

# GigaScience

## Integrating deep mutational scanning and low-throughput mutagenesis data to predict the impact of amino acid variants

--Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-23-00040R2	
<b>Full Title:</b>	Integrating deep mutational scanning and low-throughput mutagenesis data to predict the impact of amino acid variants	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	National Health and Medical Research Council (116955)	Professor Anthony Troy Papenfuss
	National Human Genome Research Institute (RM1HG010461)	Dr Alan F. Rubin
	National Human Genome Research Institute (UM1HG011969)	Dr Alan F. Rubin
	Lorenzo and Pamela Galli Medical Research Trust	Professor Anthony Troy Papenfuss
	Stafford Fox Medical Research Foundation	Professor Anthony Troy Papenfuss
	Melbourne Research Scholarship	Mr Yunfan Fu
<b>Abstract:</b>	<p>Background: Evaluating the impact of amino acid variants has been a critical challenge for studying protein function and interpreting genomic data. High-throughput experimental methods like deep mutational scanning (DMS) can measure the effect of large numbers of variants in a target protein, but because DMS studies have not been performed on all proteins, researchers also model DMS data computationally to estimate variant impacts by predictors.</p> <p>Results: In this study, we extended a linear regression-based predictor to explore whether incorporating data from alanine scanning (AS), a widely used low-throughput mutagenesis method, would improve prediction results. To evaluate our model, we collected 146 AS datasets, mapping to 54 DMS datasets across 22 distinct proteins.</p> <p>Conclusions: We show that improved model performance depends on the compatibility of the DMS and AS assays, and the scale of improvement is closely related to the correlation between DMS and AS results.</p>	
<b>Corresponding Author:</b>	Alan F. Rubin, PhD Walter and Eliza Hall Institute of Medical Research Parkville, VIC AUSTRALIA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Walter and Eliza Hall Institute of Medical Research	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Yunfan Fu	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Yunfan Fu	
	Justin Bedó, PhD	
	Anthony Troy Papenfuss, BSc (Hons) PhD	
	Alan F. Rubin, PhD	
<b>Order of Authors Secondary Information:</b>		
<b>Response to Reviewers:</b>	We would like to thank the reviewers again for their thoughtful comments and careful	

	consideration of the manuscript and revisions. We hope that this work contributes to the ongoing conversation in the field around modelling and analysis of protein mutagenesis data, and that the datasets we curated and present here will be useful for future studies.
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
<b>Resources</b>	Yes
<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	
<b>Availability of data and materials</b>	Yes
<p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using</p>	

a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

# 1 Integrating deep mutational scanning and low-through- 2 put mutagenesis data to predict the impact of amino acid 3 variants

4

## 5 Authors:

6 Yunfan Fu<sup>1,2</sup>, Justin Bedó<sup>1,2,\*</sup>, Anthony T. Papenfuss<sup>1,2,3,\*\*</sup>, Alan F. Rubin<sup>1,2,\*,\*\*</sup>

7

## 8 Affiliations:

9 <sup>1</sup>The Walter and Eliza Hall Institute of Medical Research, Parkville, VIC 3052, Australia.

10 <sup>2</sup>Department of Medical Biology, The University of Melbourne, Parkville, VIC 3010, Australia.

11 <sup>3</sup>Peter MacCallum Cancer Centre, Melbourne, VIC 3000, Australia.

12

13 \* Contributed equally

14 \*\* To whom correspondence should be addressed (papenfuss@wehi.edu.au & alan.rubin@wehi.edu.au)

15

## 16 Abstract

17 **Background:** Evaluating the impact of amino acid variants has been a critical challenge for  
18 studying protein function and interpreting genomic data. High-throughput experimental meth-  
19 ods like deep mutational scanning (DMS) can measure the effect of large numbers of variants  
20 in a target protein, but because DMS studies have not been performed on all proteins, research-  
21 ers also model DMS data computationally to estimate variant impacts by predictors.

