

# 1 Determining protein structures using genetics

2 Jörn M. Schmiedel<sup>1</sup>, Ben Lehner<sup>1-3\*</sup>

3 <sup>1</sup> Systems Biology Program, Centre for Genomic Regulation (CRG), The Barcelona Institute of  
4 Science and Technology, Dr. Aiguader 88, Barcelona 08003, Spain

5 <sup>2</sup> Universitat Pompeu Fabra (UPF), Barcelona 08003, Spain

6 <sup>3</sup> Institució Catalana de Recerca i Estudis Avançats (ICREA), Pg. Lluís Companys 23, 08010  
7 Barcelona, Spain

8 \*email: [ben.lehner@crg.eu](mailto:ben.lehner@crg.eu)

9

## 10 Summary

11 Determining the three dimensional structures of macromolecules is a major goal of biological  
12 research because of the close relationship between structure and function. Structure  
13 determination usually relies on physical techniques including x-ray crystallography, NMR  
14 spectroscopy and cryo-electron microscopy. Here we present a method that allows the high-  
15 resolution three-dimensional structure of a biological macromolecule to be determined only from  
16 measurements of the activity of mutant variants of the molecule. This genetic approach to  
17 structure determination relies on the quantification of genetic interactions (epistasis) between  
18 mutations and the discrimination of direct from indirect interactions. This provides a new  
19 experimental strategy for structure determination, with the potential to reveal functional and *in*  
20 *vivo* structural conformations at low cost and high throughput.

21

## 22 Introduction

23 Mutations within a protein or RNA can have non-independent effects on fitness<sup>1-4</sup>. Indeed, the  
24 effects of double mutants have long been used to probe the energetic couplings between  
25 positions in a protein to understand determinants of protein folding and stability<sup>5,6</sup>. Early work  
26 revealed that at least some strongly interacting positions within a protein are in direct structural  
27 contact<sup>5-8</sup>. Deep mutagenesis of proteins<sup>9-12</sup> and RNAs<sup>13-16</sup> has further confirmed this  
28 conclusion that some – but by no means all – genetic (or epistatic) interactions occur between  
29 structurally proximal mutations.

30 Further support for the idea that non-independence between mutations provides structural  
31 information comes from the analysis of amino acid and nucleotide sequence evolution. Here,  
32 correlated pairs of amino acids or nucleotides in multiple sequence alignments identify co-  
33 evolving positions within proteins and RNAs<sup>17-19</sup>. These patterns of co-evolution have been  
34 used to identify energetically-coupled positions and ‘sectors’ within proteins<sup>20,21</sup>. Moreover,  
35 when very large numbers of homologous proteins and RNAs are available in sequence  
36 databases, the application of global statistical models has proven sufficient to discriminate direct  
37 structural contacts from patterns of co-evolution<sup>22-24</sup>, allowing the prediction of macromolecular  
38 structures and interactions<sup>25-34</sup>.

39 Could epistatic interactions quantified from deep mutational scanning experiments be used to  
40 determine macromolecular structures? If successful, structure determination by deep  
41 mutagenesis would offer a number of advantages over established techniques. First, it requires  
42 no specialized equipment or expertise beyond the ability to mutate a molecule, select functional  
43 variants, and quantify enrichments by sequencing. Appropriate *in vitro* and *in vivo* selection  
44 assays already exist for very many molecules of interest and generic assays based on folding,  
45 stability, and physical interactions have also been developed<sup>9,35-38</sup>. Second, it could be applied  
46 to molecules whose structures are difficult to determine by physical techniques such as  
47 intrinsically disordered and membrane proteins. Third, unlike evolutionary coupling analysis  
48 there is no requirement for large numbers of homologous sequences and so it could be applied  
49 to fast-evolving, recently-evolved and *de novo* designed proteins and RNAs<sup>26,32,39</sup>. Finally, and  
50 perhaps most importantly, it would provide a general strategy to determine the physiologically  
51 relevant structures of molecules whilst they are performing particular functions that can be  
52 selected for, including *in vivo* within cells. A cheap and straightforward approach for studying  
53 macromolecular structures *in vivo* would be a very exciting new frontier for cell biology.

54 Here we show that deep mutational scanning (DMS) of proteins can provide sufficient  
55 information to determine their high-resolution three-dimensional structures. Our statistical  
56 approach quantifies how often mutations between positions interact epistatically and how such  
57 epistatic interaction patterns correlate. These metrics accurately identify individual tertiary  
58 structure contacts as well as secondary structure elements within a protein. The same  
59 approach also identifies contacts between protein interaction partners. DMS data alone  
60 suffices to determine protein structures with accuracies down to 1.9Å backbone root mean  
61 square deviation (RMSD) compared to known reference structures. Moreover, we show that  
62 deep learning can further improve prediction performance, allowing the use of much sparser  
63 and lower quality DMS datasets for structure determination. This approach therefore provides a  
64 new experimental strategy for structure determination that can reveal functional and *in vivo*  
65 structural conformations at low cost and high throughput.

66

67

## 68 Results

### 69 Epistasis is enriched in but not exclusive to structural contacts

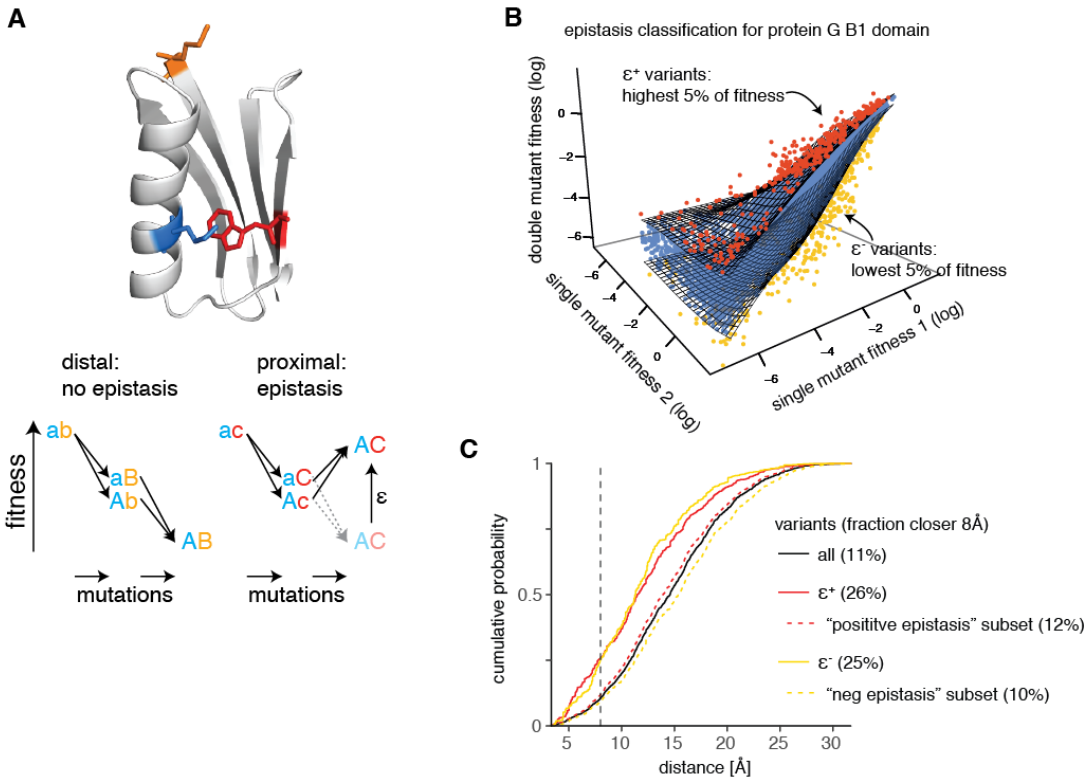
70 To investigate how genetic – or epistatic – interactions between mutations in a protein relate to  
71 structure we first used deep mutational scanning data for the immunoglobulin-binding protein G  
72 B1 domain (GB1) generated by Olson, et al.<sup>11</sup>. This dataset is the most complete double mutant  
73 deep mutagenesis of a protein domain reported to date and was generated by replacing each of  
74 55 residues of the wild-type domain with all 19 alternative amino acids both individually and as  
75 double mutant pairwise combinations, resulting in a library of more than half a million variants  
76 ( $55 \times 19 = 1,045$  single mutants plus nearly  $55 \times 54 / 2 \times 19 \times 19 = 536,085$  double mutants). mRNA  
77 display was used to combine an *in vitro* immunoglobulin G binding assay with a sequencing  
78 readout to determine protein fitness via changes in variant frequencies in the library before and  
79 after binding (Extended Data Figure 1, steps 1-4); resulting in a two orders of magnitude  
80 measurement range with a median relative error of fitness estimates of 2.8% (see Figure 2A,  
81 Table 1 and Methods).

82 We first computed which double mutant variants show epistatic fitness effects, i.e. non-  
83 independent fitness effects of the constituting single mutant variants (Figure 1B). Non-specific  
84 dependencies between mutants might be introduced by non-linearities in the fitness assay,  
85 systematic biases in error magnitudes as well as non-specific epistatic behavior, e.g. from  
86 thermodynamic stability effects<sup>1,9</sup>. We thus applied a non-parametric null model - the running  
87 median of double mutant fitness values given the constituting single mutant fitness values - for  
88 the independence of mutations. Equivalently, we calculated 5th and 95th percentile fitness  
89 surfaces; and classified double mutants with fitness lower than the 5<sup>th</sup> percentile as negative  
90 epistatic and double mutants with fitness higher than the 95<sup>th</sup> percentile as positive epistatic. We  
91 restricted the evaluation of positive or negative epistasis, however, to specific subsets of the  
92 data, where measurement errors do not impede epistasis classification (Extended Data Figure  
93 2C, see Methods), which results in about 80% and 55% of double mutants being suitable for  
94 positive or negative epistasis classification, respectively, with a lot of variability across the  
95 position matrix (Extended Data Figures 2D-F and Table 1).

96 Consistent with previous observations<sup>10-12</sup>, both positive and negative epistatic double mutants  
97 are enriched for proximal variants, for example, more than 2-fold at 8Å distance (side-chain  
98 heavy atom minimal distance, Figure 1C). However, about 75% of epistatic interactions are

99 between positions that are not in direct contact in the tertiary protein structure (as judged by an  
 100 8Å distance cutoff), suggesting that indirect effects often underlie epistatic interactions within a  
 101 molecule<sup>20,21</sup>. The challenge for structure determination therefore becomes how to infer direct  
 102 structural contacts from the mixture of direct and indirect effects that must underlie epistasis.

**Figure 1**



103

104 **Figure 1: Extracting epistatic mutational effects from deep mutational**  
 105 **scanning of a protein domain**

106 **A.** Premise: If epistatic interactions relate to structural contacts then quantifying epistatic  
 107 interactions should suffice to predict a molecule's structure. Structure: protein G B1  
 108 domain (PDB entry: 1pga) with residues a, b, and c colored.

109 **B.** Classifying epistatic variants based on deviations from expected fitness (based on  
 110 quantile fitness surface approach). Variants above the 95th or below the 5th percentile of  
 111 double mutant fitness given their single mutant fitness values were classified as positive  
 112 ( $\epsilon^+$ ) or negative ( $\epsilon^-$ ) epistatic, respectively. Shown is a random sample of  
 113  $10^4$  variants in GB1 domain<sup>11</sup>.

114 C. Distance distribution of epistatic variants separated by more than 5 amino acids in the  
115 linear sequence. (side-chain heavy atom distance in reference structure). Positive and  
116 negative epistasis subsets refer to the sets of variants applicable for epistasis analysis  
117 (see Extended Data Figure 2C). All variants,  $n = 400647$ ; positive epistatic variants  $\varepsilon^+$ ,  $n$   
118  $= 14127$ ; positive epistasis subset,  $n = 315862$ ; negative epistatic variants  $\varepsilon^-$ ,  $n = 9837$ ;  
119 negative epistasis subset,  $n = 208442$ .

120

121

## 122 Aggregated epistatic interactions predict tertiary structure 123 contacts

124 To distill direct contacts from a list of thousands of epistatic double mutants we first aggregated  
125 epistatic information on the amino acid position-pair level by calculating the fraction of positive  
126 or negative epistatic double mutant variants per position pair (Figure 2A).

127 In the GB1 epistasis dataset even moderate enrichments for positive and negative epistatic  
128 variants are mutually exclusive (Extended Data Figure 3A). Moreover, the strongest positive and  
129 negative epistatic enrichments are separated in two clusters of proximal positions in the protein  
130 that exhibit mostly either positive or negative interactions among themselves, but hardly any  
131 epistatic interactions between clusters (Figures 2B,D and Extended Data Figure 3B), as also  
132 noted before by Olson, et al. <sup>11</sup>.

133 Consistent with epistatic interaction clusters forming a dense network of proximal positions, we  
134 find that, of the top 55 epistatic pairs, 42% and 35% are direct contacts (connected by one edge  
135 smaller than 8 Ångström (Å), 3.9 and 3.2-fold over expectation) and another 45% and 55%  
136 share a common neighbor (connected via two edges  $< 8\text{Å}$ ), for positive and negative epistatic  
137 interactions respectively, while interactions across more edges are depleted (Figure 2C;  
138 throughout the manuscript we only consider position pairs spaced by more than 5 amino acids  
139 in the linear sequence; closer positions are trivially also close in 3D space, and their proximity  
140 contributes little to successful structure prediction <sup>28</sup>).

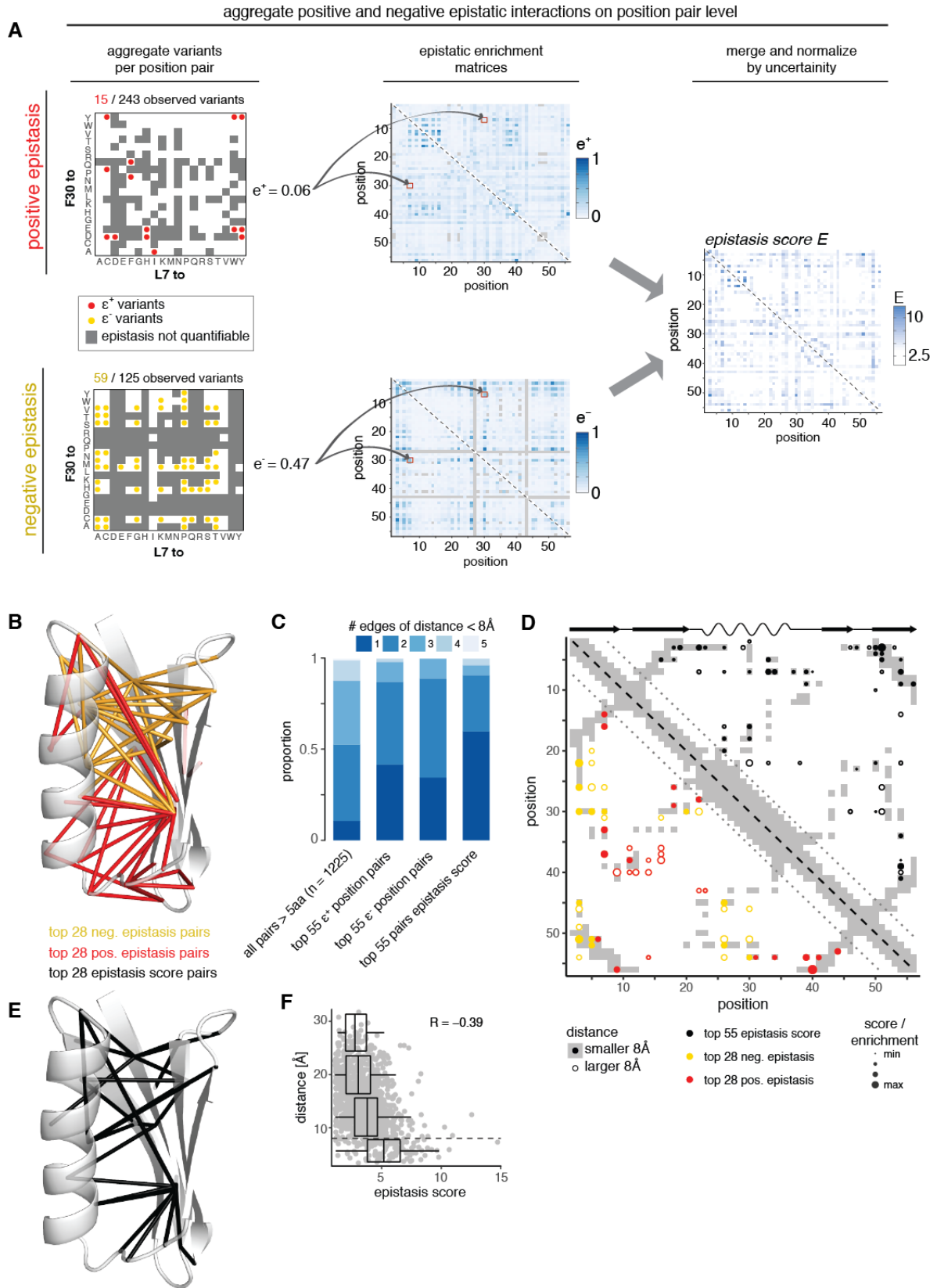
141 While aggregation of epistatic information between position pairs thus better discriminates  
142 structural contacts than individual epistatic interactions, positive and negative epistatic  
143 interactions still contain disparate structural information of the protein domain. We therefore  
144 merged positive and negative epistatic information by computing the weighted averages of

145 epistatic fractions per position pair given their uncertainty due to fitness measurement errors  
146 and the finite number of observed double mutant variants via a resampling approach (Extended  
147 Data Figure 1 and Methods). A final *epistasis score* per position pair was obtained by  
148 normalizing these weighted averages by their uncertainty (a z-score), thus giving priority to  
149 position pairs with high confidence enrichments.

150 The position pairs with highest *epistasis scores* are well distributed across the domain (Figures  
151 2D,E), and the number of direct contacts (one edge  $< 8\text{\AA}$ ) among the top 55 *epistasis score*  
152 pairs increases to 60% (Figure 2D), thus showing that the *epistasis score* successfully  
153 incorporates information from both positive and negative epistasis to discriminate direct  
154 contacts. Moreover, direct contacts as a whole are enriched for high *epistasis scores*, while  
155 further away position pairs show a gradual decrease of *epistasis scores* (Pearson correlation  
156 coefficient  $R = -0.39$ ,  $p < 10^{-6}$ , Figure 2F).

157 Thus, although many interactions are indirect, physical contacts are an important determinant of  
158 epistasis and aggregating information on position pairs and merging positive and negative  
159 epistasis information better discriminates these direct structural contacts across the protein  
160 domain.

**Figure 2**





## 162 Figure 2: Aggregated epistasis scores enrich for direct structural contacts

- 163 A. Workflow for aggregating positive and negative epistatic interactions on the position-pair  
164 level and merging them into a final *epistasis score*.
- 165 B. Top 28 position pairs (> 5 amino acids in linear sequence) each with highest positive  
166 (red) and negative (yellow) epistatic fractions marked on the reference structure (PDB  
167 entry 1pga).
- 168 C. Minimal number of edges (contact with distance < 8Å in reference structure) connecting  
169 position pairs. One edge – positions are direct contacts, two edges – positions have a  
170 common contact and so forth.
- 171 D. Interaction score map for top 28 position pairs each with highest positive (red) and  
172 negative (yellow) epistatic fractions (lower left triangle) and top 55 position pairs with  
173 highest *epistasis score* (upper right triangle). Dot size indicates relative epistatic  
174 enrichments or score; dot fill indicates distance below 8Å. Underlying in grey is the  
175 contact map of the reference structure (PDB entry 1pga, distance < 8Å) and shown on  
176 top its secondary structure elements (wave – alpha helix, arrow – beta strand).
- 177 E. Top 28 position pairs (> 5 amino acids in linear sequence) with highest *epistasis scores*  
178 marked on the crystal structure (PDB entry 1pga).
- 179 F. Distance of position pairs as a function of *epistasis scores*. Boxplots are spaced in  
180 distance intervals [0,8), [8,16), [16,24) and [24,32) Å. Dashed horizontal line indicates  
181 8Å. Pearson correlation coefficient is indicated.

182

## 183 Tertiary structure neighborhood leads to correlated epistatic 184 patterns

185 If epistasis arises mainly from structural interactions, a position's epistatic interaction profile with  
186 all other positions in the protein should provide a signature of its structural location (Figure 3A).  
187 Comparing these signatures between positions should thus reveal structurally close positions -  
188 similar to how correlated epistasis profiles in genetic interaction networks serve to identify  
189 physical and functional interaction partners<sup>40</sup>.

190 To test the idea that pattern correlation should reveal structural proximity, we calculated the  
191 correlations between the epistatic enrichment vectors for all position pairs (Figure 3B).  
192 Consistently, pair-distances and similarity of epistasis patterns between positions are strongly  
193 correlated (Pearson correlation coefficient = -0.43,  $p < 10^{-6}$ ,  $n = 1225$ , Figure 3D). Top  
194 correlated pairs from positive or negative interaction patterns do, however, form mutually  
195 exclusive clusters within the protein domain that are nearly identical to the clusters observed for  
196 direct positive and negative interactions (Figure 3C and E, c.f. Figure 2B). Thus, while  
197 correlations of epistatic interaction patterns are a good indicator of distance within the protein  
198 structure, they suffer from the same issues as epistatic enrichments, namely poor discrimination  
199 of direct and indirect interactions and disparate structural information.

