J. B. Inferring the shape of global epistasis.
846 *Proc Natl Acad Sci USA* **115**, E7550-E7558 (2018).
- 847 26. Gelman, S., Fahlberg, S. A., Heinzelman, P., Romero, P. A. & Gitter, A. Neural networks to
848 learn protein sequence-function relationships from deep mutational scanning data. *Proc Natl*
849 *Acad Sci USA* **118**, e2104878118 (2021).
- 850 27. Faure, A. J. *et al.* Global mapping of the energetic and allosteric landscapes of protein
851 binding domains. *bioRxiv* doi:10.1101/2021.09.14.460249 (2021)
- 852 28. Tonner, P. D., Pressman, A. & Ross, D. Interpretable modeling of genotype-phenotype
853 landscapes with state-of-the-art predictive power. *bioRxiv* doi:10.1101/2021.06.11.448129
854 (2021).
- 855 29. Jones, M. C. & Faddy, M. J. A skew extension of the t-distribution, with applications. *J Roy*
856 *Stat Soc B* **65**, 159-174 (2003).

- 857 30. Kinney, J. B. & Atwal, G. S. Equitability, mutual information, and the maximal information
858 coefficient. *Proc Natl Acad Sci USA* **111**, 3354-3359 (2014).
- 859 31. Cover, T. M. & Thomas, J. A. *Elements of information theory*. (Wiley, 2006).
- 860 32. Barber, D. & Agakov, F. The IM algorithm: a variational approach to information
861 maximization. *Advances in neural information processing systems* **16**. (2004).
- 862 33. Alemi, A. A., Fischer, I., Dillon, J. V. & Murphy, K. Deep variational information bottleneck.
863 *arXiv:1612.00410* (2016).
- 864 34. Chalk, M., Marre, O. & Tkačik, G. Relevant sparse codes with variational information
865 bottleneck. *arXiv:1605.07332* (2016).
- 866 35. Olson, C. A., Wu, N. C. & Sun, R. A comprehensive biophysical description of pairwise
867 epistasis throughout an entire protein domain. *Curr Biol* **24**, 2643-2651 (2014).
- 868 36. Seuma, M., Faure, A., Badia, M., Lehner, B. & Bolognesi, B. The genetic landscape for
869 amyloid beta fibril nucleation accurately discriminates familial Alzheimer's disease mutations.
870 *eLife* **10**, e63364 (2021).
- 871 37. Bolognesi, B. *et al.* The mutational landscape of a prion-like domain. *Nat Commun* **10**, 4162
872 (2019).
- 873 38. Wong, M. S., Kinney, J. B. & Krainer, A. R. Quantitative activity profile and context
874 dependence of all human 5' splice sites. *Mol Cell* **71**, 1012-1026.e3 (2018).
- 875 39. Bintu, L. *et al.* Transcriptional regulation by the numbers: models. *Curr Opin Genet Dev* **15**,
876 116-124 (2005).
- 877 40. Sherman, M. S. & Cohen, B. A. Thermodynamic state ensemble models of cis-regulation.
878 *PLoS Comput Biol* **8**, e1002407 (2012).
- 879 41. Wong, F. & Gunawardena, J. Gene Regulation in and out of equilibrium. *Annu Rev Biophys*
880 **49**, 199-226 (2020).
- 881 42. Otwinowski, J. Biophysical inference of epistasis and the effects of mutations on protein
882 stability and function. *Mol Biol Evol* **35**, 2345-2354 (2018).
- 883 43. Manhart, M. & Morozov, A. V. Protein folding and binding can emerge as evolutionary
884 spandrels through structural coupling. *Proc Natl Acad Sci USA* **112**, 1797-1802 (2015).
- 885 44. Nisthal, A., Wang, C. Y., Ary, M. L. & Mayo, S. L. Protein stability engineering insights
886 revealed by domain-wide comprehensive mutagenesis. *Proc Natl Acad Sci USA* **116**, 16367-
887 16377 (2019).
- 888 45. Pokusaeva, V. O. *et al.* An experimental assay of the interactions of amino acids from
889 orthologous sequences shaping a complex fitness landscape. *PLoS Genet* **15**, e1008079
890 (2019).

- 891 46. Tareen, A. & Kinney, J. B. Logomaker: beautiful sequence logos in Python. *Bioinformatics*
892 **36**, 2272-2274 (2020).
- 893 47. Steeg, G. V. Non-Parametric Entropy Estimation Toolbox (NPEET).
894 <https://www.isi.edu/~gregv/npeet.html> (2014).
- 895 48. Vasicek, O. A test for normality based on sample entropy. *J Roy Stat Soc B* **38**, 54-59
896 (1976).
- 897 49. Kraskov, A., Stögbauer, H. & Grassberger, P. Estimating mutual information. *Phys Rev E*
898 **69**, 066138 (2004).
- 899 50. Sharpee, T., Rust, N. C. & Bialek, W. Analyzing neural responses to natural signals:
900 maximally informative dimensions. *Neural Comput* **16**, 223-250 (2004).
- 901 51. Efron, B. Bootstrap methods: another look at the jackknife. *Ann Stat* **7**, 1-26 (1979)..
- 902 52. Efron, B. & Tibshirani, R. Bootstrap methods for standard errors, confidence intervals, and
903 other measures of statistical accuracy. *Stat Sci* **1**, 54-75 (1986).
- 904