22 **Results:** In this study, we extended a linear regression-based predictor to explore whether in-  
23 corporating data from alanine scanning (AS), a widely used low-throughput mutagenesis

24 method, would improve prediction results. To evaluate our model, we collected 146 AS da-  
25 taset, mapping to 54 DMS datasets across 22 distinct proteins.

26 **Conclusions:** We show that improved model performance depends on the compatibility of the  
27 DMS and AS assays, and the scale of improvement is closely related to the correlation between  
28 DMS and AS results.

29

30 **Keywords:** deep mutational scanning, alanine scanning, machine learning, predictor

31

## 32 **1 Introduction**

33 Deep mutational scanning (DMS) is a functional genomics method that can experimentally  
34 measure the impact of many thousands of protein variants by combining high-throughput se-  
35 quencing with a functional assay [1]. In a typical DMS, a cDNA library of genetic variants of  
36 a target gene is generated, containing all possible single amino acid substitutions. This variant  
37 library is then expressed in a functional assay system where the DMS variants can be selected  
38 based on their properties. The change in variant frequency in the pre- and post-selection popu-  
39 lations is determined by high-throughput sequencing which is then used to calculate a multi-  
40 plexed functional score that captures the variant's impact [2–4]. The versatility of DMS assays  
41 makes it possible to measure variant impact on a wide range of protein properties, including  
42 protein binding affinity [5,6], protein abundance [7–9], enzyme activity [10,11] and cell sur-  
43 vival [12–14]. So far, hundreds of DMS studies covering tens of thousands of nucleotides have  
44 been published [15], and experiments targeting over a hundred additional genes are underway  
45 according to MaveRegistry [16].

46

47 Computational studies have used DMS data to build predictive models of variant impact. These  
48 predictors use supervised or semi-supervised learning models trained on experimental DMS  
49 data and various protein features to make predictions [17–23]. Envision is one such method  
50 that used protein structural, physicochemical, and evolutionary features to predict variant effect  
51 scores and was trained on DMS data from 8 proteins using gradient boosting [17]. Another  
52 method, DeMaSk, predicted DMS scores by combining two evolutionary features (protein po-  
53 sitional conservation and variant homologous frequency) with a DMS substitution matrix and  
54 was trained on data from 17 proteins using a linear model [19]. Deep learning algorithms have  
55 also been applied to build protein fitness predictors [18,20], which are usually based only on  
56 variant sequences. These variant effect predictors can also be benchmarked using DMS exper-  
57 imental results and assist in the interpretation of experimental data [20,24,25].

58

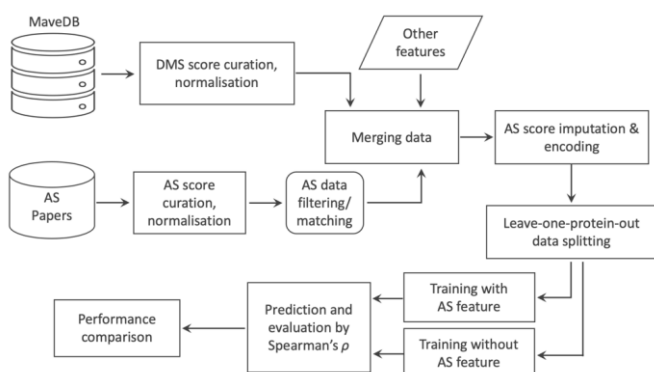
59 Low-throughput mutagenesis experiments that measure tens of variants at a time have also  
60 been used extensively to study diverse protein properties, including substrate binding affinity  
61 [26,27], protein stability [28,29], and protein-specific activities [30,31]. Alanine scanning (AS)  
62 is a widely-used low-throughput mutagenesis method [32,33], and AS data are available for  
63 many proteins. In this method, each targeted protein residue is substituted with alanine, and the  
64 impacts of these variants are measured by a functional assay [34]. AS experiments are typically  
65 used to identify functional hot spots or critical residues in the target protein [35,36] and have  
66 been used as a source of independent validation for DMS studies [31,37–39].