200 We reasoned that partial correlations - the association between two positions after accounting  
201 for the global correlation structure - might provide the possibility to eliminate the dependencies  
202 observed in the epistasis pattern structure and thus help to distinguish direct from indirect  
203 contacts; similar to how mean-field approaches can help discriminate direct from indirect  
204 evolutionary couplings in multiple sequence alignments<sup>22,28,41</sup>. We derived partial correlations  
205 by inversion of the correlation matrices, merged values from positive and negative epistatic  
206 patterns by their estimated uncertainty, and ranked these merged values by their z-scores,  
207 which we refer to as *association scores* (Figure 3B and Methods).

208 In contrast to the correlation of epistasis patterns, partial correlation of epistasis patterns for  
209 both positive and negative epistasis display no clustering but are well distributed across the  
210 whole protein domain, consistent with partial correlations removing dependencies between  
211 correlated pairs (Figure 3C). Moreover, the merged *association scores* are less well correlated  
212 with pair-distance (Pearson correlation coefficient  $R = -0.26$ ,  $p < 10^{-6}$ ,  $n = 1225$ ) and show a  
213 more binary all-or-none response, with most distant position pairs having an *association score*  
214 around 0 and only proximal pairs systematically deviating to higher values (Figure 3D).  
215 Moreover, the top pairs involve many different individual positions and are well distributed  
216 across the protein domain (Figure 3C and E). Thus, *association scores* are able to prioritize  
217 direct over indirect structural contacts across the whole protein domain.

218

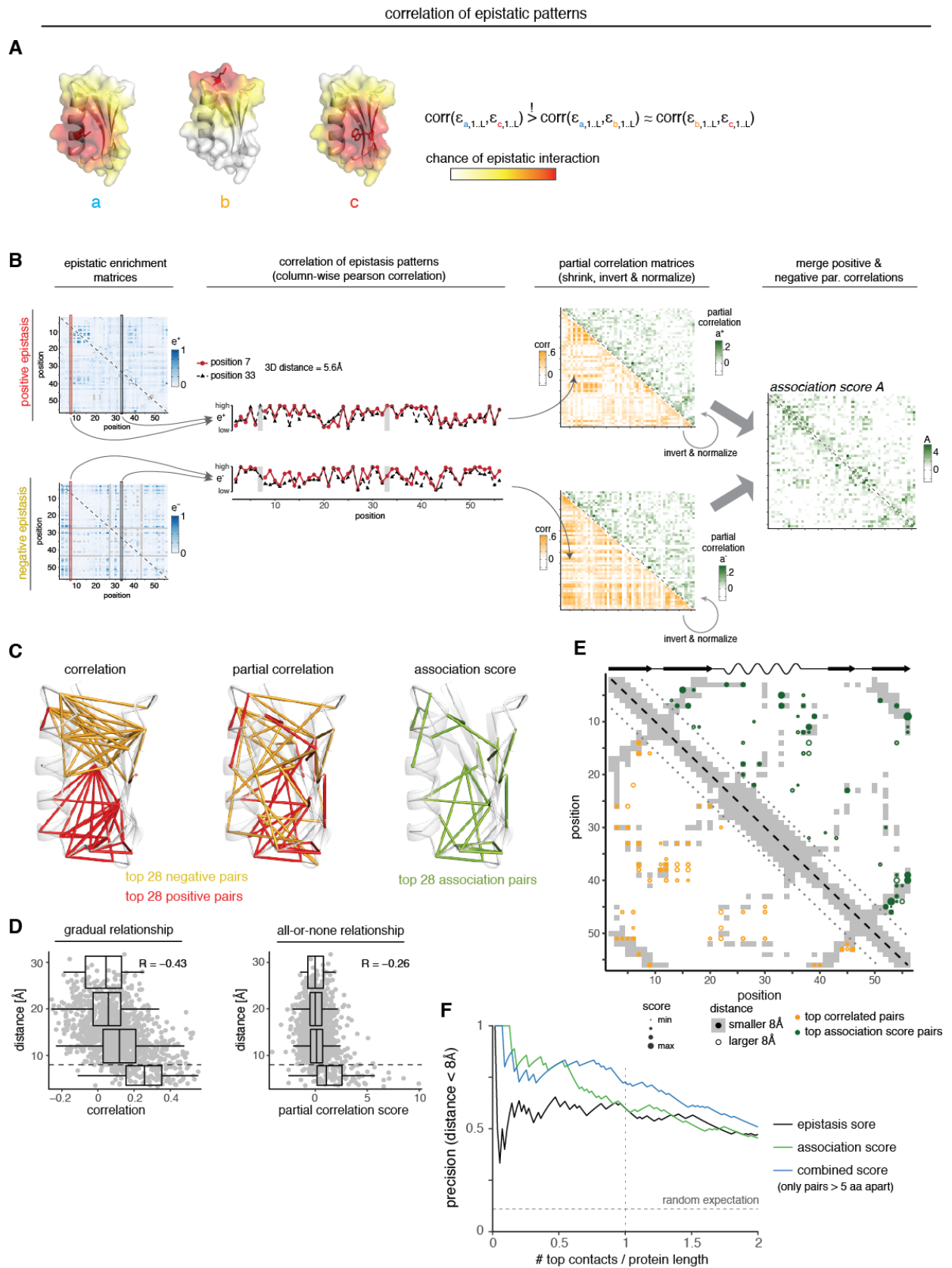
## 219 Combining *epistasis* and *association* scores better discriminates 220 structural contacts

221 We derived a *combined score* by summing the standardized *epistasis* and *association scores*,  
222 to explore whether combining information from individual epistatic interactions and epistasis  
223 interaction patterns can improve proximity estimates; thereby prioritizing position pairs that are  
224 both enriched for direct epistatic interactions and have correlated epistasis patterns.

225 We evaluated the precision of the three interaction scores in predicting direct contacts in the  
226 protein domain. For all pairs separated by more than 5 amino acids in the linear sequence, the  
227 *epistasis score* has a roughly constant precision of around 60% across the first  $2*L$  predicted  
228 contacts (L being the mutated length of the protein i.e. 55 amino acids). The *association score*  
229 has higher precision than the *epistasis score* up to the first L predicted contacts, with a precision  
230 of 79% at L/2 top contacts. Finally, the *combined score* has similar precision to the *association*  
231 *score* for the first L/2 contacts, but then remains at higher precision, with an improvement of  
232 about 10-15% over the individual scores at more predicted contacts (73% at L contacts);  
233 showing that combining information from epistatic interactions and interaction patterns further  
234 improves the discrimination of direct structural contacts.

235 Together, the derivation of the interaction scores demonstrates that it is possible to discriminate  
236 direct three-dimensional structural contacts from a mainly non-proximal set of epistatic  
237 interactions within a protein domain.

**Figure 3**



### 239 Figure 3: Tertiary structure neighborhood leads to correlated epistatic 240 patterns

- 241 A. Mutations in directly contacting residues should interact similarly with all other mutations  
242 in the protein. Thus, the similarity of epistasis patterns of two positions with all other  
243 positions in the protein should inform about their structural proximity.
- 244 B. Workflow for quantification of correlated epistasis patterns. Pairs of columns from  
245 epistatic enrichment matrices (here columns 7 and 33) are compared and their Pearson  
246 correlation coefficients are calculated, which constitute entries in the correlation matrix  
247 (here entries 7:33 and 33:7, due to matrix symmetry). Correlation matrices are inverted  
248 to yield the partial correlation matrices. Finally, entries of the positive and negative  
249 partial correlation matrices are merged (weighted average by uncertainty) and z-  
250 normalized to yield *association scores* (see Methods).
- 251 C. Top 28 position pairs (> 5 amino acids in linear sequence) marked on reference  
252 structure. Left: Top pairs from positive (red) or negative (yellow) epistasis pattern  
253 correlations. Middle: Top pairs after partial correlation transformation. Right: Top  
254 *association score* pairs (merge positive and negative partial correlations).
- 255 D. Distance of position pairs as a function of merged correlation (left) or *association scores*.  
256 Boxplots are spaced in intervals of 8Å. Dashed horizontal line indicates 8Å. Pearson  
257 correlation coefficient is indicated.
- 258 E. Interaction score map for top 55 position pairs with highest merge correlation (positive  
259 and negative correlations merged, lower left triangle, orange) and *association scores*  
260 (upper right triangle, green). Dot size indicates relative correlations or scores; dot fill  
261 indicates distance below 8Å. Underlying in grey is the contact map of the reference  
262 structure (PDB entry 1pga, distance < 8Å) and shown on top its secondary structure  
263 elements (wave – alpha helix, arrow – beta strand).
- 264 F. Precision of interaction scores to predict direct contacts (distance < 8Å in crystal  
265 structure 1pga) as a function of top scoring position pairs. Only position pairs with linear  
266 chain distance greater than 5 amino acids are considered (n = 1225 pairs, n = 131 direct  
267 contacts in reference structure). Horizontal dashed line indicates random expectation.

## 268 Periodic epistatic patterns reveal secondary structure 269 arrangements

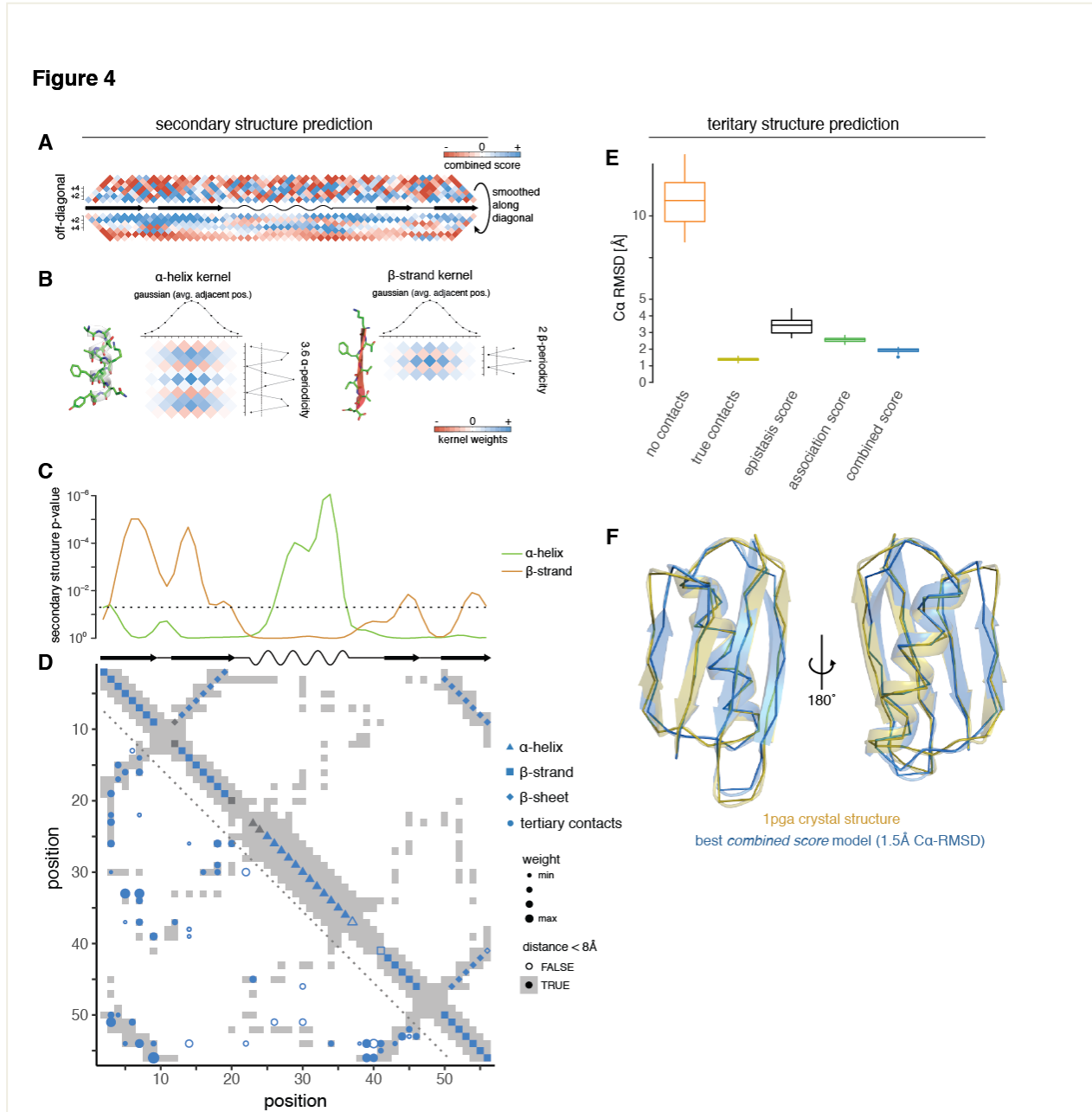
270 We investigated whether the periodic geometrical arrangement of amino acid residues in  
271 secondary structures results in periodic epistasis patterns<sup>26,42</sup>. Within an alpha helix with 3.6  
272 residues per helical turn, a helical position would be predicted to interact epistatically with the  
273 third or fourth-over position along the linear amino acid chain (Figure 4A). Equivalently, within a  
274 beta strand, positions should interact epistatically with the next-but-one position.

275 We used a two-dimensional kernel smoothing approach to estimate the positions of alpha  
276 helices and beta strands from the deep mutational scanning data (Figure 4B). Here, the  
277 propensity of a position to belong to an alpha helix or a beta strand depends on whether it  
278 shows the expected periodicity in its interaction with neighboring positions, as well as whether  
279 neighboring positions display similar propensities for the same secondary structure element,  
280 and how strong these interactions are compared to those found in randomized data sets (see  
281 Methods).

282 We found that, while secondary structure element predictions derived from direct interaction-  
283 based *epistasis scores* are somewhat inaccurate and underpowered, predictions derived from  
284 correlation-based *association scores* (as well as *combined scores*) coincide very well with  
285 secondary structure elements in the reference structure (Figure 4C and Extended Data Figure  
286 4C), with precision and recall values of about 90% (Extended Data Figure 4D). This suggests  
287 that the correlated profiles of epistatic interactions are informative about side chain orientations  
288 and also that eliminating transitive interactions is important for (secondary) structure prediction.

289 We further used two-dimensional kernel smoothing to detect parallel and anti-parallel beta sheet  
290 interactions, by applying beta strand kernels to off-diagonal entries on the interaction score  
291 matrices (Extended Data Figure 4A, see Methods). Several stretches of position pairs show the  
292 expected alternating interaction profiles for either parallel or anti-parallel beta sheets (Extended  
293 Data Figure 4B), with the top predictions corresponding to the known beta-sheet interactions in  
294 the reference structure (Figure 4D). Furthermore, updating beta strand predictions according to  
295 inferred beta sheet pairings can further improved beta strand prediction itself, notably  
296 introducing a correct split of beta strand 1 and 2 and adjusting the length of beta strands 3 and 4  
297 (Figure 4C,D).

298 Together this shows that epistatic interaction data contains information on the periodic  
 299 secondary structure of a protein domain and, *vice versa*, that secondary structure strongly  
 300 influences genetic interactions.



301

302 Figure 4: Secondary and tertiary structure prediction from deep mutational  
303 scanning data

- 304 A. Local interactions reveal signatures of secondary structure elements. Middle line is  
305 diagonal of interaction score map (rotated by 45 degree) and shows secondary structure  
306 elements of reference structure (PDB entry 1pga). Data above diagonal shows  
307 *combined score* data close to the diagonal, i.e. local interactions. Below the diagonal,  
308 the same data are smoothed with a Gaussian kernel along the direction of the diagonal  
309 (i.e. horizontally, length of Gaussian kernel as for kernels in panel b) to reveal  
310 periodicities in local interactions.
- 311 B. Two-dimensional kernels for alpha helix and beta strand detection. Kernel has a  
312 sinusoidal or alternating profile in the off-diagonal direction to detect alpha helices and  
313 beta strands propensities, respectively and a Gaussian profile along the diagonal, to  
314 average over propensities of adjacent positions.
- 315 C. Secondary structure propensity derived from kernel smoothing (orange – beta strand,  
316 green – alpha helix). P-values were derived by comparison to randomized datasets (see  
317 Methods). Dashed line indicates  $p = 0.05$ .
- 318 D. All structural predictions derived from *combined score* data. Lower left: Top 55 non-local  
319 (>5 aa in linear sequence) position pairs, i.e. tertiary contacts (circles); fill indicates  
320 correct prediction at 8Å, size of circles indicates relative score. Upper right: Predicted  
321 secondary structure elements (triangle – alpha helix, square – beta strand, diamond –  
322 beta sheet interaction). Fill indicates correct prediction. Note that beta strand predictions  
323 are derived by intersection of beta strand propensity (as shown in panel C) and results  
324 from beta sheet prediction (Extended Data Figure 4B, see Methods). Underlying in grey  
325 is the contact map of the reference structure (PDB entry 1pga, distance < 8Å) and  
326 shown on top are its secondary structure elements (wave – alpha helix, arrow – beta  
327 strand).
- 328 E. Accuracy ( $C\alpha$  root-mean-square deviation) of top 5% structural models generated from  
329 deep mutational scanning data derived restraints compared to GB1 reference structure.  
330 Structural models were generated in XPLOR-NIH by simulated annealing with restraints  
331 derived from top 55 top scoring position pairs, secondary structure element prediction  
332 and beta sheet pairing predictions from the indicated interaction scores. No contacts –  
333 negative control with restraints only for secondary structure (predicted by PSIPRED)<sup>43</sup>.



334 True contacts – positive control with 55 contacts (random subset), secondary structure  
335 elements and beta sheet interactions restraints derived from reference structure.  
336 F. Overlay of top structural model generated with restraints from *combined score* (blue) and  
337 crystal structure (gold, PDB entry 1pga). Shown is backbone ribbon and secondary  
338 structure cartoon generated in PyMOL <sup>44</sup>.

339

## 340 Protein structure determination by deep mutagenesis

341 Together, these findings show that deep mutagenesis data contain substantial information about  
342 a protein's secondary and tertiary structure. We therefore tested whether the data would suffice  
343 to determine *ab initio* the structure of protein G domain B1. We performed structural simulations  
344 by simulated annealing using the XPLOR-NIH modeling suite <sup>45</sup>, with structural restraints  
345 derived from the deep mutational scanning data (see Methods). In particular, we defined  
346 distance restraints (distance < 8Å between Cβ atoms) for the top scoring position pairs; we  
347 found that using the top L (L = 55) contacts gave best results (Extended Data Figure 4F).  
348 Furthermore, we defined dihedral angle restraints for predicted secondary structure elements.  
349 Finally, we defined restrictive distance restraints (distance smaller than 2.1Å for N-H : C=O atom  
350 pairs) for beta sheet positions that form hydrogen bonds with each other.

351 We evaluated the top 5% of structural models (25/500, evaluation based on XPLOR internal  
352 energy terms) generated against the known crystal structure of protein G domain B1 (PDB entry  
353 1pga) (Figures 4E and Extended Data Figure 4F). Models predicted from *combined score* data  
354 performed best, with an average Cα-root mean squared deviation of the top models ( $\langle C\alpha -$   
355  $RMSD \rangle$ ) of 1.9Å and an average template modeling score of 0.71, which is very close to the  
356 optimum achievable with our simulation protocol (using contacts, secondary structure elements  
357 and beta sheet interactions from the reference structure,  $\langle C\alpha - RMSD \rangle = 1.4\text{Å}$  and TM score =  
358 0.8); and the top evaluated *combined score* structural model has a  $C\alpha - RMSD$  of only 1.5Å  
359 (Figure 4F). Consistent with somewhat lower precision of contact and secondary structure  
360 predictions, models generated with restraints from *epistasis* or *association scores* have on  
361 average a lower accuracy ( $\langle C\alpha - RMSD \rangle = 3.4\text{Å}$  and  $\langle C\alpha - RMSD \rangle = 2.6\text{Å}$ , respectively), with  
362 *association score* models performing consistently better (Figures 4E and Extended Data 4F).

363 Together, this shows that deep mutation scanning alone is sufficient to accurately determine the  
364 structure of a protein domain.

## 365 Contact prediction in additional protein domains

366 To test the generality of our approach, we analyzed two additional, incomplete deep mutational  
367 scanning datasets. First, a mutational scan of the 75 amino acid Pab1 RRM2 domain (Figure  
368 5A), for which fitness was assessed in a complementation assay<sup>10</sup>. Second, a mutational scan  
369 of the hYAP65 WW domain (Figure 5C), in which 33 out of 50 amino acids were mutated and  
370 fitness was assayed by binding to a polyproline peptide ligand in a phage display assay<sup>46</sup>. Both  
371 datasets were created by ‘doped’ oligonucleotide synthesis and thus consist primarily of amino  
372 acid changes elicited by just one nucleotide change, which results in only 10% of possible  
373 double mutants being present. Additionally, their selection assays have smaller measurement  
374 ranges than that of the GB1 domain, which results in higher relative errors of fitness estimates  
375 as well as in negative epistasis being quantifiable for a smaller fraction of double mutants, as  
376 low as 0.8% in the case of the WW domain (Extended Data Figure 5A, see Table 1 for  
377 comparison of dataset properties).

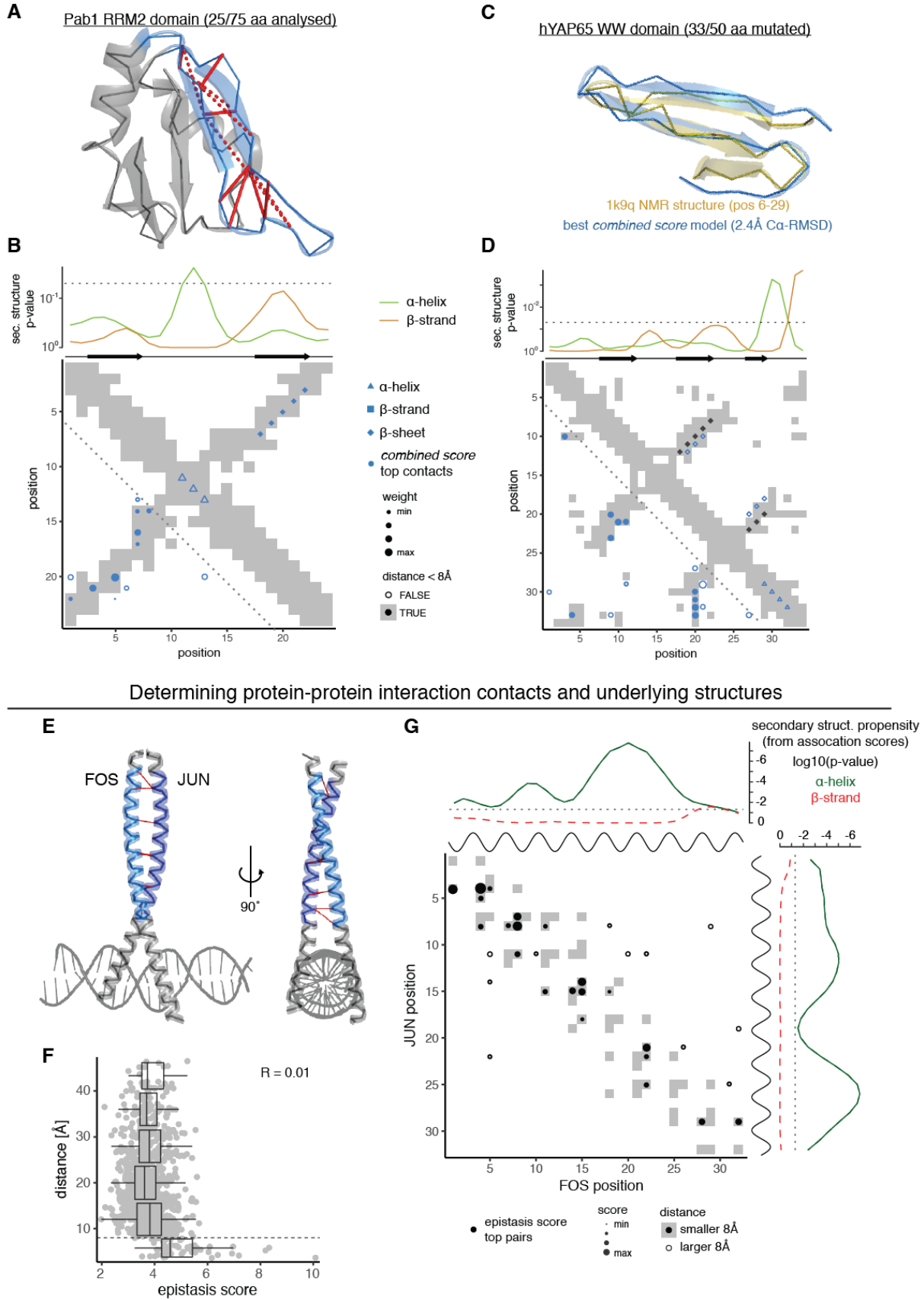
378 For the RRM domain, three 25 amino acid segments were mutated separately and we restricted  
379 our analysis to the central one, as it is the only segment that exhibits a reasonable number of  
380 intra-segment contacts in the reference structure (Figure 5A). We find that predicted tertiary  
381 contacts fall on or very close to known contacts in the region of the anti-parallel beta sheet and  
382 the intervening loop region (Figure 5B), with a precision of 57% for the top L/2 and 50% for the  
383 top L position pairs of the *combined score* (3-fold and 2.7-fold over expectation, respectively).  
384 Predicted beta strand propensities peak at the correct positions, albeit with low statistical  
385 significance; additionally, an alpha helical propensity is detected in the intervening loop region.  
386 Nonetheless, the correct anti-parallel beta sheet conformation at the correct position pairs is  
387 predicted.

388 For the WW domain, we find that top predicted tertiary contacts fall on or very close to known  
389 interactions between the beta strands and the N-terminal and C-terminal loop regions (Figure  
390 5D), with a precision of 59% for the top L/2 and 38% for the top L position pairs of the *combined*  
391 *score* (3.9-fold and 2.5-fold over expectation, respectively). Secondary structure elements are  
392 not well predicted, but beta sheet interactions are predicted in the right anti-parallel  
393 conformation of  $\beta_1 - \beta_2 - \beta_3$ , though the exact pairing between positions is off by one to two  
394 positions. We determined the three dimensional structure of the secondary structure-rich central  
395 part of the domain (positions 6 to 29, 24 amino acids), using restraints derived from top  
396 *combined score* pairs and PSIPRED-predicted secondary structure elements (see Methods).

397 The top 5% of structural models have an average accuracy of  $3.3\text{\AA}$  ( $C\alpha - RMSD$ ) compared to  
398 the reference structure, which is on par with simulations using a set of ‘true’ contacts ( $\langle C\alpha -$   
399  $RMSD \rangle = 3.6\text{\AA}$ ) (Extended Data Figure 5C). Moreover, the structural model with the best  
400 XPLOR-NIH energy has an accuracy of  $2.4\text{\AA}$   $C\alpha - RMSD$  (or  $2.0\text{\AA}$  over 22 of the 24 residues)  
401 (Figure 5C). Despite similar precision of predicted contacts, *association* and *combined score-*  
402 *derived* WW domain structural models are more accurate than *epistasis score-* derived models  
403 (Extended Data Figure 5C).

404 Together these results strongly support the generality of our approach for extracting structural  
405 information from deep mutagenesis data, including from sparser and lower quality data.

**Figure 5**



407 Figure 5 – Predicting structural contacts in two additional proteins and a  
408 protein-protein interaction

- 409 A. Pab1 RRM2 domain structure (PDB entry 1cvj) with 25/75 positions analyzed here  
410 highlighted in blue. Top 12 *combined score* position pairs are connected with red lines,  
411 solid if distance < 8Å, dashed otherwise.
- 412 B. Structural predictions derived from *combined scores* in RRM domain. Upper plot shows  
413 secondary structure propensities from kernel smoothing ( $p = 0.05$  indicated as dashed  
414 line). Just below are shown the secondary structure elements in the reference structure.  
415 Map shows top 12 *combined score* position pairs in lower left and secondary structure  
416 predictions in upper right triangle. Shape indicates type of prediction, fill indicates correct  
417 prediction. Underlying is the contact map of the reference structure (grey if < 8Å).
- 418 C. Overlay of top structural model of hYAP65 WW domain (positions 6-29) generated with  
419 restraints from *combined score* (blue) and solution NMR structure (gold, PDB entry  
420 1k9q).
- 421 D. Structural predictions derived from *combined scores* in WW domain. Upper plot shows  
422 secondary structure propensities from kernel smoothing ( $p = 0.05$  indicated as dashed  
423 line). Just below are shown the secondary structure elements in the reference structure.  
424 Map shows top 17 *combined score* position pairs in lower left and secondary structure  
425 predictions in upper right triangle. Shape indicates type of prediction, fill indicates correct  
426 prediction. Underlying is the contact map of the reference structure. Black diamonds  
427 indicate positions of beta sheet pairing in reference structure. Crystal structure of the  
428 leucine zipper domains of FOS and JUN with a DNA strand (PDB entry 1fos). The  
429 mutated regions (32 amino acids each) are highlighted in light blue (FOS) and dark blue  
430 (JUN). Top 10 epistasis score pairs are shown with red dashes.
- 431 E. Distance of position pairs as a function of interaction scores. Boxplots are spaced in  
432 distance intervals of 8Å. Dashed horizontal line indicates 8Å. Pearson correlation  
433 coefficient is indicated.
- 434 F. FOS-JUN trans interaction score map for top 32 position pairs with highest *epistasis*  
435 *scores*. Note that protein-protein interaction maps are not symmetric. Dot size indicates  
436 relative score; dot fill indicates distance below 8Å; underlying in grey is the contact map  
437 of the reference structure (PDB entry 1fos, distance < 8Å). Shown on top and to the right

438 of the contact map are the known alpha helices and secondary structure propensities  
439 derived from *association scores* of FOS and JUN, respectively (black – known alpha  
440 helix; green – predicted alpha helix propensity, orange - predicted beta strand  
441 propensity; see Extended Data Figures 5F,G).

## 442 Contact prediction in a protein-protein interaction

443 Genetic interactions do not only occur between mutations within individual proteins but also  
444 between molecules that physically interact <sup>2</sup>. We investigated a deep mutational scanning  
445 dataset of the coiled-coil interaction between the proteins encoded by the proto-oncogenes *FOS*  
446 and *JUN* (Figure 5E) <sup>9</sup>. In this experiment, all possible single amino acid changes were made in  
447 each of 32 positions of each protein and the physical interaction of all single and (trans-)double  
448 mutants was quantified using a deep sequencing-based protein complementation assay. After  
449 filtering, the dataset contains 43% of all possible double mutants and has a median relative  
450 error of fitness measurements of 3.6% (Table 1).

451 When assessing the enrichment of epistatic interactions between positions in the two interaction  
452 partners we find a striking all-or-none relationship between *epistasis scores* and pair-distances  
453 (Figure 5F), with all distant pairs contained in a low *epistasis score* peak and only proximal  
454 interactions enriched for high *epistasis scores* (Pearson correlation coefficient  $R = 0.01$ ,  $n =$   
455  $1024$ ). Indeed, the top 11 *epistasis score* pairs are all proximal interactions, and the precision of  
456 contact prediction is 75% for the top  $L/2$  contacts and 66% for the top  $L$  contacts (12-fold and  
457 10.5-fold over expectation). Moreover, top *epistasis score* pairs are evenly distributed across  
458 the interaction surface (Figures 5E and 5G).

459 When correlating epistatic patterns between columns of the epistatic enrichment matrices, one  
460 is comparing the epistatic interactions that two positions in FOS have with all positions in JUN.  
461 Therefore, the similarity of column-wise epistatic patterns reveals the *cis* relationships between  
462 positions in FOS (Extended Data Figure 5F). Similarly, correlating epistasis patterns across  
463 combinations of rows of the epistatic enrichment matrices reveals *cis* relationships between  
464 positions in JUN. We find that *cis*-interaction maps from *association scores* for both FOS and  
465 JUN are highly enriched for strong local interactions with an alpha helical periodicity; and  
466 applying our secondary structure prediction algorithms to the *cis*-interaction maps reveals strong  
467 alpha helix propensities across the full lengths of both FOS and JUN (Figures 5G and Extended  
468 Data Figure 5G).

469 This shows that deep mutagenesis of protein interaction partners can accurately predict direct  
470 contact across the interaction surface as well as the underlying structures of the interaction  
471 partners themselves.

472

## 473 Deep learning improves contact prediction

474 Evolutionary coupling-based structural predictions have been successfully improved by machine  
475 learning approaches that transform the two-dimensional interaction score maps after learning  
476 the stereotypical patterns between evolutionary coupling-predicted contact maps and the actual  
477 contact maps of the known structures<sup>47,48</sup>.

478 We tested whether such an approach can also improve deep mutagenesis-derived contact  
479 predictions. We applied a convolutional neural network approach called *DeepContact*,  
480 developed by Liu, et al.<sup>47</sup>. The basic *DeepContact* architecture takes as a sole input a two-  
481 dimensional interaction score map that it then transforms based on the structural patterns it has  
482 previously learned on evolutionary coupling-derived contact predictions for representative  
483 families of the SCOPe database<sup>49</sup> (Figure 6A and Methods). When transforming evolutionary  
484 coupling-derived contact predictions of proteins not contained in the training set, this basic  
485 *DeepContact* architecture has been shown to improve contact prediction precision by about 10-  
486 20%<sup>47</sup>.

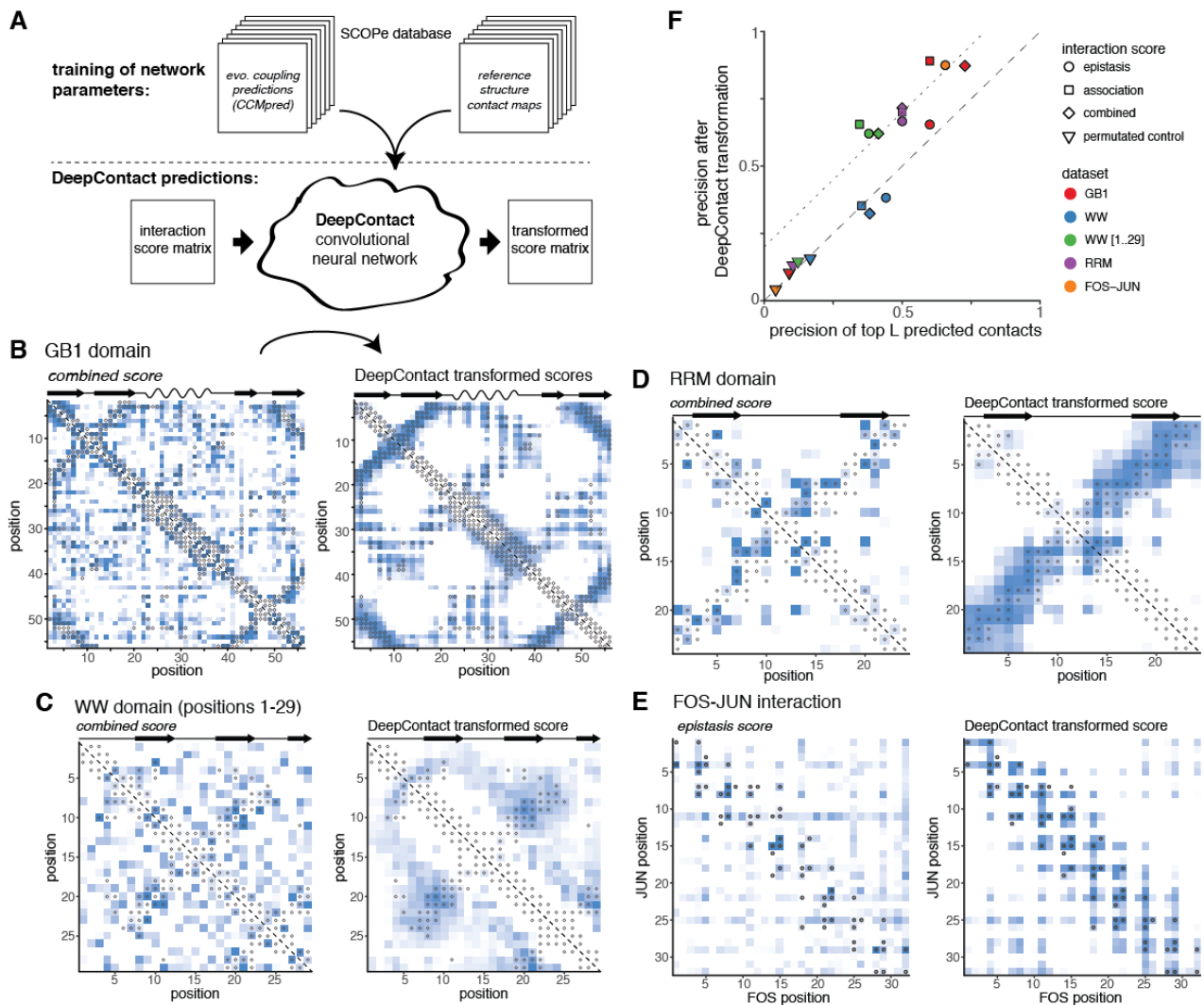
487 We first transformed the GB1 domain *combined score* interaction map with the *DeepContact*  
488 network (Figure 6B). These transformations take as sole input our deep mutational scanning-  
489 derived predictions and include no evolutionary coupling or otherwise-derived structural  
490 predictors for GB1. The scores on the transformed map are much less noisy, with high scores  
491 exclusively focused in areas of known contacts, especially those of secondary structure element  
492 interactions, and areas devoid of contacts showing homogenously low scores. Moreover, the  
493 transformed scores are highly correlated with pair-distances in the reference structure (Pearson  
494 correlation coefficient  $R = -0.68$ ,  $p < 10^{-6}$ ,  $n = 1225$ , Extended Data Figure 6A). The precision of  
495 top predicted contacts improves from 82% to 96% for L/2 and from 73% to 87% for L predicted  
496 contacts (Figure 6F). *Epistasis score*-derived predictions improve by about 5%, while  
497 *association score*-derived predictions improve by 29% at L predicted contacts. In contrast,  
498 randomized interaction score maps show no changes in prediction performance over random  
499 expectation after transformation with *DeepContact*.

500 Interaction score maps for the other datasets show similar improvements to GB1 both in terms  
 501 of cleaner interaction score maps that resemble the reference contact maps as well as  
 502 increases in contact prediction precision of up to 30% (Figure 6 C-F).

503 This shows that machine learning can substantially improve contact map prediction from deep  
 504 mutational scanning data, allowing the use of sparser and lower quality data for accurate  
 505 prediction.

506

**Figure 6**



507



508 Figure 6: Deep learning improves contact prediction from deep  
509 mutagenesis data

- 510 A. DeepContact convolutional neural network transforms DMS-derived interaction score  
511 maps based on learned structural patterns<sup>47</sup>. The particular DeepContact architecture  
512 used here takes as only input the DMS-derived interaction score map and transforms it  
513 based on structural patterns previously learned on an orthogonal and independent  
514 training set (in which it compared evolutionary coupling-derived contact predictions with  
515 contacts in known structures of representative protein families in the SCOPe database).
- 516 B. GB1 domain *combined score* interaction map before (left panel) and after (right panel)  
517 transformation with *DeepContact* convolutional neural network. Heat maps show scores  
518 (normalized to have similar range). Grey open circles show contacts (side-chain heavy  
519 atom distance < 8Å) in reference structure.
- 520 C. WW domain *combined score* interaction map before (left) and after (right) DeepContact  
521 transformation. Note that the maps shown here lack the 5 c-terminal positions (see  
522 Extended Data Figure 6B for full map).
- 523 D. RRM domain *combined score* interaction map before (left) and after (right) *DeepContact*  
524 transformation.
- 525 E. FOS-JUN trans-interaction *epistasis score* interaction map before (left) and after (right)  
526 *DeepContact* transformation.
- 527 F. Precision of top L predicted contacts of different interaction scores before and after  
528 *DeepContact* transformation for the four datasets. Color indicates dataset, shape  
529 indicates interaction score. Permuted control score is average over three random  
530 permutations of *combined score* matrices (in case of FOS-JUN *epistasis score*  
531 matrices). Dashed diagonal line indicates no changes in precision, dotted diagonal line  
532 shows precision improvement of 20% after DeepContact learning.

533

534

## 535 Minimal data quality requirements for successful protein structure 536 prediction

537 We further investigated how robust our prediction strategy is to changes in data quality by  
538 artificially down-sampling the GB1 domain dataset, thus assessing the minimal requirements for  
539 deep mutational scanning datasets to be useful for structure prediction.

540 First, we considered the sequencing read coverage. The GB1 domain dataset consists of about  
541 600 million sequencing reads<sup>11</sup>. We find that artificially down-sampling the sequencing read  
542 coverage of the dataset to 25% or 10% hardly affects the precision of predicted tertiary contacts  
543 (Figures 7A). Only when using just 2.5% of sequencing reads (15 million) does the precision of  
544 top L contacts drop below 50%.

545 Next, we simulated a ‘doped’ dataset, by only considering amino acid mutations that can be  
546 reached by one mutation in the nucleotide sequence - thus reducing the coverage of double  
547 mutants to ~10% (similar to RRM and WW domain datasets). The doped dataset with full  
548 sequencing read coverage exhibits a drop in precision of predicted tertiary contacts of about  
549 20%. Moreover, the doped dataset shows an increased sensitivity to lower sequencing read  
550 coverage.

551 We also tested the effect of reducing the signal-to-noise ratio (i.e. the measurement range of  
552 selection assay relative to the median error of fitness estimates), which results in non-  
553 quantifiably of negative epistasis (Extended Data Figures 2D-F and 5A). We thus tested how  
554 our prediction strategy performs on the GB1 domain dataset when only positive epistasis  
555 information is available; and find that it results in a drop of precision of about 20%, comparable  
556 to that observed for a doped dataset. In contrast, only using negative epistasis information  
557 resulted in a drop to ~35% precision, as low as a doped dataset with 10% sequencing  
558 coverage.

559 We evaluated secondary structure prediction performance of the various down-sampled GB1  
560 domain datasets. Beta strand and alpha helix predictions are hardly affected by lowered data  
561 quality or partial epistasis information (Extended Data Figure 7A). In contrast, precision and  
562 recall of beta sheet positional pairing is strongly affected by dataset quality, although often the  
563 correct overall conformation of beta sheets is still recovered (Extended Data Figure 7B).