67

68 In this study, we explore whether a predictive model can be improved by incorporating low-  
69 throughput mutagenesis data (Fig 1). We find that AS data can increase prediction accuracy

70 and that the improvement is related to the similarity of the functional assays and the correlation  
71 of DMS and AS results.

72



73

74 **Fig 1. Workflow for model training and testing.** DMS and AS datasets are collected from online resources and  
75 are normalized. DMS and AS datasets targeting the same protein are then matched, filtered and merged. Two  
76 predictors are constructed and tested: the first uses DMS data, AS data and other protein features, and the second  
77 uses only DMS data and the same other protein features.

78

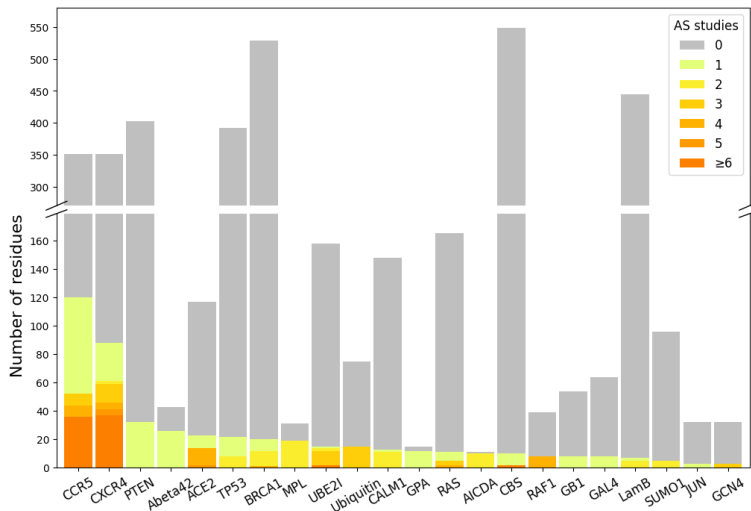
## 79 2 Results

### 80 2.1 Overview of DMS and alanine scanning (AS) data

81 To build the predictive model, 130 DMS datasets were collected from MaveDB [40,41] (Sup-  
82 plementary table 1). We searched the literature and found 146 AS datasets targeting the same  
83 proteins as 54 of the DMS datasets. In total, we obtained both DMS and AS data for 22 different  
84 proteins: 17 human proteins, three yeast proteins, and two bacterial proteins. Most DMS ex-  
85 periments were highly complete, with a mean coverage of 95.0% of all possible single amino  
86 acid substitutions assayed in the target region, comprising 373,219 total protein variant meas-  
87 urements. AS data were only available on a small number of protein residues (Fig 2), and we

88 were able to curate 1,480 alanine substitution scores from the 146 studies. Variant scores from  
89 collected DMS and AS studies were linearly normalized to a common scale (see Methods) to  
90 make them comparable across datasets (Fig S1).

91



92

93 **Fig 2. DMS data generally cover more protein residues than AS data.** Each bar shows the number of residues  
94 assayed by DMS studies on given target proteins. Colour indicates the number of AS studies available for the  
95 DMS-tested residues.

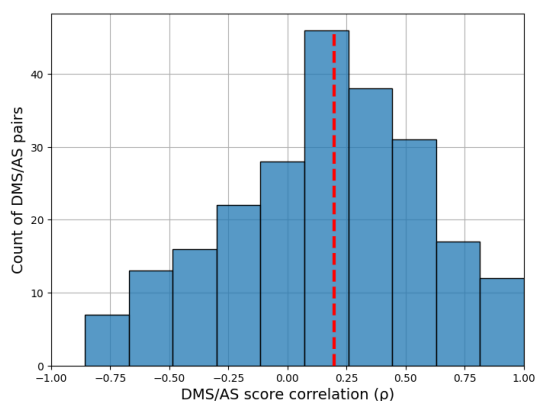
96

## 97 2.2 The correlation of DMS and AS scores is related to assay compatibility

98 To evaluate the similarity of AS and DMS scores, we calculated Spearman's correlation ( $\rho$ )  
99 between the AS scores and DMS scores for the same alanine substitutions. Since each protein  
100 may have results from several AS and DMS experiments, we calculated  $\rho$  between each possi-  
101 ble pair. The median  $\rho$  over DMS and AS data (DMS/AS) pairs was 0.2, indicating that the  
102 experimental scores were poorly correlated overall (Fig 3).

103





104

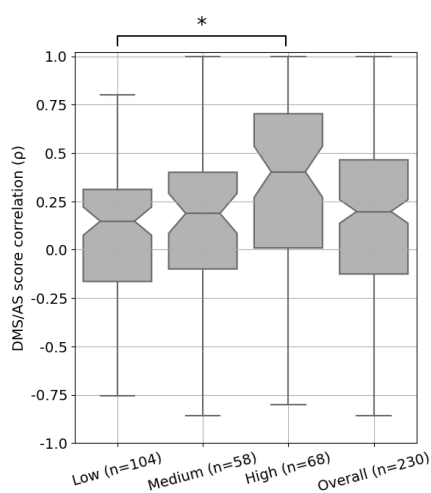
105 **Fig 3. Correlation between DMS and AS data shows substantial variation.** We calculated Spearman's  $\rho$  be-  
 106 tween alanine substitution scores in each pair of AS and DMS data. The results for pairs with less than three  
 107 alanine substitutions are not shown. The red dashed line shows the median  $\rho$ .

108

109 We then considered if differences between AS and DMS assay designs might contribute to this  
 110 low agreement between scores. To explore this, we developed a decision tree (Fig S2) to clas-  
 111 sify whether DMS/AS pairs had low, medium, or high assay compatibility, which we defined  
 112 as a similarity measurement of the functional assays performed. For example, the DMS assay  
 113 measuring the binding affinity of a cell surface protein, CXCR4, to its natural ligand [42] has  
 114 high compatibility with the AS experiment also measuring this ligand binding but has low  
 115 compatibility with the study on CXCR4's ability to facilitate virus infection [43]. A full assay  
 116 compatibility table can be found in Supplementary Table 1 with the compatibility classifica-  
 117 tions and justification for each pair. We then compared DMS and AS score correlation for each  
 118 compatibility class and found that score correlations were closely related to assay compatibility.  
 119 Data from low compatibility assays had a median correlation of 0.15, rising to 0.19 for medium  
 120 compatibility assays and 0.40 for high compatibility assays (Fig 4). This trend of increased  
 121 correlation for high compatibility assay pairs holds across secondary structures (Table S1).

122 This link between assay compatibility and score correlation indicates that our decision tree  
123 approach was able to capture the similarity between assay systems.

124



125

126 **Fig 4. DMS and AS data pairs with high assay compatibility show a higher score correlation.** Each box  
127 shows the Spearman's  $\rho$  between DMS and AS data pairs for each level of assay compatibility or overall. The  
128 correlation coefficients were calculated between alanine substitution scores in each pair of AS and DMS datasets.  
129 Results for pairs with less than three alanine substitutions were removed. P-values calculated using Welch's test  
130 and corrected using Holm-Šidák, \*:  $p < 0.05$ ; notches show 95% confidence interval around median, and whiskers  
131 show the full value range.

132

### 133 2.3 Compatible AS data improve DMS score prediction accuracy

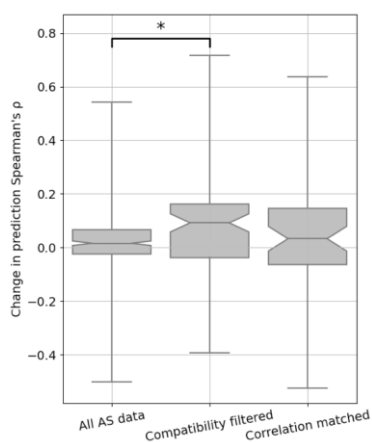
134 To test if incorporating AS data into DMS score models would improve prediction accuracy,  
135 we decided to build a new model based on DeMaSk [19]. We chose DeMaSk because it showed  
136 better performance compared to similar methods and was straightforward to modify. The pub-  
137 lished DeMaSk model predicts DMS scores using protein positional conservation, variant ho-