564 We next tested whether DeepContact could also improve prediction performance on these  
565 down-sampled datasets. Similar to interaction scores derived from full datasets, DeepContact

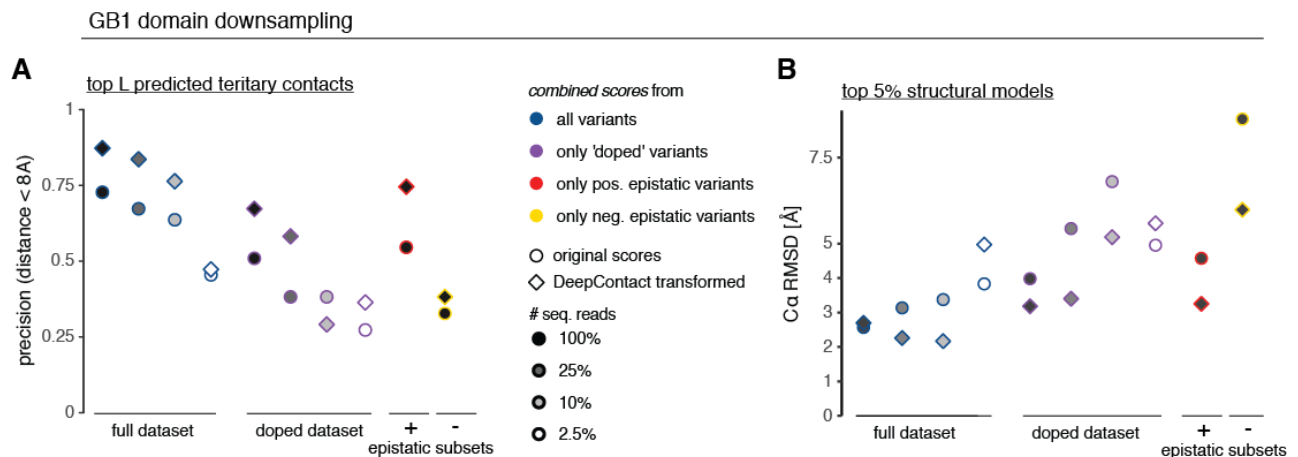
566 transformation of *combined scores* derived from down-sampled GB1 datasets improves the  
 567 precision of predicted contacts by about 10-25% even for quite low quality datasets, i.e. the  
 568 complete datasets with at least 10% read coverage, doped datasets with at least 25% read  
 569 coverage or the dataset with only positive epistasis information (Figure 7A).

570 Finally, we evaluated how differences in prediction performance of tertiary contacts affect  
 571 structural modeling. We find that changes in accuracy of the top structural models roughly scale  
 572 with changes in contact prediction performance (Figure 7B). Down-sampling of sequencing  
 573 reads in the complete dataset from 100% to 2.5% leads to a drop in accuracy from 2.5Å to 4Å  
 574 ( $\langle C\alpha - RMSD \rangle$ ), which is roughly also the accuracy of top structural models from the doped  
 575 dataset and the dataset using only positive epistasis information. Accuracies for lower quality  
 576 datasets range from 5Å to 9Å ( $\langle C\alpha - RMSD \rangle$ ). DeepContact increases the accuracy of the top  
 577 structural models by up to 2.6Å. For the complete datasets with only 25% or 10% of sequencing  
 578 reads, the top structural models have better accuracy than those from the complete dataset with  
 579 full sequencing read coverage but untransformed scores. Also, structural models based on  
 580 DeepContact transformed scores from the doped dataset with full or 25% sequencing coverage  
 581 and those from the dataset using only positive epistasis information reach average accuracies  
 582 of 3.2Å ( $\langle C\alpha - RMSD \rangle$ ). Only for the two datasets with 2.5% sequencing read coverage do  
 583 structural simulations based on DeepContact transformed scores not improve model accuracy.

584 Together these findings show that contact and structure prediction from deep mutational  
 585 scanning data can also work for lower quality datasets and that the use of deep learning allows  
 586 the use of much sparser and lower quality datasets.

587

**Figure 7**



588

589 Figure 7: Deep learning allows contact and structure prediction from  
590 sparser and lower quality datasets

- 591 A. Precision of top L *combined score* position pairs for different down-sampled versions of  
592 GB1 dataset. Color indicates type of dataset (blue – full dataset, purple – ‘doped’  
593 dataset, red – only positive epistasis information, yellow – only negative epistasis  
594 information), fill indicates number of sequencing reads used in analysis, shape indicates  
595 whether DeepContact has been used to transform the interaction score matrix.
- 596 B. Accuracy ( $C\alpha - RMSD$ ) of top 5% structural models derived with tertiary contact  
597 restraints from down-sampled GB1 datasets compared to reference structure. Note that  
598 for better comparability, for these structural simulations only distance restraints were  
599 derived from *combined scores* but the same secondary structure restraints predicted  
600 from PSIPRED and no beta sheet pairing restraints were used for all simulations. Colors  
601 and fills as in panel A.

## 602 Discussion

603 We have shown here that simply quantifying the activity of a large number of single and double  
604 mutant variants of a macromolecule can provide enough information to determine its high-  
605 resolution 3D structure.

606 We found that although most epistasis within a protein occurs between positions that are not  
607 direct structural contacts, aggregation on position pairs, merging of positive and negative  
608 epistasis information and partial correlation analysis of epistasis patterns can successfully  
609 discriminate direct from indirect structural contacts. Thus, mostly indirect epistatic couplings can  
610 be transformed to predict accurate structural contacts and elements. We have shown that this  
611 approach works robustly across multiple protein domains and a protein interaction. Moreover,  
612 we have demonstrated that the application of a convolutional neural network previously trained  
613 on patterns of co-evolution in proteins of known structure both improves structure prediction and  
614 allows the use of much lower quality deep mutation datasets.

615 Determining structures by mutagenesis requires an *in vitro* or *in vivo* selection assay. For many  
616 important molecules and drug targets, specific selection assays based on known functions or  
617 interaction partners already exist <sup>11,14,15,46,50-55</sup>. Additionally, many generic selection assays have  
618 recently been developed that should allow the stability or functional activity of many proteins to  
619 be assayed *in vivo* without the need for much prior knowledge about the protein under  
620 investigation <sup>35-38</sup>. Moreover, many molecules have known interaction partners – proteins, DNA,  
621 RNA, or small ligands – for which *cis*- and *trans*-epistasis can thus be assessed by binding  
622 assays <sup>9,55</sup>. We have shown here how *trans*-epistasis, for which library design is relatively easy,  
623 can lead to information about direct contacts in interaction surfaces as well as in the individual  
624 molecules.

625 Although structural information exists in the epistasis maps, our analyses and previous work <sup>5-16</sup>  
626 has shown that many epistatic interactions occur between positions that are not in direct  
627 structural contact. Indeed, in the GB1 domain the interactions are strikingly modular, with two  
628 mutually exclusive clusters of positive and negative epistatic interactions. This is consistent with  
629 many interactions being due to functional or energetic couplings between positions. The cluster  
630 of mostly positive epistatic interactions corresponds to a dynamic region involved in IgG binding  
631 <sup>11</sup>. In contrast, the cluster of negative epistatic interactions identifies positions important for the  
632 thermodynamic stability of the domain <sup>56</sup>, and the periodicity of local negative epistatic

633 interactions provide evidence for a shift towards an alternative three-helical conformation that  
634 has been previously reported for this sequence family (Extended Data Figure 4E)<sup>57</sup>. This  
635 modular organization of epistasis is thus reminiscent of the concept of energetically-coupled  
636 protein sectors identified from patterns of sequence co-evolution<sup>20,21,58</sup>.

637 For macromolecules with very large numbers of homologs available in sequence databases,  
638 correlated changes in sequence can provide sufficient information for structure determination<sup>25-</sup>  
639<sup>34</sup>. However, for many proteins and RNAs insufficient numbers of homologs are available, and  
640 for fast evolving, recently-evolved or *de novo* designed molecules this is a fundamental  
641 limitation<sup>26,32,37,39</sup>. Moreover, co-evolutionary analysis provides information on the average  
642 structure across a large set of homologs, whereas it is easy to envisage how deep mutagenesis  
643 could be used to directly determine alternative conformations of macromolecules when they are  
644 performing particular selectable functions. The success of evolutionary coupling analysis for  
645 predicting the structures of diverse folds and macromolecules does, however, strongly support  
646 the generality of the approach outlined here. The demonstration that a deep learning approach  
647 previously trained on evolutionary couplings dramatically improves the prediction of contacts  
648 from deep mutagenesis data further supports this.

649 As a proof-of-principle we have shown that information from deep mutational scanning  
650 experiments alone is sufficient for accurate structure prediction. In practice, however, integration  
651 with other structural information is likely to further boost performance. As a first step, we used a  
652 deep learning approach, DeepContact, that was trained on evolutionary couplings to learn  
653 stereotypical structural patterns in contact maps<sup>47</sup>. DeepContact improved DMS-derived  
654 contact prediction precision by up to 30% for individual proteins (GB1, WW and RRM domains).  
655 Moreover, even though it had only been trained on data from individual proteins, it also  
656 improved DMS-derived contact predictions for the FOS-JUN protein-protein interaction.  
657 Integration with other structural predictors<sup>47,59,60</sup> and homology-driven structure modeling<sup>61,62</sup> is  
658 likely to further improve accuracy and lower the data quality requirements for structure  
659 determination by deep mutagenesis.

660 An analysis of incomplete and down-sampled variants of the GB1 dataset suggests that a high  
661 signal-to-noise ratio (measurement range relative to experimental error of fitness estimates),  
662 which allows both positive and negative epistasis to be quantified, is an important factor for  
663 generating datasets with a quality sufficient for protein structure prediction. For datasets with  
664 complete epistasis information, however, sequencing coverage hardly affected prediction

665 performance. In contrast, prediction performance on the incomplete, ‘doped’ dataset was  
666 sensitive to sequencing coverage. Down to 25% sequencing coverage, however, performance  
667 could be recovered by deep learning. Together, these analyses suggest that our approach  
668 should be easily applicable to longer molecules. For example, with the experimental effort  
669 undertaken to create and sequence the 55 amino acid GB1 library, a protein of length ~350  
670 amino acids should be assayable at similar prediction performance (using a doped library with  
671 25% sequencing coverage, i.e. 2.5% of the data:  $55aa/\sqrt{0.025} \approx 350aa$ ). Such libraries for  
672 longer proteins could be created via fragment-based ligation approaches<sup>52</sup> or via random  
673 mutagenesis and barcode-variant linking<sup>35</sup>.

674 Taken together, the results presented here establish deep mutagenesis as a new experimental  
675 strategy for structure determination. The approach that we have outlined is not the only one  
676 that can be envisaged to predict direct structural contacts from deep mutagenesis data, and  
677 other related approaches are also likely to work<sup>63</sup>. The determination of macromolecular  
678 structures by physical techniques requires access to very expensive scientific infrastructure. In  
679 contrast, deep mutagenesis only requires techniques familiar to many molecular biologists and  
680 access to sequencing that is increasingly low cost and available to all. Most importantly,  
681 however, deep mutagenesis allows the structures of macromolecules to be studied whilst they  
682 are performing particular functions *in vitro* as well as *in vivo* in the cell. As such, deep  
683 mutagenesis opens up the possibility of low cost and high throughput determination of *in vivo*  
684 macromolecular structures by many molecular biology and genomics labs. A large-scale project  
685 to systematically determine the structures of proteins and protein domains should therefore be  
686 possible using the existing infrastructure of genomics institutes.

687

688 Table 1: Dataset properties

Dataset	Mutated positions	% double mutants <sup>§</sup>	% doubles quantifiable <sup>#</sup>		# input reads per double mutant (median) <sup>*</sup>	measurement range (log fitness units) <sup>+</sup>	relative error (median) <sup>&amp;</sup>
			positive epistasis	negative epistasis			
Protein G B1 domain	55	97	80	55	248	5.8	2.8%
hYAP WW domain	33	10	8.3	0.8	73	0.8	8.6%
Pab1 RRM2 domain	25	11	8.3	3.9	209	3.1	3.7%
FOS-JUN	2 x 32	43	37	31	124	8.6	3.6%

689 <sup>§</sup> median percentage of all possible double mutants (361 per position pair) that passed read quality thresholds per  
690 position pair

691 <sup>#</sup> median percentage of all possible double mutants (361 per position pair) that passed read quality thresholds and  
692 are deemed suitable for epistasis quantification per position pair

693 <sup>\*</sup> summed number of reads across all input replicates for double mutants that passed read quality thresholds

694 <sup>+</sup> measurement range of selection assay: log fitness range between peak of lethal mutants and the wild-type variant

695 <sup>&</sup> median error of fitness estimates of double mutant variants relative to measurement range of selection assay

## 696 Acknowledgements

697 We are grateful to Yang Liu and Jian Peng for making their DeepContact code available and for  
698 their advice. We thank members of the Lehner lab, T. Gross, G. Mönke, M. Bolognesi and C.  
699 Camilloni for discussions and feedback. This work was supported by a European Research  
700 Council (ERC) Consolidator grant (616434), the Spanish Ministry of Economy and  
701 Competitiveness (BFU2011-26206 and SEV-2012-0208), the AXA Research Fund, the  
702 Bettencourt Schueller Foundation, Agencia de Gestio d'Ajuts Universitaris i de Recerca  
703 (AGAUR, SGR-831), the EMBL-CRG Systems Biology Program, and the CERCA  
704 Program/Generalitat de Catalunya. J.M.S. was supported by an EMBO Long-Term Fellowship  
705 (ALTF 857-2016). This project has received funding from the European Union's Horizon 2020  
706 research and innovation programme under the Marie Skłodowska-Curie grant agreement No  
707 752809 (J.M.S).

708



## 709 Author Contributions

710 Conceptualization, B.L. and J.M.S.; Methodology, J.M.S.; Investigation, J.M.S.; Writing, J.M.S.  
711 and B.L.; Supervision, B.L.

712

## 713 Competing Interests

714 The authors declare no competing interests.

715

## 716 Additional Information

717 **Correspondence and requests for materials** should be addressed to B.L.

718

## 719 References

- 720 1 Tokuriki, N. & Tawfik, D. S. Stability effects of mutations and protein evolvability. *Current*  
721 *Opinion in Structural Biology* **19**, 596-604, doi:10.1016/j.sbi.2009.08.003 (2009).
- 722 2 Lehner, B. Molecular mechanisms of epistasis within and between genes. *Trends in*  
723 *genetics : TIG* **27**, 323-331, doi:10.1016/j.tig.2011.05.007 (2011).
- 724 3 Fowler, D. M. & Fields, S. Deep mutational scanning: a new style of protein science.  
725 *Nature Methods* **11**, 801-807, doi:10.1038/nmeth.3027 (2014).
- 726 4 Starr, T. N. & Thornton, J. W. Epistasis in protein evolution. *Protein science : a*  
727 *publication of the Protein Society* **25**, 1204-1218, doi:10.1002/pro.2897 (2016).
- 728 5 Horovitz, A. & Fersht, A. R. Strategy for analysing the co-operativity of intramolecular  
729 interactions in peptides and proteins. *J Mol Biol* **214**, 613-617, doi:10.1016/0022-  
730 2836(90)90275-Q (1990).
- 731 6 Carter, P. J., Winter, G., Wilkinson, A. J. & Fersht, A. R. The use of double mutants to  
732 detect structural changes in the active site of the tyrosyl-tRNA synthetase (*Bacillus*  
733 *stearothermophilus*). *Cell* **38**, 835-840 (1984).
- 734 7 Ackermann, E. J., Ang, E. T., Kanter, J. R., Tsigelny, I. & Taylor, P. Identification of  
735 pairwise interactions in the alpha-neurotoxin-nicotinic acetylcholine receptor complex  
736 through double mutant cycles. *J Biol Chem* **273**, 10958-10964 (1998).
- 737 8 Chen, J. & Stites, W. E. Energetics of side chain packing in staphylococcal nuclease  
738 assessed by systematic double mutant cycles. *Biochemistry* **40**, 14004-14011 (2001).

- 739 9 Diss, G. & Lehner, B. The genetic landscape of a physical interaction. *eLife* **7**, 594,  
740 doi:10.7554/eLife.32472 (2018).
- 741 10 Melamed, D., Young, D. L., Gamble, C. E., Miller, C. R. & Fields, S. Deep mutational  
742 scanning of an RRM domain of the *Saccharomyces cerevisiae* poly(A)-binding protein.  
743 *RNA* **19**, 1537-1551, doi:10.1261/rna.040709.113 (2013).
- 744 11 Olson, C. A., Wu, N. C. & Sun, R. A comprehensive biophysical description of pairwise  
745 epistasis throughout an entire protein domain. *Current biology : CB* **24**, 2643-2651,  
746 doi:10.1016/j.cub.2014.09.072 (2014).
- 747 12 Sahoo, A., Khare, S., Devanarayanan, S., Jain, P. C. & Varadarajan, R. Residue  
748 proximity information and protein model discrimination using saturation-suppressor  
749 mutagenesis. *eLife* **4**, 371, doi:10.7554/eLife.09532 (2015).
- 750 13 Li, C. & Zhang, J. Multi-environment fitness landscapes of a tRNA gene. *Nature Ecology*  
751 *& Evolution* **15**, 1, doi:10.1038/s41559-018-0549-8 (2018).
- 752 14 Li, C., Qian, W., Maclean, C. J. & Zhang, J. The fitness landscape of a tRNA gene.  
753 *Science* **352**, 837-840, doi:10.1126/science.aae0568 (2016).
- 754 15 Domingo, J., Diss, G. & Lehner, B. Pairwise and higher-order genetic interactions during  
755 the evolution of a tRNA. *Nature* **558**, 117-121, doi:10.1038/s41586-018-0170-7 (2018).
- 756 16 Puchta, O. *et al.* Network of epistatic interactions within a yeast snoRNA. *Science* **352**,  
757 840-844, doi:10.1126/science.aaf0965 (2016).
- 758 17 Göbel, U., Sander, C., Schneider, R. & Valencia, A. Correlated mutations and residue  
759 contacts in proteins. *Proteins* **18**, 309-317, doi:10.1002/prot.340180402 (1994).
- 760 18 Altschuh, D., Lesk, A. M., Bloomer, A. C. & Klug, A. Correlation of co-ordinated amino  
761 acid substitutions with function in viruses related to tobacco mosaic virus. *Journal of*  
762 *molecular biology* **193**, 693-707 (1987).
- 763 19 Gloor, G. B., Martin, L. C., Wahl, L. M. & Dunn, S. D. Mutual information in protein  
764 multiple sequence alignments reveals two classes of coevolving positions. *Biochemistry*  
765 **44**, 7156-7165, doi:10.1021/bi050293e (2005).
- 766 20 Halabi, N., Rivoire, O., Leibler, S. & Ranganathan, R. Protein sectors: evolutionary units  
767 of three-dimensional structure. *Cell* **138**, 774-786, doi:10.1016/j.cell.2009.07.038 (2009).
- 768 21 Lockless, S. W. & Ranganathan, R. Evolutionarily conserved pathways of energetic  
769 connectivity in protein families. *Science* **286**, 295-299 (1999).
- 770 22 Morcos, F. *et al.* Direct-coupling analysis of residue coevolution captures native contacts  
771 across many protein families. *Proceedings of the National Academy of Sciences* **108**,  
772 E1293-1301, doi:10.1073/pnas.1111471108 (2011).
- 773 23 Weigt, M., White, R. A., Szurmant, H., Hoch, J. A. & Hwa, T. Identification of direct  
774 residue contacts in protein-protein interaction by message passing. *Proceedings of the*  
775 *National Academy of Sciences* **106**, 67-72, doi:10.1073/pnas.0805923106 (2009).
- 776 24 Burger, L. & van Nimwegen, E. Disentangling Direct from Indirect Co-Evolution of  
777 Residues in Protein Alignments. *PLoS Computational Biology* **6**, e1000633,  
778 doi:10.1371/journal.pcbi.1000633 (2010).
- 779 25 Weinreb, C. *et al.* 3D RNA and Functional Interactions from Evolutionary Couplings. *Cell*  
780 **165**, 963-975, doi:10.1016/j.cell.2016.03.030 (2016).