138 mologous frequency, and substitution score matrix, and we incorporated AS data as an addi-  
139 tional feature. Our new predictor was modelled with all 130 DMS we collected and we applied  
140 a leave-one-protein-out cross-validation approach to training and testing, avoiding information  
141 leakage for variants of the same protein target [17]. Prediction performance was evaluated us-  
142 ing the Spearman's correlation ( $\rho$ ) between the experimentally-derived DMS scores and the  
143 predicted scores for each pair of DMS and AS studies. The performance of our DMS/AS model  
144 was compared with a model trained only on DMS data, equivalent to retrained DeMaSk (Fig  
145 S3), by calculating the change of prediction  $\rho$  (see Methods).

146

147 We trained our model with either all or a subset of AS data we collected (Fig 5, Table S2). We  
148 first integrated all 146 AS data collected for training and evaluation but observed only a modest  
149 improvement of prediction  $\rho$  (Fig 5 left box, and Fig S4). We then retrained and evaluated our  
150 model on filtered AS data with only high compatibility assays, and observed a median increase  
151 in prediction Spearman's  $\rho$  of 0.1 compared to the results with no AS data (Fig 5 middle box,  
152 and Fig S4). However, training with both high and medium compatibility pairs reduced the  
153 performance improvement (Fig S5). These results indicate that medium and low compatibility  
154 pairs might provide inconsistent training data, degrading model performance. We also evalu-  
155 ated the impact of including high compatibility AS data in an alternative model based on En-  
156 vision [17], and found similar results (Fig S6). To differentiate between high assay compatibil-  
157 ity and high DMS/AS score correlation, we trained the model using the most highly correlated  
158 AS result for each DMS dataset (see Methods). Although the upper quartile was high, the me-  
159 dian performance change of this predictor was lower than the high assay compatibility model,  
160 suggesting that matching with the highest score correlation alone is insufficient (Fig 5 right  
161 box). However, when applying a stricter threshold, the correlation matched models still show

162 limited improvement (Fig S7). Additionally, to ensure the models performance is not biased  
163 by pseudo-replication of multiple datasets, we averaged DMS and AS scores that were part of  
164 the same study and type of assay, and saw similar results (Fig S8).  
165



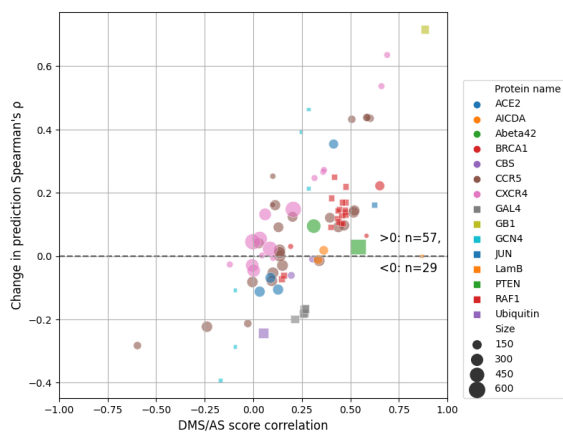
166  
167 **Fig 5. Performance of variant impact prediction is improved using AS data with high assay compatibility.**

168 The change in prediction  $\rho$  achieved by including the AS data feature for each DMS and AS data pair is shown as  
169 box plots. A higher value represents higher prediction accuracy achieved for using AS data. Different approaches  
170 to filtering/matching the data are shown on the x-axis: "All AS data" used all available data; "Compatibility fil-  
171 tered" used only data of high assay compatibility; "Correlation matched" used only data with the highest regular-  
172 ised correlation for each DMS dataset. Results for data pairs with only one residue are not shown. P-values were  
173 calculated using Welch's test and jointly corrected using Holm-Šidák (Methods), \*:  $p < 0.05$ . Notches show the  
174 95% confidence interval around the median, and whiskers show the full value range.