- 781 26 Tóth-Petróczy, A. *et al.* Structured States of Disordered Proteins from Genomic  
782 Sequences. *Cell* **167**, 158-170.e112, doi:10.1016/j.cell.2016.09.010 (2016).
- 783 27 Hopf, T. A. *et al.* Three-Dimensional Structures of Membrane Proteins from Genomic  
784 Sequencing. *Cell* **149**, 1607-1621, doi:10.1016/j.cell.2012.04.012 (2012).
- 785 28 Marks, D. S. *et al.* Protein 3D structure computed from evolutionary sequence variation.  
786 *PLoS ONE* **6**, e28766, doi:10.1371/journal.pone.0028766 (2011).
- 787 29 Jones, D. T., Buchan, D. W. A., Cozzetto, D. & Pontil, M. PSICOV: precise structural  
788 contact prediction using sparse inverse covariance estimation on large multiple  
789 sequence alignments. *Bioinformatics* **28**, 184-190, doi:10.1093/bioinformatics/btr638  
790 (2012).
- 791 30 De Leonardis, E. *et al.* Direct-Coupling Analysis of nucleotide coevolution facilitates RNA  
792 secondary and tertiary structure prediction. *Nucleic Acids Research* **43**, 10444-10455,  
793 doi:10.1093/nar/gkv932 (2015).
- 794 31 Sułkowska, J. I., Morcos, F., Weigt, M., Hwa, T. & Onuchic, J. N. Genomics-aided  
795 structure prediction. *Proceedings of the National Academy of Sciences* **109**, 10340-  
796 10345, doi:10.1073/pnas.1207864109 (2012).
- 797 32 Ovchinnikov, S. *et al.* Protein structure determination using metagenome sequence data.  
798 *Science* **355**, 294-298, doi:10.1126/science.aah4043 (2017).
- 799 33 Ovchinnikov, S. *et al.* Large-scale determination of previously unsolved protein  
800 structures using evolutionary information. *eLife* **4**, e09248, doi:10.7554/eLife.09248  
801 (2015).
- 802 34 Ovchinnikov, S., Kamisetty, H. & Baker, D. Robust and accurate prediction of residue-  
803 residue interactions across protein interfaces using evolutionary information. *eLife* **3**,  
804 e02030, doi:10.7554/eLife.02030 (2014).
- 805 35 Matreyek, K. A. *et al.* Multiplex assessment of protein variant abundance by massively  
806 parallel sequencing. *Nature Genetics* **50**, 874-882, doi:10.1038/s41588-018-0122-z  
807 (2018).
- 808 36 Weile, J. *et al.* A framework for exhaustively mapping functional missense variants.  
809 *Molecular Systems Biology* **13**, 957, doi:10.15252/msb.20177908 (2017).
- 810 37 Rocklin, G. J. *et al.* Global analysis of protein folding using massively parallel design,  
811 synthesis, and testing. *Science* **357**, 168-175, doi:10.1126/science.aan0693 (2017).
- 812 38 Kim, I., Miller, C. R., Young, D. L. & Fields, S. High-throughput analysis of in vivo protein  
813 stability. *Molecular & cellular proteomics : MCP* **12**, 3370-3378,  
814 doi:10.1074/mcp.O113.031708 (2013).
- 815 39 Marks, D. S., Hopf, T. A. & Sander, C. Protein structure prediction from sequence  
816 variation. *Nature Biotechnology* **30**, 1072-1080, doi:10.1038/nbt.2419 (2012).
- 817 40 Costanzo, M. *et al.* The genetic landscape of a cell. *Science* **327**, 425-431,  
818 doi:10.1126/science.1180823 (2010).
- 819 41 Stein, R. R., Marks, D. S. & Sander, C. Inferring Pairwise Interactions from Biological  
820 Data Using Maximum-Entropy Probability Models. *PLoS Computational Biology* **11**,  
821 e1004182, doi:10.1371/journal.pcbi.1004182 (2015).
- 822 42 Andreani, J. & Söding, J. bbcontacts: prediction of  $\beta$ -strand pairing from direct coupling  
823 patterns. *Bioinformatics* **31**, 1729-1737, doi:10.1093/bioinformatics/btv041 (2015).

- 824 43 Jones, D. T. Protein secondary structure prediction based on position-specific scoring  
825 matrices. *Journal of molecular biology* **292**, 195-202, doi:10.1006/jmbi.1999.3091  
826 (1999).
- 827 44 The PyMOL Molecular Graphics System, V. S., LLC. *The PyMOL Molecular Graphics*  
828 *System, Version 1.8* (2015).
- 829 45 Schwieters, C. D., Kuszewski, J. J., Tjandra, N. & Clore, G. M. The Xplor-NIH NMR  
830 molecular structure determination package. *Journal of magnetic resonance (San Diego,*  
831 *Calif. : 1997)* **160**, 65-73 (2003).
- 832 46 Araya, C. L. *et al.* A fundamental protein property, thermodynamic stability, revealed  
833 solely from large-scale measurements of protein function. *Proceedings of the National*  
834 *Academy of Sciences* **109**, 16858-16863, doi:10.1073/pnas.1209751109 (2012).
- 835 47 Liu, Y., Palmedo, P., Ye, Q., Berger, B. & Peng, J. Enhancing Evolutionary Couplings  
836 with Deep Convolutional Neural Networks. *Cell systems*, doi:10.1016/j.cels.2017.11.014  
837 (2017).
- 838 48 Schaarschmidt, J., Monastyrskyy, B., Kryshtafovych, A. & Bonvin, A. M. J. J.  
839 Assessment of contact predictions in CASP12: Co-evolution and deep learning coming  
840 of age. *Proteins* **86 Suppl 1**, 51-66, doi:10.1002/prot.25407 (2018).
- 841 49 Fox, N. K., Brenner, S. E. & Chandonia, J. M. SCOPe: Structural Classification of  
842 Proteins--extended, integrating SCOP and ASTRAL data and classification of new  
843 structures. *Nucleic Acids Res* **42**, D304-309, doi:10.1093/nar/gkt1240 (2014).
- 844 50 Starita, L. M. *et al.* Massively Parallel Functional Analysis of BRCA1 RING Domain  
845 Variants. *Genetics* **200**, 413-422, doi:10.1534/genetics.115.175802 (2015).
- 846 51 Starita, L. M. *et al.* Activity-enhancing mutations in an E3 ubiquitin ligase identified by  
847 high-throughput mutagenesis. *Proceedings of the National Academy of Sciences* **110**,  
848 E1263-1272, doi:10.1073/pnas.1303309110 (2013).
- 849 52 Poelwijk, F. J., Socolich, M. & Ranganathan, R. Learning the pattern of epistasis linking  
850 genotype and phenotype in a protein. *bioRxiv*, doi:10.1101/213835 (2017).
- 851 53 Starr, T. N., Picton, L. K. & Thornton, J. W. Alternative evolutionary histories in the  
852 sequence space of an ancient protein. *Nature* **5**, e16965, doi:10.1038/nature23902  
853 (2017).
- 854 54 Fowler, D. M. *et al.* High-resolution mapping of protein sequence-function relationships.  
855 *Nature Methods* **7**, 741-746, doi:10.1038/nmeth.1492 (2010).
- 856 55 McLaughlin Jr, R. N., Poelwijk, F. J., Raman, A., Gosal, W. S. & Ranganathan, R. The  
857 spatial architecture of protein function and adaptation. *Nature* **491**, 138-142,  
858 doi:10.1038/nature11500 (2012).
- 859 56 Otwinowski, J. Biophysical inference of epistasis and the effects of mutations on protein  
860 stability and function. *Mol Biol Evol*, doi:10.1093/molbev/msy141 (2018).
- 861 57 Alexander, P. A., He, Y., Chen, Y., Orban, J. & Bryan, P. N. A minimal sequence code  
862 for switching protein structure and function. *Proceedings of the National Academy of*  
863 *Sciences* **106**, 21149-21154, doi:10.1073/pnas.0906408106 (2009).
- 864 58 Salinas, V. H. & Ranganathan, R. Coevolution-based inference of amino acid  
865 interactions underlying protein function. *eLife* **7**, 41, doi:10.7554/eLife.34300 (2018).

- 866 59 Jones, D. T., Singh, T., Kosciolk, T. & Tetchner, S. MetaPSICOV: combining  
867 coevolution methods for accurate prediction of contacts and long range hydrogen  
868 bonding in proteins. *Bioinformatics* **31**, 999-1006, doi:10.1093/bioinformatics/btu791  
869 (2015).
- 870 60 Wang, S., Sun, S., Li, Z., Zhang, R. & Xu, J. Accurate De Novo Prediction of Protein  
871 Contact Map by Ultra-Deep Learning Model. *PLoS Computational Biology* **13**, e1005324,  
872 doi:10.1371/journal.pcbi.1005324 (2017).
- 873 61 Rohl, C. A., Strauss, C. E., Misura, K. M. & Baker, D. Protein structure prediction using  
874 Rosetta. *Methods Enzymol* **383**, 66-93, doi:10.1016/S0076-6879(04)83004-0 (2004).
- 875 62 Yang, J. *et al.* The I-TASSER Suite: protein structure and function prediction. *Nat*  
876 *Methods* **12**, 7-8, doi:10.1038/nmeth.3213 (2015).
- 877 63 Rollins, N. J. *et al.* 3D protein structure from genetic epistasis experiments.  
878 doi:10.1101/320721 (2018).
- 879 64 Rubin, A. F. *et al.* A statistical framework for analyzing deep mutational scanning data.  
880 *Genome biology* **18**, 741, doi:10.1186/s13059-017-1272-5 (2017).
- 881 65 Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.  
882 *Bioinformatics* **26**, 2460-2461, doi:10.1093/bioinformatics/btq461 (2010).
- 883 66 Barlow, R. *Statistics : a guide to the use of statistical methods in the physical sciences.*  
884 (Wiley, 1989).
- 885 67 Schäfer, J. & Strimmer, K. A shrinkage approach to large-scale covariance matrix  
886 estimation and implications for functional genomics. *Statistical applications in genetics*  
887 *and molecular biology* **4**, Article32, doi:10.2202/1544-6115.1175 (2005).
- 888 68 Gallagher, T., Alexander, P., Bryan, P. & Gilliland, G. L. Two crystal structures of the B1  
889 immunoglobulin-binding domain of streptococcal protein G and comparison with NMR.  
890 *Biochemistry* **33**, 4721-4729 (1994).
- 891 69 Pires, J. R. *et al.* Solution structures of the YAP65 WW domain and the variant L30 K in  
892 complex with the peptides GTPPPYTVG, N-(n-octyl)-GPPPY and PLPPY and the  
893 application of peptide libraries reveal a minimal binding epitope. *Journal of molecular*  
894 *biology* **314**, 1147-1156, doi:10.1006/jmbi.2000.5199 (2001).
- 895 70 Deo, R. C., Bonanno, J. B., Sonenberg, N. & Burley, S. K. Recognition of polyadenylate  
896 RNA by the poly(A)-binding protein. *Cell* **98**, 835-845 (1999).
- 897 71 Glover, J. N. & Harrison, S. C. Crystal structure of the heterodimeric bZIP transcription  
898 factor c-Fos-c-Jun bound to DNA. *Nature* **373**, 257-261, doi:10.1038/373257a0 (1995).
- 899 72 Adhikari, B., Bhattacharya, D., Cao, R. & Cheng, J. CONFOLD: Residue-residue  
900 contact-guided ab initio protein folding. *Proteins* **83**, 1436-1449, doi:10.1002/prot.24829  
901 (2015).
- 902 73 Zhang, Y. & Skolnick, J. Scoring function for automated assessment of protein structure  
903 template quality. *Proteins* **57**, 702-710, doi:10.1002/prot.20264 (2004).
- 904 74 Seemayer, S., Gruber, M. & Söding, J. CCMpred--fast and precise prediction of protein  
905 residue-residue contacts from correlated mutations. *Bioinformatics* **30**, 3128-3130,  
906 doi:10.1093/bioinformatics/btu500 (2014).

907  
908

## 909 Methods:

### 910 Datasets and preprocessing

#### 911 Protein G B1 domain

912 Protein G B1 domain (GB1) deep mutational scanning data was obtained from the  
913 supplementary information of Olson, et al. <sup>11</sup>. The data consists of summed read counts of three  
914 replicate experiments assaying the binding affinity of GB1 variants to immunoglobulin G (IgG).

915 Read frequencies of each single or double mutant variant in input library and output library (after  
916 binding affinity assay) were calculated as variant read counts relative to wild-type variant read  
917 counts. A variant's fitness was calculated as the natural logarithm of the ratio of output to input  
918 read frequency, i.e.  $f_i = \log\left(\frac{n_i^{out}/n_{wt}^{out}}{n_i^{in}/n_{wt}^{in}}\right)$ , with  $n$  as read counts, superscripts denoting input or  
919 output sequencing library and subscripts denoting variant  $i$  or wild-type variant.

920 The standard error of fitness estimates was calculated from read counts under Poissonian  
921 assumptions, i.e.  $\sigma_i = \sqrt{\frac{1}{n_i^{in}} + \frac{1}{n_i^{out}} + \frac{1}{n_{wt}^{in}} + \frac{1}{n_{wt}^{out}}}$  <sup>64</sup>. We note that this is a lower bound estimate of  
922 the actual error, due to the lack of replicate information.

923 Each measurement assay has a lower measurement limit due to unspecific background effects  
924 (Extended Data Figure 2A). In the case of the IgG-binding assay for GB1, this is presumably  
925 mainly due to unspecific carry-over on beads <sup>11</sup>. The fitness values derived from the  
926 measurement are therefore a convolution of the actual binding affinities to IgG and nonspecific  
927 carryover, i.e.  $\exp(f_i^{measured}) = \exp(f_i^{binding}) + \exp(f^{carryover})$ . Fitness values of variants  
928 close to the lower measurement limit of the assay are therefore dominated by unspecific  
929 carryover effects. The lower measurement limit of the assay was estimated by two approaches  
930 that yielded identical estimates. One, from a kernel density estimate of the single mutant fitness  
931 distribution (R function *density* with parameter *bw* set to 0.15), where the position of the lower  
932 mode of the data corresponded to  $f^{carryover} = -5.85$  (~0.3% on linear scale). Two, from  
933 examining the fitness distribution of double mutants with expected fitness lower than -8 log-  
934 units, i.e. double mutants resulting from two lethal or nearly lethal single mutant variants, whose  
935 fitness values are thus expected to be dominated by background effects. The median of this  
936 background fitness distribution yielded an estimate of  $f^{carryover} = -5.85$ .

937 7% of double mutant variants were discarded due to too low sequencing coverage in input or  
938 output libraries (Extended Data Figure 2B). That is, variants that had less than 200 input reads  
939 and no output reads were discarded, because it is not possible to determine their fitness. Above  
940 200 input reads, variants without output reads are certain to be dominated by nonspecific  
941 carryover effects. These variants were retained and their fitness was calculated by setting their  
942 output read count to 0.5.

## 943 GB1 down-sampling

944 Down-sampling of the full GB1 dataset was performed in three different ways.

945 For the 'doped' datasets, we only allowed amino acid changes created by one nucleotide  
946 mutation from the wild-type sequence (ENA entry M12825). To down-sample the sequencing  
947 read coverage, for each variant we picked as a down-sampled read count the draw of  
948 successes from a binomial distribution with the number of sequencing reads in the full datasets  
949 as trials and the target down-sampling rate (25%, 10% or 2.5%) as chance of success. For the  
950 read down-sampled and doped datasets (and combinations of both), the analysis workflow for  
951 the full dataset was repeated.

952 For the down-sampled datasets taking only positive or negative epistatic information into  
953 account, we calculated *epistasis* and *association scores* from epistatic enrichment matrices and  
954 partial correlation matrices of only positive or negative epistasis information. Instead of merging  
955 positive and negative matrices and then calculating z-scores, we only calculated z-scores with  
956 the individual errors from positive or negative epistasis information. The *combined scores* (for  
957 which results are reported) for each set was then calculated as for the full dataset by summing  
958 standardized *epistasis* and *association scores*.

## 959 hYAP WW domain

960 hYAP WW domain data was obtained from Sequence Read Archive (SRA) entry SRP015751<sup>46</sup>.  
961 Paired-end reads were merged with USearch<sup>65</sup> and merged reads with any base having a  
962 Phred base quality score below 20 were discarded. Read counts from the two technical  
963 sequencing replicates were merged and read counts for the same amino acid variants with at  
964 most one synonymous mutation in another codon were summed up. The dataset consists of an  
965 input library and three output libraries after consecutive rounds of selection in a phage display  
966 assay. Fitness was estimated as the slope of log frequency (variant counts divided by wild-type  
967 counts) changes over the rounds of selection experiment<sup>46</sup>. For each variant at each selection

968 step a Poissonian error of  $\sigma_{i,x} = \sqrt{\frac{1}{n_i^x} + \frac{1}{n_{wt}^x}}$  was calculated, with x denoting the selection step.

969 Slopes were calculated as weighted straight line least square fits<sup>66</sup>. Comparison of library-wide  
970 changes in variant frequencies between selection rounds suggested differential selection  
971 pressures across the rounds. We thus applied a non-equidistant spacing of 0.6, 1.17 and 1.22  
972 between selection rounds when calculating slopes. Only variants that have more than 10 reads  
973 in the input library and at least one read after the first selection were retained for further analysis  
974 (45% of constructed double mutants). The lower fitness limit was calculated as the weighted  
975 mean fitness of all variants containing STOP codons (-0.78 in log-fitness units).

## 976 Pab1 RRM2 domain

977 Pab1 RRM2 domain data was obtained from the Supplementary Table 5 of Melamed, et al.<sup>10</sup>.  
978 Reported enrichment scores were log-transformed to obtain fitness values. Output reads per  
979 variant were deduced from the number of input reads times the enrichment score and used to  
980 calculate a Poissonian error of the fitness estimate. Single mutant count data is not provided  
981 and we thus estimated the error of single mutant fitness estimates to be 0.01. Lower bound of  
982 fitness assay was estimated as weighted mean fitness of all double mutant variants containing  
983 STOP codons (-3.1 log-fitness units).

## 984 FOS-JUN interaction

985 Raw count tables were provided by Guillaume Diss<sup>9</sup>. The dataset consists of input and output  
986 sequencing libraries after selection for physical interaction between the two proteins in a protein  
987 complementation assay in three biological replicates. Per sequencing library, read counts from  
988 all synonymous variants were summed up. Only variants that had more than 10 reads in each of  
989 the three input libraries were used for further analysis (43% of double mutants). Per input/output  
990 replicate, fitness of each variant was calculated as the log change in frequency compared to the  
991 wild-type variant (as for GB1). A Poissonian error for each variant's fitness estimate was  
992 derived.

993 The dataset has a large dynamic range, thus many low-fitness variants with low input read  
994 coverage have very low or no output read counts (per replicate ~1/3 of variants have below 3  
995 output counts, ~15% of variants have zero output counts), effectively reducing the dynamic  
996 range of the assay for low input variants and distorting the estimate of the overall fitness  
997 distribution (see Extended Data Figure 5E). To overcome this, we implemented a Bayesian  
998 estimator of fitness. For each double mutant variant, we first identified the 1000 nearest



999 neighbors in single mutant fitness space (i.e. those double mutants whose respective single  
1000 mutant fitness values are similar to the single mutant fitness values of the variant under  
1001 consideration) with sufficient input coverage (more than 100 reads in the input library). From this  
1002 set of 1000 nearest neighbors we calculated the expected distribution of double mutant fitness  
1003 values, which served as a prior distribution. For the variant under consideration we calculated  
1004 the likelihood distribution of fitness values given its input and output read counts under  
1005 Poissonian assumptions. Fitness was then estimated as the mean of the distribution resulting  
1006 from the multiplication of prior and likelihood distributions. Error of fitness estimate was  
1007 estimated as the standard deviation of the resulting distribution. Estimated fitness from the three  
1008 replicate experiments were merged by weighted averaging.

1009

## 1010 Epistasis classification

1011 Epistasis was calculated from a non-parametric null model (Figure 1B) in order to account for  
1012 nonlinearities close to the lower limit of the fitness assay measurement range, non-specific  
1013 epistatic behavior resulting from e.g. thermodynamic stability thresholds as well as differential  
1014 uncertainty of fitness measurements across the fitness landscape, due to lower read counts in  
1015 the output for low fitness variants.

1016 First, double mutant fitness values were corrected by subtracting the average local fitness  
1017 computed using a two-dimensional local polynomial regression (using R function *loess* with  
1018 span = 0.2). This was necessary to avoid boundary effects of quantile-based fits in boundary  
1019 regions with non-zero slopes. 5th and 95th percentile surfaces were then fit to these corrected  
1020 double mutant fitness values, by computing for each double mutant variant the 5th and 95th  
1021 percentile of the fitness distribution made up of the 1% closest neighbors in single mutant  
1022 fitness space. Double mutant variants with fitness values below the 5th or above the 95th  
1023 percentile were categorized as negative or positive epistatic, respectively (Figure 1B).

1024 The evaluation of positive or negative epistasis was, however, restricted to specific subsets of  
1025 the data where measurement errors do not impede epistasis classification (Extended Data  
1026 Figure 2C). The data subset deemed suitable for positive epistasis classification is limited to  
1027 regions where

- 1028 • the 95th percentile fitness surface is below wild-type fitness

1029       • at least one single mutant fitness value is significantly smaller than wild-type fitness at  
1030           two standard errors

1031       • the expected fitness (sum of both single mutant fitness values) is not significantly higher  
1032           than wild-type at two standard errors

1033 The rationale for these criteria is to avoid double mutants from two neutral single mutants,  
1034 because these are dominated by measurement noise of overabundant wild-type like variants.  
1035 No restrictions were instead applied to the lower limits of the measurement range, because  
1036 otherwise no/little epistasis quantification would have been available for several positions with  
1037 very strong detrimental effects as well as because strong positive epistatic effects are observed  
1038 in these regions, despite the dominance of background measurement effects.

1039 The data subset in which variants were potentially classified as negative epistatic is limited to  
1040 data regions where

1041       • the 5th percentile fitness surface is above the 95th percentile of the background effect  
1042           distribution; this value is derived from the 95th percentile of double mutant fitness values  
1043           with expected fitness below -8 (analogous to lower fitness limit estimation, see above).

1044       • both single mutant fitness values are significantly higher than the lower limit of the  
1045           fitness assay measurement range at two standard errors

1046       • the expected fitness (sum of both single mutant fitness values) is not significantly higher  
1047           than wild-type at two standard errors

1048 The rationale for criteria 1 and 2 is to avoid background measurement effects that make  
1049 negative epistasis quantification unreliable.

1050 As a result of these restrictions as well as differences in initial coverage, the number of double  
1051 mutant variants that can be used to assess positive and negative epistasis varies substantially  
1052 across position pairs and datasets (see Table 1, Extended Data Figures 2D-F and 5A&D).

1053

## 1054 Epistatic interactions (*epistasis scores*)

1055 We derived several interaction scores to estimate which position pairs are in close contact in the  
1056 tertiary structure (see Figures 2A and 3A,B and Extended Data Figure 1). These scores are  
1057 based on summarizing epistasis information on the position pair-level and accounting for the

1058 uncertainty inherent in the summarized estimates due to differential error of fitness estimates  
1059 across the measurement range as well as varying numbers of double mutants amenable to  
1060 epistasis classification (see Table 1, Extended Data Figures 2D-F and 5A&D).

1061 To summarize epistasis information on the position pair-level, we calculated the fraction of  
1062 positive or negative epistatic variants per position pair. The fraction of positive epistatic variants  
1063 per position pair is the number of positive epistatic variants divided by the number of all variants  
1064 that lie in the double mutant space amenable to positive epistasis classification (Extended Data  
1065 Figure 1, step 5b, equivalent calculation for negative epistasis fraction). Because enrichments  
1066 with positive and negative epistatic variants per position are anti-correlated (Extended Data  
1067 Figure 3A), we treated both separately and only aggregated them to derive the final interaction  
1068 scores.

1069 To estimate the uncertainty in epistatic fractions per position pair for downstream analyses we  
1070 implemented a re-sampling approach (Extended Data Figure 1, step 5, described here for  
1071 positive epistatic variants, but equivalent for negative epistatic variants). In each of 10,000 re-  
1072 sampling runs:

- 1073 • each variant's fitness was drawn from a normal distribution with the measured fitness as  
1074 mean and the uncertainty due to sequencing coverage as standard deviation  $f_{ij}^{sampled} =$   
1075  $\mathcal{N}(f_{ij}, \sqrt{\sigma_{ij}^2 + s_i^2 * \sigma_i^2 + s_j^2 * \sigma_j^2})$ , with  $s_x$  as the local slope of the median fitness  
1076 landscape in direction of the respective single mutant (step 5a)
- 1077 • positive epistasis of variants was re-classified given the drawn fitness values (also step  
1078 5a)
- 1079 • each position pair's fraction of positive epistatic variants was drawn from the posterior  
1080 probability distribution of how likely an underlying true fraction of epistatic variants  $E_{xy}^+$  is  
1081 to generate the observed fraction of epistatic variants given the finite number of overall  
1082 variants, i.e.  $e_{xy}^+ \sim p(E_{xy}^+ | \# \varepsilon_{xy}^+, \# \text{ variants}_{xy})$  (step 5b). The posterior probability  
1083 distribution is the product of a prior probability distribution – the kernel density estimate  
1084 of the expected epistatic fractions across all position pairs (calculated using R function  
1085 *density* with parameter *bw* set to 0.05) – and the likelihood function for the underlying  
1086 true fraction of epistatic variants given the observed fraction of epistatic variants and the  
1087 overall number of variants under binomial sampling assumptions

1088 To derive an interaction score from the epistatic fractions per position pair, mean positive and  
1089 negative epistatic fractions across resampling runs were combined by weighted averaging, with

1090 weights as the inverse variances of epistatic fractions across resampling runs, i.e.  $\langle e_{xy} \rangle =$   
1091 
$$\frac{\langle e_{xy}^+ \rangle * \sigma_{e_{xy}^+}^{-2} + \langle e_{xy}^- \rangle * \sigma_{e_{xy}^-}^{-2}}{\sigma_{e_{xy}^+}^{-2} + \sigma_{e_{xy}^-}^{-2}}.$$

1092 To arrive at the final *epistasis score*, the mean epistatic fractions were further normalized by  
1093 their combined uncertainty,  $E_{xy} = \langle e_{xy} \rangle / \sigma_{xy}$ , with  $\sigma_{xy} = (\sigma_{e_{xy}^+}^{-2} + \sigma_{e_{xy}^-}^{-2})^{-1/2}$  (step 6).

1094

## 1095 Epistasis pattern correlations (*association scores*)

1096 In addition to the *epistasis score* we derived an interaction score from the partial correlation of  
1097 epistasis patterns between position pairs, termed *association score*. The rationale behind this  
1098 score is that proximal positions in the protein should have similar distances and geometrical  
1099 arrangements towards all other positions in the protein and should therefore also have similar  
1100 patterns of epistatic interactions with all other positions.

1101 In each re-sampling run we constructed a symmetric matrix of the drawn positive epistatic  
1102 fractions  $e_{xy}^+$  with number of rows and columns as the number of mutated positions (Extended  
1103 Data Figure 1, step 5c, equivalent for negative epistatic fractions). Missing values (positions  
1104 pairs without observed variants) were imputed by drawing a random value from the overall  
1105 distribution of epistatic fractions. A pseudo count equal to the first quartile of the epistatic  
1106 fraction distribution was added to each epistatic fraction. Diagonal elements (epistatic fractions  
1107 of a position with itself) were set to 1. The matrix values were transformed by the natural  
1108 logarithm (to make distributions more symmetric, thus correlations are not dominated by few  
1109 position pairs with large epistatic fractions) and for each pair of columns the Pearson correlation  
1110 coefficient was calculated to arrive at the correlation matrix (step 5d).

1111 A shrinkage approach was used to improve the estimate of the correlation matrix<sup>67</sup>. In brief, the  
1112 empirical correlation matrix is shrunk towards the identity matrix in order to minimize the mean-  
1113 squared error between estimated and true correlation matrix. Additionally, this yields a positive  
1114 definite and well-conditioned correlation matrix, suitable for inversion. All computations on  
1115 correlation matrices, shrinkage and matrix inversion were performed with the R package  
1116 *corpcor*, functions *cor.shrink* and *pcor.shrink*<sup>67</sup>.

1117 Partial correlations of epistatic patterns between each position pair were calculated by inverting  
1118 the shrunk correlation matrix and normalizing each off-diagonal entry of the inverted matrix by

1119 the geometric mean of the two respective diagonal entries, i.e.  $a_{xy}^+ = \frac{r_{xy}^{-1}}{\sqrt{r_{xx}^{-1} * r_{yy}^{-1}}}$ , with  $r_{xy}^{-1}$  as the  
1120 (x,y)-entry of the inverted correlation matrix (Extended Data Figure 1, step 5d). Equivalent to the  
1121 *epistasis score*, positive and negative partial correlation estimates were merged by calculating  
1122 weighted averages of their mean estimates across re-sampling runs, with weights as the inverse  
1123 variances across resampling runs, i.e.  $\langle a_{xy} \rangle = \frac{\langle a_{xy}^+ \rangle * \sigma_{a_{xy}^+}^{-2} + \langle a_{xy}^- \rangle * \sigma_{a_{xy}^-}^{-2}}{\sigma_{a_{xy}^+}^{-2} + \sigma_{a_{xy}^-}^{-2}}$ , and the final *association*  
1124 *score* normalized by the combined uncertainty,  $A_{xy} = \langle a_{xy} \rangle / \sigma_{xy}$ , with  $\sigma_{xy} = (\sigma_{a_{xy}^+}^{-2} +$   
1125  $\sigma_{a_{xy}^-}^{-2})^{-1/2}$  (step 6).

1126

## 1127 Aggregating *epistasis* and *association scores* (combined scores)

1128 We further derived a *combined score* by summing the standardized *epistasis* and *association*  
1129 *scores*, i.e.  $C_{xy} = \frac{E_{xy} - \langle E \rangle}{\sigma_E} + \frac{A_{xy} - \langle A \rangle}{\sigma_A}$ . We note that this is a naïve approach to combining the  
1130 information from these two complementary sources, and surely more sophisticated approaches  
1131 that further improve proximity estimates can be developed.

1132

## 1133 Secondary structure prediction

1134 We used a two-dimensional kernel smoothing approach to predict secondary structure elements  
1135 from interaction score matrices (Figure 4A-C). For a given position along the linear chain (on the  
1136 diagonal of the interaction score matrix), interaction scores (off the diagonal) are integrated with  
1137 distance-specific weighting according to the kernel, which reflects the known geometry of  
1138 secondary structures.

1139 The alpha kernel takes on a sinusoidal profile perpendicular to the diagonal to weight  
1140 interactions according to whether the position pair considered should have congruent side-chain  
1141 orientations (see diagonal and perpendicular profiles in Figure 4B). The kernel was defined as

1142  $k_\alpha(d, p) = \left( \cos\left(p * \frac{2\pi}{3.6}\right) + 1/3 \right) * e^{-\frac{d^2}{c^2}}$ , with  $d = |2x - i - j|$  as the diagonal distance of the  
1143 interaction  $ij$  (off the diagonal) to the reference position  $x$  (on the diagonal) and  $p = |i - j|$  as the  
1144 perpendicular distance of the interaction off the diagonal. The kernel weight for positions with  $p$

1145 > 5 was set to 0, thus only including interactions across little more than the first helical turn.  
1146 Finally,  $c = 4$  is the integration scale for the Gaussian kernel along the diagonal. While smaller  
1147 integration scales do yield noisier results and longer integration scales can lead to non-detection  
1148 of shorter secondary structure stretches, we found that in practice our whole approach  
1149 (including the actual detection algorithm described below) is robust to alterations of the  
1150 integration length.

1151 The kernel smoothed alpha value for a given position  $x$  on the diagonal is then calculated as the  
1152 sum over all interaction scores times their kernel weights  $K_{\alpha,x} = \sum_i \sum_j k_{\alpha}(d,p) * S_{ij}$ , where  $S_{ij}$  is  
1153 one of the interaction scores (*epistasis*, *association* or *combined score*) at position pair  $ij$ .

1154 The beta kernel takes an alternating profile perpendicular to the diagonal to weight interactions  
1155 according to alternating side-chain orientations in a beta strand and was defined as  $k_{\beta}(d,p) =$   
1156  $\left( (p + 1) \bmod 2 - 1/3 \right) * e^{-\frac{d^2}{c^2}}$ , with  $c = 4$ . Only interactions with perpendicular distances equal  
1157 or smaller than two (i.e.  $k_{\beta}(d,p > 3) = 0$ ) were included.

1158 To calculate whether kernel-weighted interaction scores of a specific position are larger than  
1159 expected, they were compared to kernel-weighted scores obtained from  $10^4$  randomized control  
1160 datasets. Randomization was performed by shuffling all interaction scores, while preserving  
1161 matrix symmetry, and kernel-weighted interaction scores from randomized control datasets  
1162 were calculated for each position independently to control for possible boundary effects in  
1163 positions close to the borders of the protein chain. A p-value for each position was calculated as  
1164 the fraction of random controls with kernel smoothed values above that of the real data.

1165 Secondary structure elements were identified by searching for continuous stretches of positions  
1166 with high propensities to belong to either alpha helices or beta strands. The following workflow  
1167 was implemented:

- 1168 1. calculate a combined p-value for seeds of length 3 by combining position-wise p-values  
1169 using Fisher's method for both alpha and beta kernel smoothed interaction scores
- 1170 2. separate positions according to whether combined p-values of seeds from alpha or beta  
1171 kernels are more significant, i.e.
  - 1172 2.1. for alpha helical propensity calculations only consider stretches of at least 5 consecutive  
1173 positions for which the combined p-value of seeds for alpha kernel smoothing is smaller  
1174 than that from beta kernel smoothing (thus setting the lower size limit of alpha helical  
1175 elements to five)

1176 2.2. for beta strand propensity calculations only consider stretches of at least 3 consecutive  
1177 positions for which the combined p-value of seeds for beta kernel smoothing is smaller  
1178 than that from alpha kernel smoothing (thus setting the lower size limit of beta strands  
1179 to three)

1180 For alpha helices and beta strands separately and while combined p-values of seeds < 0.05

- 1181 3. select the most significant seed
- 1182 4. test whether extension to any side yields lower combined p-value
  - 1183 4.1. if yes: extend seed in this direction and repeat step 4
  - 1184 4.2. else: repeat step 4 once to see whether further extension in same direction yields lower  
1185 combined p-value
    - 1186 4.2.1. if yes: extend and repeat step 4
    - 1187 4.2.2. else: proceed to step 5
- 1188 5. fix as secondary structure element and delete all 'used' p-values (and combined seed p-  
1189 values), such that other elements cannot incorporate them
- 1190 6. check whether other already fixed elements are adjacent or at most one position away
  - 1191 6.1. if yes: merge both elements
- 1192 7. repeat steps 3-6 until no more seeds with combined p-value < 0.05 are left

1193 This yields a list of predicted locations of secondary structure elements. We note that the  
1194 secondary structure elements predicted from deep mutational scanning data could be compared  
1195 to and combined with predictions derived from other tools, such as PSIPRED (Jones, 1999), to  
1196 further improve reliability.

1197 To detect beta sheet interactions a modified beta strand kernel was used. In contrast to beta  
1198 strand detection, the beta sheet interaction kernel is centered on each off-diagonal position. For  
1199 beta sheet kernels diagonal and perpendicular distances are therefore modified as  $d =$   
1200  $|x + y - i - j|$  and  $p = |x - i - (y - j)|$ . The kernels to detect parallel and anti-parallel beta  
1201 sheets differ in which is their 'diagonal' direction, i.e. the direction at which consecutive position  
1202 pairs interact in the beta sheet (Extended Data Figure 4A). Therefore, parameters  $d$  and  $p$  were  
1203 swapped for the anti-parallel beta sheet kernel. Also, because these positions can be deemed  
1204 the most crucial for deciding whether a position participates in a beta sheet interaction or not,  
1205 we up-weighted these positions (those with  $p = 0$ ) in the kernel by a factor of two, i.e.  $K_{\beta}(d, 0) =$

1206  $4/3 * e^{-\frac{d^2}{c^2}}$ .

1207 Beta sheet interactions were identified by searching for the most significant stretches of parallel  
1208 and anti-parallel interactions (similar to workflow for alpha helices and beta strands), then  
1209 identifying the set of most significant interactions that is consistent with previously predicted  
1210 secondary structure elements.

1211 In particular, step 1 & 3-7 from the above-described workflow were performed for the parallel  
1212 beta sheet kernel on each sub-diagonal (parallel to the main diagonal) of the interaction score  
1213 matrix separately; and for the anti-parallel beta sheet kernel on each perpendicular diagonal of  
1214 the interaction score matrix separately.

1215 The steps were modified as follows:

- 1216 ● for anti-parallel beta sheet interactions, only positions with a distance greater than 1 from  
1217 the main diagonal were used to calculate seed p-values (assuming anti-parallel beta  
1218 sheet interactions need a turn of at least length two to be connected)
- 1219 ● for parallel beta sheet interactions, only sub-diagonals with a distance greater than 4  
1220 from the main diagonal were considered (assuming parallel beta sheet interactions of  
1221 two adjacent beta strands need a linker region)

1222 We extended the workflow with the following steps to predict beta sheet interactions within the  
1223 protein domain:

- 1224 8. compute association of seeds with known beta strands (e.g. seed positions overlap strand 1  
1225 on one side and coincide with strand 3 on the other side)
- 1226 9. while there are seeds with  $p < 0.05$ : pick most significant seed from either the parallel or  
1227 anti-parallel sheet subset
- 1228 10. check consistency with secondary structure elements
  - 1229 10.1. discard the seed and jump back to step 9 if
    - 1230 10.1.1. it is overlapping or too close to an alpha helix or the linker region between two  
1231 beta strands that interact (minimal distance smaller one)
    - 1232 10.1.2. at least one of the two strands it is associated with already has two other beta  
1233 sheet interactions or the total number of beta sheet interactions exceeds  $2 * (\# \text{beta}$   
1234  $\text{strands} - 1)$
  - 1235 10.2. modify secondary structure elements and start anew from step 3 if
    - 1236 10.2.1. one side of the seed is not associated to a known beta strand: create this beta  
1237 strand
    - 1238 10.2.2. if both sides of the seed are associated with the same known strand: split the  
1239 strand and create a linker region in-between the strands



1240 10.3. else fix the beta sheet interaction and delete all other interactions that are  
1241 associated with the same strands and haven't been fixed yet, jump back to step 9

1242 11. if no more seeds with  $p < 0.001$ , finish

1243 12. update beta strands: keep only those positions that are part of a beta sheet interaction

1244

1245 For beta sheet pairing detection in the GB1 domain (as reported in Figure 4D and Extended  
1246 Data Figures 4A-C and 7B) we used as input the secondary structure element predictions  
1247 derived from the deep mutational scanning data (as shown in Figure 4A-C and Extended Data  
1248 Figure 4D). For the RRM and WW domains, we used as input PSIPRED predicted secondary  
1249 structure elements, due to the insufficient signal from secondary structure element predictions  
1250 from deep mutational scanning data.

1251

## 1252 Protein distance metrics

1253 The minimal side chain heavy atom distance, i.e. the minimal distance between any two side-  
1254 chain heavy atoms of a position pair (in case of glycine,  $C\alpha$ ), was used as the general distance  
1255 measure. A direct contact was defined as minimal side-chain heavy atom distance  $< 8\text{\AA}$ . For all  
1256 evaluations of predicted contact precision we only considered position pairs with non-trivial  
1257 tertiary contacts (those with a linear sequence separation of greater than 5 positions).

1258 We do find that, while using all heavy atoms to calculate distances increase the true positive  
1259 rates of predicted contacts by about 10%, side-chain heavy atom distances display much higher  
1260 true positive rates over random expectation, thus suggesting that side-chain interactions are  
1261 more informative for epistatic interactions (Extended Data Figure 7C).

1262 The Floyd-Warshall algorithm (implemented as custom script in R) was used to calculate the  
1263 minimal number of edges  $< 8\text{\AA}$  that connect any two positions in the protein.

1264 Reference structures used as comparison were

- 1265 • GB1 domain: PDB entry 1pga, X-ray diffraction structure<sup>68</sup>
- 1266 • WW domain: PDB entry 1k9q, solution NMR structure<sup>69</sup>
- 1267 • RRM domain: PDB entry 1cvj (chain A), X-ray diffraction structure of human Pab1<sup>70</sup>;  
1268 note that the central section of the yeast RRM domain analysed is one nucleotide longer  
1269 than the corresponding homologous region in the human RRM domain. We thus

1270 arbitrarily removed position 14 (in the loop region) when comparing the DMS-derived  
1271 predictions to the human Pab1 structure.

1272 • FOS-JUN interaction: PDB entry 1fos (chains E and F), X-ray diffraction structure <sup>71</sup>

1273 We found that precision or accuracy calculated against other reference structures differed only  
1274 marginally, thus we have limited reporting to the aforementioned PDB entries.

1275

## 1276 Protein folding

1277 To *ab initio* determine protein structures, we performed simulated annealing using the XPLOR-  
1278 NIH modeling suite <sup>45</sup> with structural restraints derived from the deep mutational scanning data.  
1279 Simulations were performed in three stages, in each of which 500 structural models were  
1280 generated. Stages 1 and 2 served to identify inconsistencies among defined structural  
1281 restraints. Additionally, in stage 2 an average structure of the best 10% of models was  
1282 calculated. Stage 3 served to refine this average structure to obtain a final set of best models.

1283 Restraints from top predicted contacts (position pairs with highest interaction scores and linear  
1284 chain separation greater than 5 positions) were implemented by setting C $\beta$ -C $\beta$  atom distances  
1285 (C $\alpha$  in case of Glycine) between positions to range between 0 and 8Å and weighting the  
1286 restraints according to their relative interaction score (interaction score divided by mean  
1287 interaction score of all predicted contacts used).

1288 Restraints from secondary structures elements were implemented as dihedral angle restraints.  
1289 Dihedral angles of both beta strands and alpha helices were set to range between values  
1290 commonly observed in crystal structures <sup>72</sup>, for alpha helices  $\Phi_{\alpha} = -63.5^{\circ} \pm 4.5^{\circ}$  and  $\Psi_{\alpha} =$   
1291  $-41.5^{\circ} \pm 5^{\circ}$  and for beta strands  $\Phi_{\beta} = -118^{\circ} \pm 10.7^{\circ}$  and  $\Psi_{\beta} = 134^{\circ} \pm 8.6^{\circ}$ .

1292 Restraints for beta sheet interactions were implemented by setting H-N:O=C hydrogen bond  
1293 distances between interacting positions to range between 1.8 and 2.1Å <sup>72</sup>, with weight one.  
1294 Predictions of beta sheet interactions derived from deep mutational scanning data yield a string  
1295 of interacting positions, but hydrogen bonding in beta sheets occurs in specific non-continuous  
1296 patterns between position pairs (between alternating positions off the interaction diagonal in  
1297 parallel beta sheets and between every second set of position pairs in anti-parallel beta sheets).  
1298 Specifically, for each set of interacting positions there are two alternative patterns of hydrogen  
1299 bonding possible. These alternative possibilities of pairing were implemented as mutually  
1300 exclusive selection pairs with the “assign ... or” syntax in Xplor-NIH.

1301 Distance restraints were implemented in XPLOR-NIH as NOE (nuclear Overhauser effect)  
1302 potential, with potential type set to “soft” for stages 1 and 2 and “hard” for the final simulation  
1303 stage. Dihedral angle restraints were implemented via the CDIH potential.

1304 After simulation stages 1 and 2 restraints were checked for their consistency with predicted  
1305 structural models. First, structural models were clustered according to their violations of  
1306 distance and dihedral angle restraints (k-means clustering,  $k = 4$ ). Clusters were ranked by the  
1307 mean total energy (from all energy potentials used) of their 50 models with lowest total energy  
1308 (or all of their structures if clusters are smaller 50 models). From the 50 models with the lowest  
1309 total energy from the top-ranked cluster (or however many top-ranked clusters were necessary  
1310 to arrive at 50 models) the fraction of models that violate specific restraints was recorded. For  
1311 the subsequent simulation stage, distance restraints were down-weighted according to the  
1312 fraction of models that violated them,  $w_{x,i} = w_{x,i-1} * (1 - f_x)^2$ , and distance restraints with a  
1313 weight below 0.1 were discarded. There is no option to weight dihedral angle restraints, thus  
1314 instead dihedral angle restraints with a ‘weight’ below 1/3 were discarded for the subsequent  
1315 simulation stages.

1316 The top 5% structural models from simulation stage 3 were evaluated against the reference  
1317 structure. The TM-score program (update 2016/03/23) was used to calculate the C $\alpha$  root mean  
1318 squared deviation and the template modeling score<sup>73</sup>.

1319 Several types of control simulations were performed to judge the predictive power of restraints  
1320 derived from deep mutational scanning data. As a negative control we performed simulations  
1321 without restraints from predicted contacts and beta sheet interactions, but with restraints from  
1322 secondary structure elements predicted by PSIPRED (version 3.3,<sup>43</sup>). As a positive control we  
1323 performed simulations with restraints derived from the reference structure. Here, L true contacts  
1324 of position pairs with linear chain distance greater than 5 amino acids were randomly sampled  
1325 and beta sheet interactions were determined by PyMOL<sup>44</sup>. These simulations serve as a  
1326 positive control and give the maximally achievable accuracy of our Xplor-NIH workflow.

1327 For the WW domain, simulations on the full mutated 33aa section gave mediocre results, both  
1328 when using combined scores with PSIPRED predicted secondary structure (5.8Å C $\alpha$ -  
1329 RMSD), as well as when using perfect information from the reference structure (4.1Å C $\alpha$ -  
1330 RMSD). Upon inspection, this seemed to be an issue of the unstructured tail regions. We thus  
1331 conducted structural simulations for a truncated version of the WW domain using only mutated  
1332 positions 6-29 (the core region including the three beta strands).

1333 For structural simulations of down-sampled GB1 datasets (and DeepContact transformed  
1334 versions thereof) we used distance restraints derived from top predicted contacts and  
1335 secondary structure restraints derived from PSIPRED predictions, but no restraints for beta  
1336 sheet pairing. This was done to avoid skewed results due to false beta sheet pairing predictions  
1337 in low quality datasets (Extended Data Figure 7B). For structural simulations from DeepContact-  
1338 transformed predictions, we find that using more tertiary contacts results in better models. We  
1339 conclude that this is because the deep learning algorithm focuses many strong predictions in  
1340 few structural features (such as interactions of secondary structure elements), which are  
1341 therefore the top contacts. Restraints in other regions of the protein are therefore only included  
1342 if more predicted contacts are used for restraint calculations, therefore improving structural  
1343 predictions. Because of this, when comparing structural simulations from scores derived before  
1344 and after deep learning, we compare the top 5% of structural models derived with the top L  
1345 predicted contacts from original scores with those derived with the top  $1.5 \cdot L$  predicted contacts  
1346 from DeepContact transformed scores.

1347

## 1348 DeepContact learning

1349 DeepContact software was obtained from GitHub (<https://github.com/largelymfs/deepcontact>)<sup>47</sup>.  
1350 We are grateful to Yang Liu and Jian Peng for also making - without any hesitation - their basic  
1351 DeepContact network architecture available on their GitHub repository and helping us with the  
1352 implementation. The DeepContact architecture used here only takes one 2D input of predicted  
1353 contact scores and returns a 2D map of transformed scores (denoted as “DeepContact  
1354 CCMpred only” in<sup>47</sup> and described in the first paragraph of the result section therein). The  
1355 DeepContact architecture employed came with a pre-trained network model that had been  
1356 trained on solved structures of the 40% homology filtered ASTRAL SCOPe 2.06 database (see  
1357 GitHub repository and Liu, et al.<sup>47</sup>), which were filtered to avoid structure and sequence  
1358 redundancy of the training data. Because CCMpred scores<sup>74</sup> are distributed in the range of 0 to  
1359 1, we pre-normalized our deep mutational scanning derived interaction scores to this range  
1360 (such that the minimum score on the interaction score matrix was 0 and the maximum score  
1361 was 1) before providing them as an input to DeepContact. As negative control, we created for  
1362 each dataset three random permutations of *combined score* matrices (while preserving matrix  
1363 symmetry; in case of FOS-JUN dataset non-symmetric *epistasis score* matrices were  
1364 permuted), which were transformed by the DeepContact algorithm. These control datasets  
1365 show no increased precision of random expectation (Figure 6F).

1366

## 1367 Code availability

1368 Data was analyzed with custom scripts written and executed in R programming language,  
1369 version 3.4.3. Structural simulations were performed with Xplor-NIH modeling suite version  
1370 2.46. Analysis scripts are available at <https://github.com/lehner-lab/DMS2structure>.

1371

## 1372 Data availability

1373 No primary data was generated in this study. Processed interaction scores for all datasets are  
1374 included in Supplementary Table 1. All intermediate steps of data processing can be  
1375 recapitulated with the scripts at <https://github.com/lehner-lab/DMS2structure>.

# 1376 Extended Data Figures

## Extended Data Figure 1

### deep mutational scanning experiment

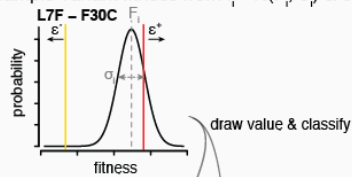
- 1) mutate protein of interest
- 2) competitive selection assay
- 3) genotype frequencies from sequencing

### computational analysis

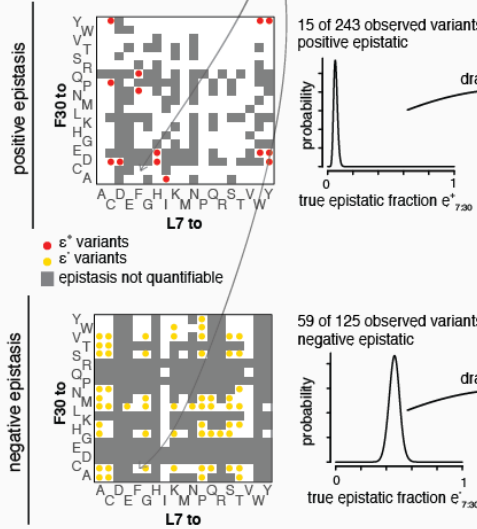
- 4) calculate fitness  $F_i$  and error  $\sigma_{F_i}$  from sequencing read counts

### 5) $10^4$ x re-sampling procedure

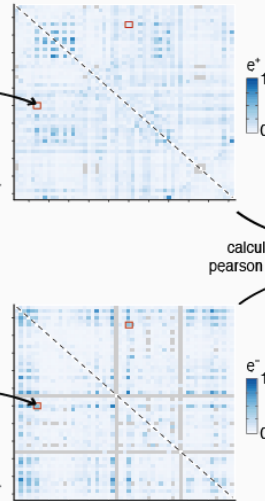
5a) sample variant fitness from  $f_i \sim N(F_i, \sigma_i)$  & classify epistasis (see panel C)



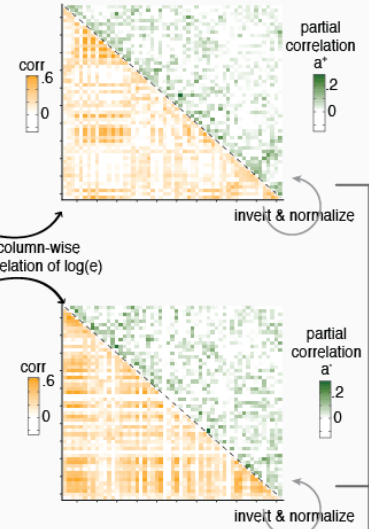
5b) aggregate per position pair & sample epistatic fractions



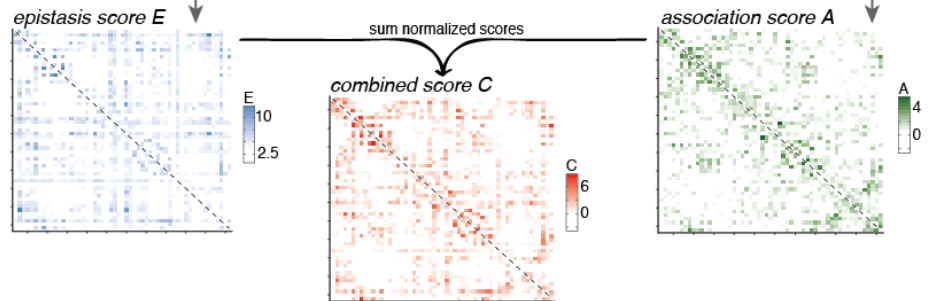
5c) epistatic fraction matrices



5d) correlation of epistasis patterns



6) compute **interaction scores** by merging positive & negative epistasis information



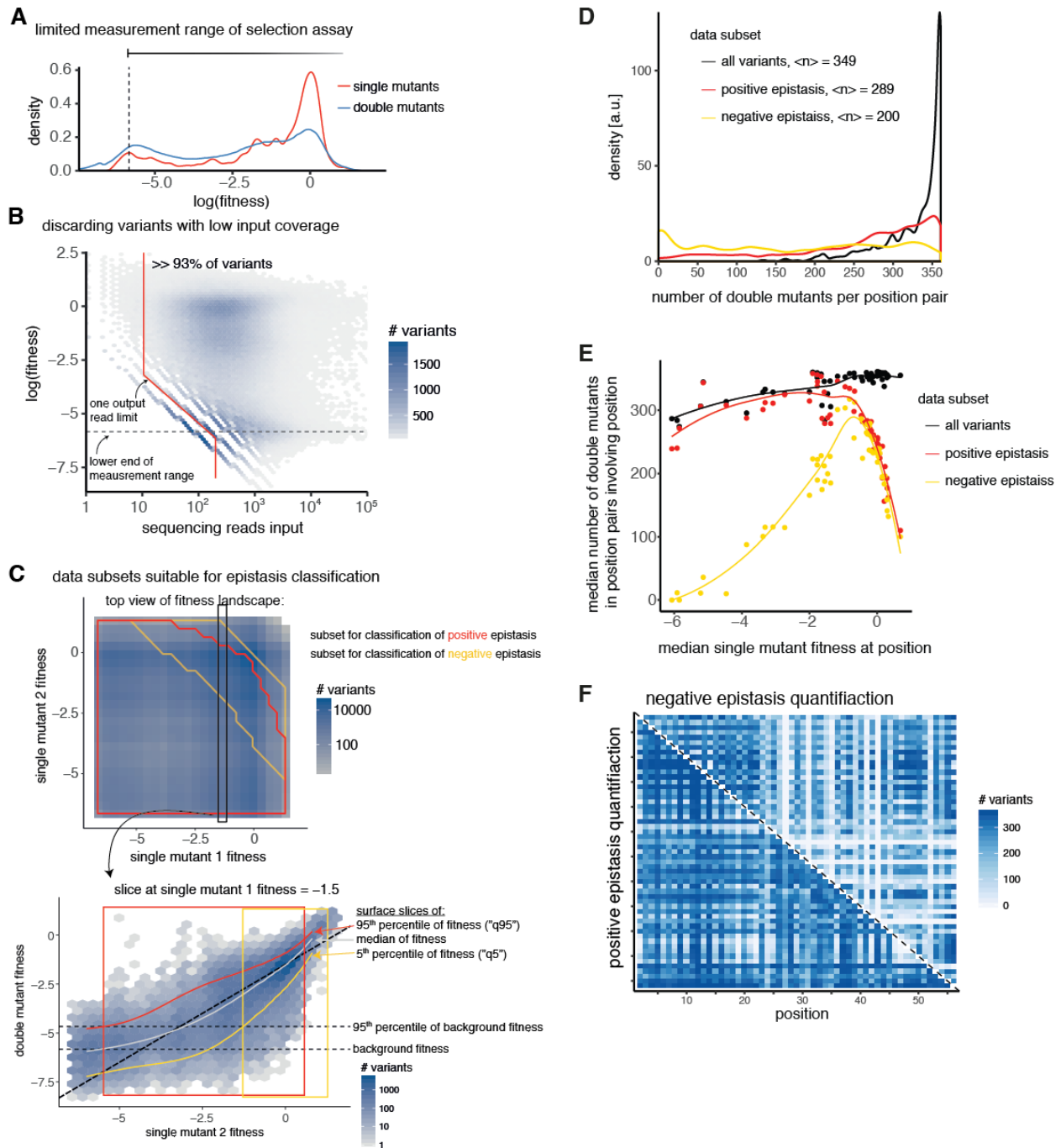
1377

1378 Extended Data Figure 1: Deep mutational sequencing data to contact  
1379 prediction workflow

1380 Overview of workflow to predict interacting position pairs from deep mutational scanning  
1381 datasets (see Methods and Results).

1382

## Extended Data Figure 2



1383

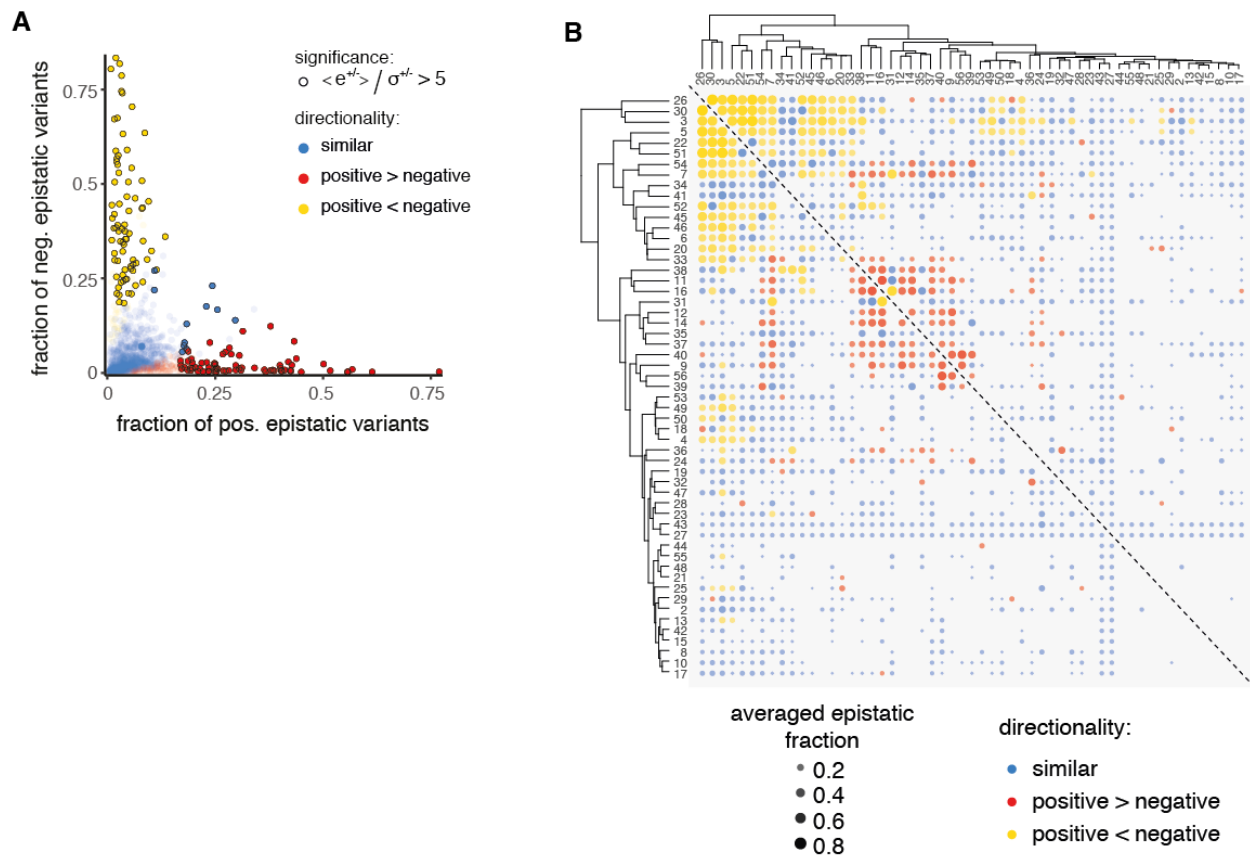
1384 Extended Data Figure 2: GB1 deep mutational scanning data processing

1385 A. Distribution of fitness values for single and double mutant variants. Lower peak in  
1386 distributions indicates lower limit of fitness assay measurement range (see Methods).



- 1387 B. Two-dimensional variant density showing dependency of fitness values on sequencing  
1388 read counts in input library. For variants with very low coverage in the input library low  
1389 fitness values cannot be accurately estimated. Red line shows sequence read cutoff  
1390 used for variant inclusion (93% of variants included for downstream analysis). Horizontal  
1391 dashed line indicates lower limit of fitness assay measurement range.
- 1392 C. Reliable quantification of positive and negative epistasis is limited to subsets of the data.  
1393 Upper plot shows two-dimensional double mutant variant density in single mutant fitness  
1394 space. Red outline shows single mutant fitness space that enables positive epistasis  
1395 quantification. Yellow outline shows single mutant fitness space that enables negative  
1396 epistasis quantification. Lower plot shows an example slice through fitness landscape at  
1397 single mutant fitness = -1.5. Red, grey and yellow curves show slices through quantile  
1398 fitness surfaces of 95<sup>th</sup> percentile, median and 5<sup>th</sup> percentile, respectively (see Figure  
1399 1B). Diagonal dashed line shows  $double\ mutant\ fitness =$   
1400  $single\ mutant\ 2\ fitness - 1.5$  (expected fitness = observed fitness). Horizontal dashed  
1401 lines give lower limit of fitness assay measurement range and the 95<sup>th</sup> percentile of  
1402 fitness values of variants dominated by background fitness effects. Red and yellow  
1403 boxes indicate the range that includes 99% of variants within the slice that are suitable  
1404 for positive or negative epistasis quantification, respectively.
- 1405 D. Distribution of number of double mutant variants across all position pair. Legend gives  
1406 median number of double mutants per position pair for different data subsets.
- 1407 E. Relationship between median single mutant fitness at a position and the median number  
1408 of double mutants observed in position pairs the position is involved in. Curves are loess  
1409 smoothed. Across all variants, positions with stronger fitness effects show lower  
1410 coverage of double mutants. Restrictions for quantification of positive epistasis  
1411 additionally reduce coverage for positions with mostly neutral or positive effects. Finally,  
1412 restrictions for quantification of negative epistasis strongly reduce coverage for positions  
1413 with strong fitness effects, due to the lower measurement limit of the fitness assay.
- 1414 F. Number of double mutants for which positive (lower left triangle) or negative (upper right  
1415 triangle) epistasis can be quantified per position pair plotted on the interaction matrix.
- 1416

Extended Data Figure 3



1417

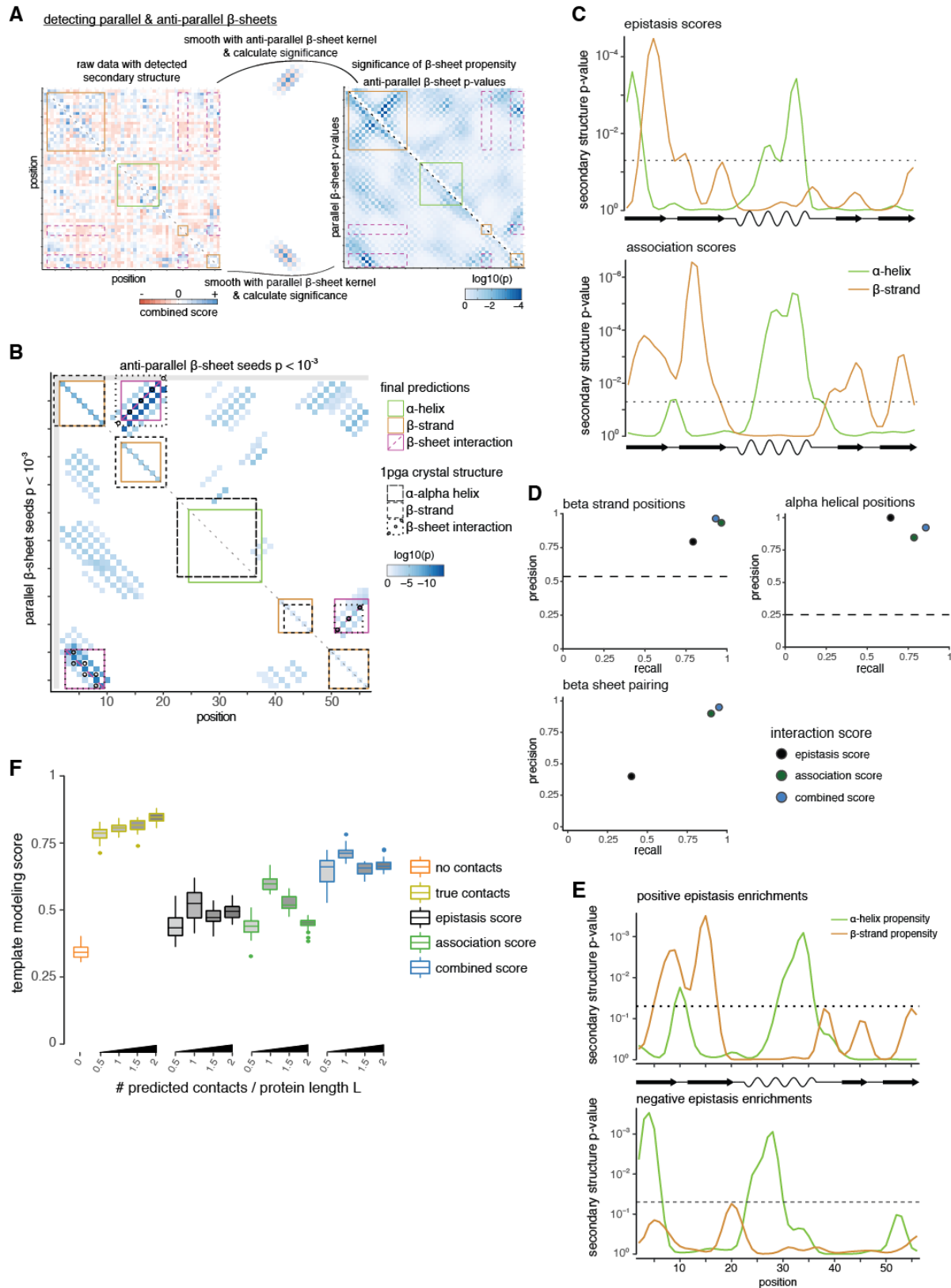
1418 Extended Data Figure 3: Positive and negative epistasis enrichments  
1419 across position pairs

1420 A. Position pair-wise fractions of positive and negative epistatic variants. Black  
1421 circles mark position pairs with highly significant fractions (either positive or  
1422 negative, or both); red dots: positive epistatic fraction significantly larger than  
1423 negative epistatic fraction; yellow dots: negative epistatic fraction significantly  
1424 larger than positive epistatic fraction; blue dots: no significant differences.

1425 B. Hierarchical cluster analysis of epistatic fraction patterns. Positions are clustered  
1426 according to the Euclidean distance of their mean epistatic fractions (weighted  
1427 average of positive and negative epistatic fractions, weights are inverse  
1428 variances of fractions in resampling runs) to all other positions. Note that  
1429 directionality of interactions (positive or negative fractions more significant) was  
1430 not used for clustering but only marked post-analysis. Clustering shows two  
1431 highly interconnected clusters of positions that interact mostly positively or

1432 negatively within each cluster but hardly any strong interactions are observed  
1433 between the two clusters (with exception of positions 7 and 54).  
1434

## Extended Data Figure 4



1436 Extended Data Figure 4: Secondary and tertiary structure prediction for  
1437 GB1 domain

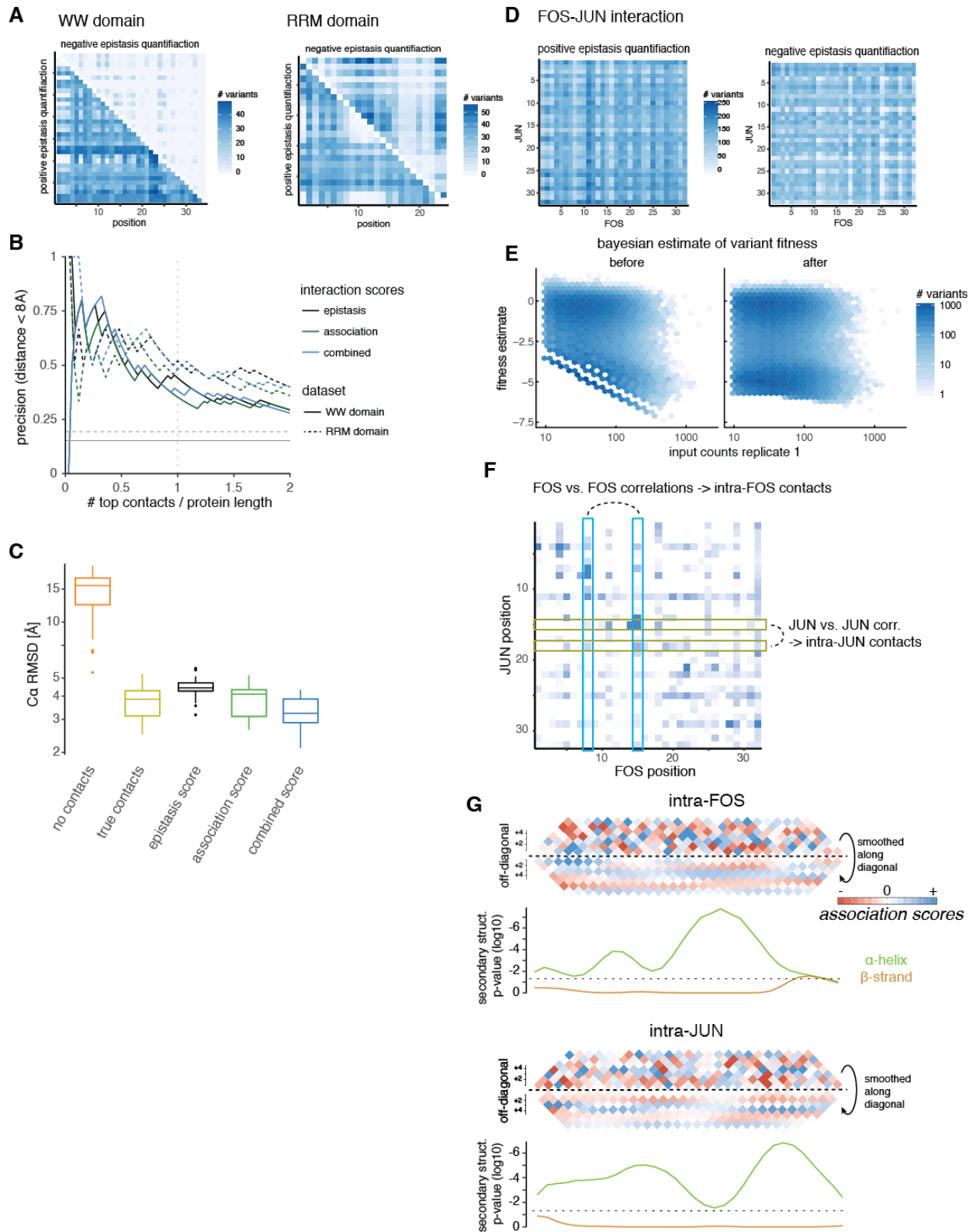
- 1438 A. Detecting beta sheet pairing with two-dimensional kernel smoothing. Left plot shows raw  
1439 combined score interaction matrix, with secondary structure element predictions (see  
1440 Figure 4A,B) marked as squares along the diagonal (red – beta strand, green – alpha  
1441 helix). Off-diagonal orange rectangles show potential regions of beta sheet pairing. Right  
1442 plot: calculation of beta sheet pairing propensity with beta sheet kernels. Upper right  
1443 triangle shows anti-parallel beta sheet propensity. Lower left shows parallel beta sheet  
1444 propensity.
- 1445 B. Matrix of aggregated propensities of beta sheet pairing stretches (upper right – anti-  
1446 parallel, lower left – parallel,  $p < 10^{-3}$ ) and the predictions for beta sheet pairing and  
1447 secondary structure elements derived from them. In brief, predictions are performed by  
1448 picking the highest propensity stretch that is consistent with predicted beta strands, if  
1449 necessary modifying beta strand predictions (e.g. introducing an initially not predicted  
1450 split between beta strands 1 and 2), then disregarding all stretches that conflict with the  
1451 picked top-stretch. This procedure is repeated until no more beta sheet stretches with  
1452 propensity  $P < 10^{-3}$  are left. Finally, beta strand predictions are updated such that only  
1453 positions involved in a beta sheet interaction are retained. Reference elements from  
1454 crystal structure are shown as comparison (lower triangle – parallel beta sheets, upper  
1455 triangle – anti-parallel beta sheet, diagonal – secondary structure elements).
- 1456 C. Secondary structure propensity derived from kernel smoothing (red – beta strand, green  
1457 – alpha helix) for *epistasis* (upper) and *association scores* (lower). P-values were  
1458 derived by comparison to randomized datasets (see Methods). Dashed line indicates  $p =$   
1459 0.05.
- 1460 D. Precision and recall for beta strand, alpha helix and beta sheet predictions from  
1461 *epistasis*, *association* and *combined scores* (in comparison to crystal structure). Dashed  
1462 lines for beta strand and alpha helical positions give random expectation. Random  
1463 expectation for beta sheet pairing precision is below 1%.
- 1464 E. Secondary structure propensities derived from local positive or negative epistatic  
1465 enrichments. The upper panel shows secondary structure propensity derived from  
1466 positive epistatic interactions, which are in line with secondary structure elements in the  
1467 GB1 crystal structure (PDB entry 1pga). The lower panel shows secondary structure  
1468 propensity derived from negative epistatic interactions, which are devoid of beta strand

1469 signals and instead show a three-helical pattern, which is reminiscent of the three-helical  
1470 structure of the protein G A domain that binds albumin <sup>57</sup>.

1471 F. Template modeling score of top 5% structural models compared to crystal structure  
1472 1pga and the dependency on number of predicted contacts used. “No contacts” – only  
1473 restraints for secondary structure predicted by PSIPRED. “True contacts” – restraints  
1474 derived from 0.5-2\*L contacts (linear sequence separation greater than 5 positions,  
1475 random subset), secondary structure elements and beta sheet interactions from crystal  
1476 structure. All other: restraints derived from top 0.5-2\*L contacts, secondary structure  
1477 element and beta sheet interaction predictions from the three interaction scores, as  
1478 indicated by color.

1479

Extended Data Figure 5



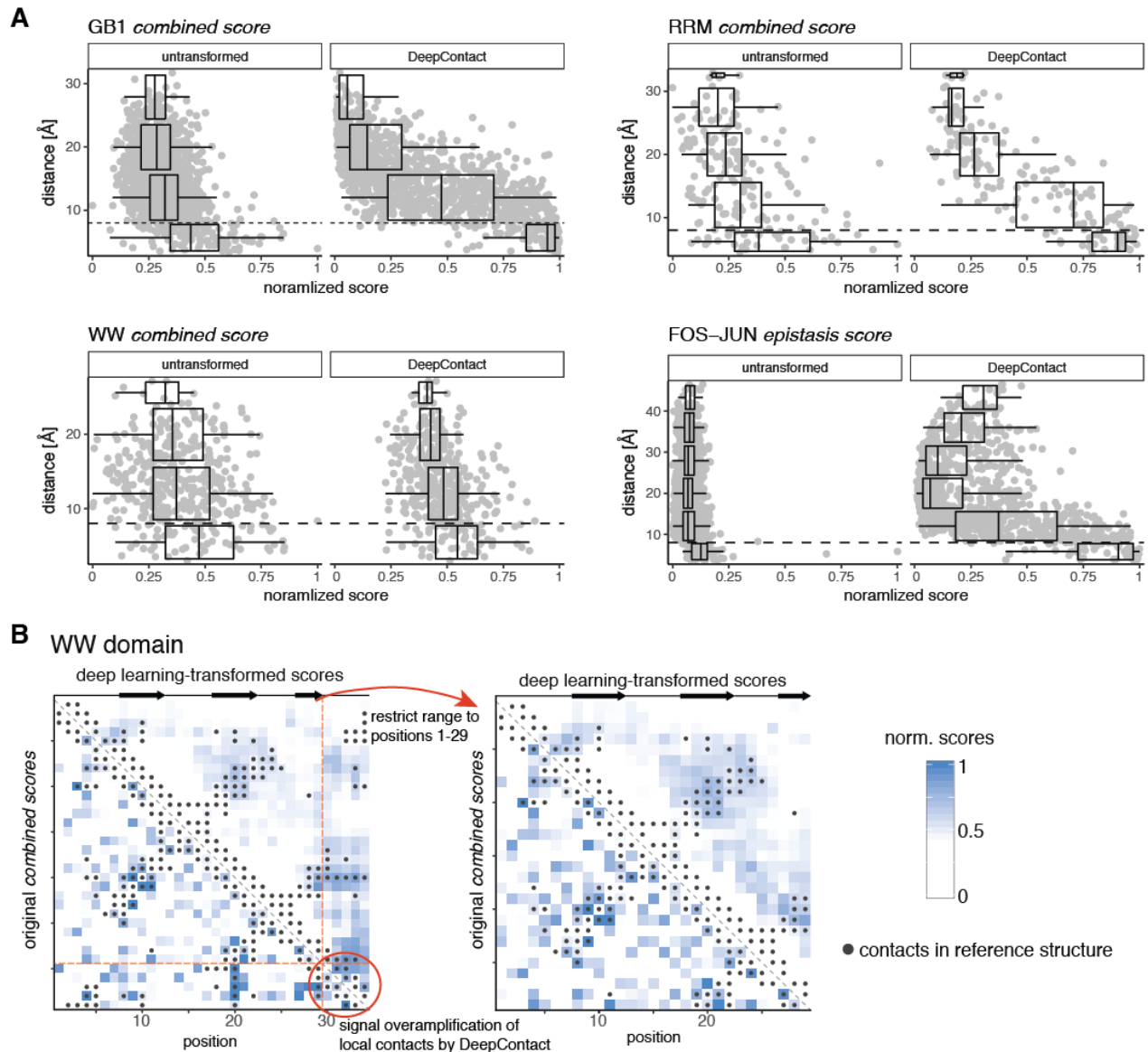
1481 Extended Data Figure 5: WW and RRM domain and FOS-JUN  
1482 interaction dataset properties

- 1483 A. Number of double mutants for which positive (lower left triangle) or negative (upper right  
1484 triangle) epistasis can be quantified per position pair plotted on the interaction matrix for  
1485 WW (left) and RRM (right) domains.
- 1486 B. Precision of interaction scores to predict direct contacts (distance  $< 8\text{\AA}$  in reference  
1487 structure) as a function of top scoring position pairs for WW and RRM domain interaction  
1488 scores. Color denotes interaction scores, solid lines for WW domain, dashed lines for  
1489 RRM domain. Grey horizontal lines give random expectation. Only position pairs with  
1490 linear sequence separation greater than 5 amino acids are considered.
- 1491 C. Accuracy ( $C\alpha$  root-mean-square deviation) of top 5% structural models of the WW  
1492 domain (core positions 6-29) generated from deep mutational scanning data derived  
1493 restraints compared to reference structure (PDB entry 1k9q). Structural models were  
1494 generated in XPLOR-NIH by simulated annealing with restraints derived from top 17  
1495 tertiary contacts and secondary structure elements predicted by PSIPRED. No beta  
1496 sheet pairing information was used.
- 1497 D. Number of double mutants for which positive (left) or negative (right) epistasis can be  
1498 quantified per position pair plotted on the trans-interaction matrix of the FOS-JUN  
1499 interaction.
- 1500 E. Bayesian estimation of fitness values in FOS-JUN interaction data. Mutants with low  
1501 input sequencing coverage display limited measurement range and many dropouts  
1502 (~15% of variants without reads in output). Left panel shows original fitness distribution  
1503 as function of input coverage in replicate 1, right panel shows Bayesian estimates of  
1504 fitness as function of input coverage in replicate 1.
- 1505 F. Learning about intra-molecular contacts in FOS or JUN from epistatic pattern  
1506 correlations. Column-wise correlation of epistatic patterns of the trans interaction score  
1507 map serve to calculate intra-FOS *association scores* and thus reveal relationships  
1508 between positions in FOS. Likewise, row-wise correlation of epistatic patterns reveal  
1509 relationships between positions in JUN.



1510 G. Local interactions in intra-molecular *association scores* reveal secondary structures of  
1511 protein interaction partners. Upper panels: Data above diagonal shows *association score*  
1512 data close to the diagonal, i.e. local interactions. Data below the diagonal is smoothed  
1513 with a Gaussian kernel to reveal interaction periodicity. Lower panels: Secondary  
1514 structure propensities derived from kernel smoothing (see Figure 4A-C). Green indicates  
1515 alpha helical propensity, orange indicates beta sheet propensity,  $p = 0.05$  is indicated by  
1516 dashed line.  
1517

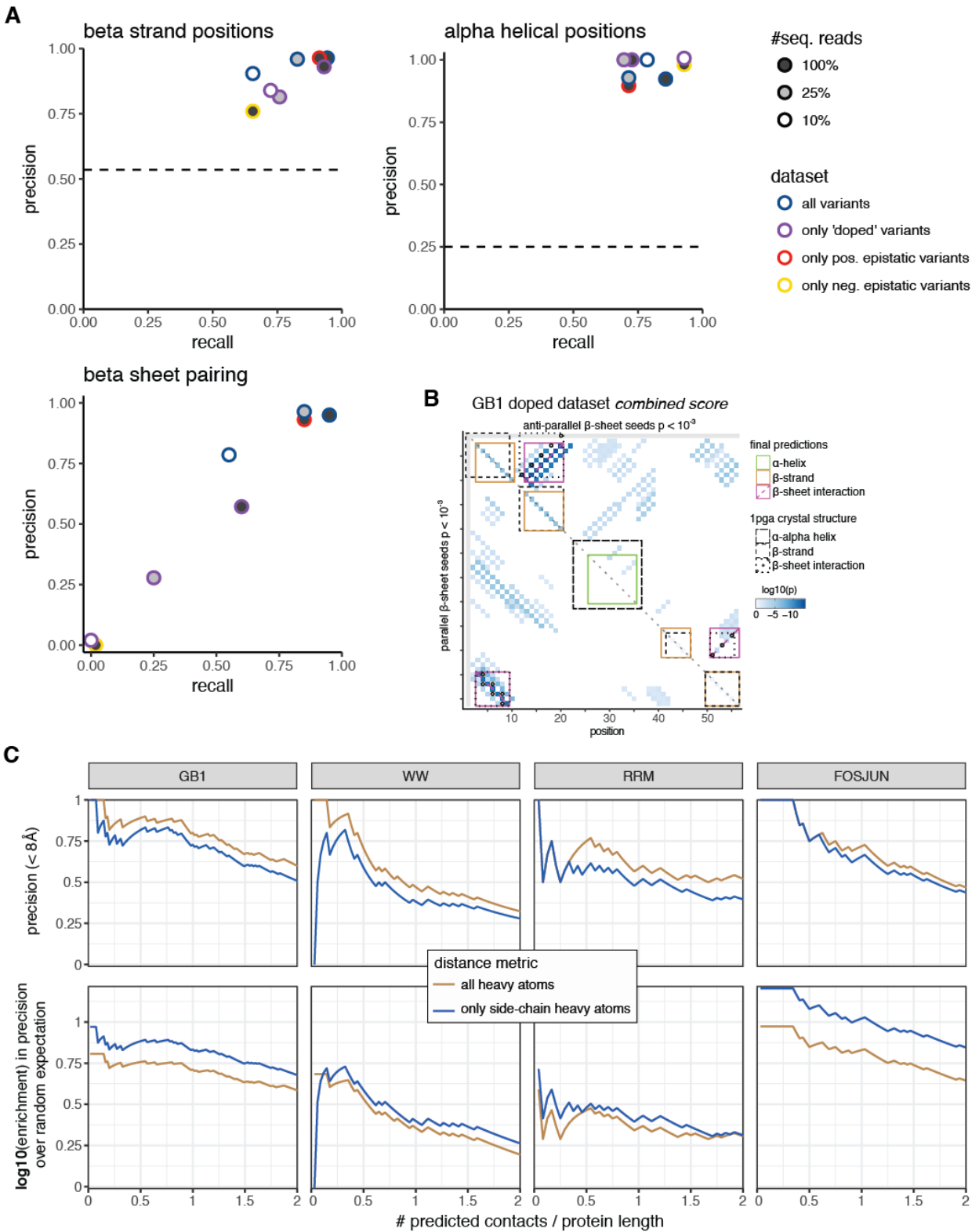
## Extended Data Figure 6



1527 contacts in the C-terminal region of the domain, thus concentrating the strongest  
1528 transformed signal in this region. Removing positions 30-34 removes this artefact (right  
1529 plot). Heat maps show interaction scores that have been normalized to have similar  
1530 range. Grey dots show contacts (distance < 8Å) in reference structure.

1531

## Extended Data Figure 7



1533 **Extended Data Figure 7: Distance metric comparison and**  
1534 **secondary structure prediction for lower data-quality GB1**  
1535 **datasets**

- 1536 A. Precision and recall for beta strand, alpha helix and beta sheet predictions derived from  
1537 combined scores of down-sampled GB1 datasets (in comparison to reference structure).  
1538 Dashed lines for beta strand and alpha helical positions give random expectation.  
1539 Random expectation for beta sheet pairing precision is below 1%. Note that some  
1540 coinciding data points were slightly moved for better identifiability.
- 1541 B. Beta sheet pairing predictions for the doped GB1 dataset with 100% sequencing read  
1542 coverage (cf. Extended Data Figure 4B). Beta sheet pairing between beta strands 1 and  
1543 2 is predicted in correct anti-parallel direction, but exact pairing of positions are off by 2;  
1544 thus precision and recall of beta sheet pairing for doped GB1 dataset drops to ~60%  
1545 (see panel B).
- 1546 C. Differences in precision and enrichment over random expectation for all heavy atom or  
1547 side-chain heavy atom distance metrics. As expected, using all heavy atoms (including  
1548 backbone heavy atoms) increases precision of predicted contacts by about 10%.  
1549 Restricting distance measurements to side-chain heavy atoms, however, increases  
1550 precision over random expectation, often by more than 2-fold (note the log<sub>10</sub>-scale),  
1551 indicating that side-chain distances are more informative for epistatic interactions. For  
1552 these calculations, only position pairs with linear sequence separation greater than 5  
1553 amino acids were considered.