175  
176 Our compatibility-filtered predictor shows improved prediction accuracy for these regions  
177 compared to not only the baseline model, but other widely used predictors as well (Fig S9). To  
178 further explore the higher performance of this compatibility-filtered predictor, we examined  
179 the relationship between prediction  $\rho$  change and score correlation for each high compatibility

180 DMS/AS pair (Fig 6). For most pairs, prediction performance was improved by using AS data,  
181 and the scale of improvement was also related to the score correlation. This relationship could  
182 also be observed for multiple DMS/AS pairs from an individual protein, such as CXCR4 and  
183 CCR5. We saw the same trend in the predictor trained with all DMS/AS pairs but noted that  
184 the performance even of highly correlated pairs was worse, likely due to the influence of low  
185 compatibility training data on the model (Fig S10).

186



187

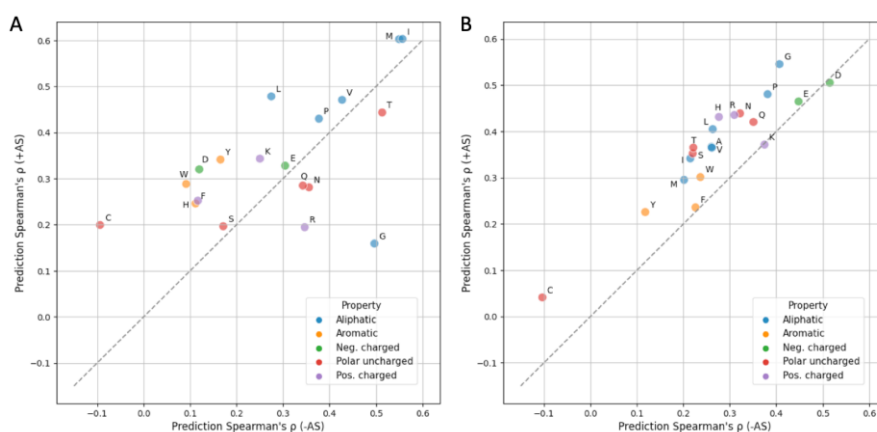
188 **Fig 6. Prediction performance change is related to DMS and AS score correlation.** Each dot represents a  
189 filtered DMS/AS data pair of high assay compatibility. The vertical axis shows the change of prediction  $\rho$  by using  
190 AS data (larger means higher performance achieved by using AS data). The horizontal axis shows the DMS/AS  
191 score correlation for *all* variants on the matched residues rather than just alanine substitutions. The colours and  
192 shapes of the dots correspond to the target protein, and size indicates the number of variants in each data pair.  
193 Results for data pairs with only one residue are not shown.

194

195 We also explored the consequences of the sparsity of AS data on our model in three ways: i)  
196 by training only with variants that have AS data available (Fig S11); ii) by using a boosting

197 approach that focuses only on residues with AS data (Fig S12) and iii) by using complete ala-  
198 nine substitution information from DMS as the AS feature (Fig S13). The first approach gave  
199 lower absolute prediction performance, presumably because the model was under-fitted due to  
200 the small number of variants. The last two approaches performed very similarly to the primary  
201 model constructed using high-compatibility DMS/AS data and simple mean score imputation.  
202  
203 To test the influence of amino acids on our predictor, we grouped the prediction results by  
204 either wild-type or variant amino acid and calculated the prediction improvement when AS  
205 data were included (Fig 7). We found that 14 of 19 wild-type amino acids performed better  
206 with the addition of AS data, with cysteine showing the largest improvement and performing  
207 worst in the model lacking AS data. 18 of 20 variant amino acids benefited from the inclusion  
208 of AS data, with marginal performance decrease on lysine and aspartic acid ( $|\Delta\rho|<0.01$ ) (Fig  
209 7). We also noticed that variants to alanine are not most improved, however we observed an  
210 overall trend showing higher improvement for amino acids that are physiochemically similar  
211 to alanine (Fig S15).

212



213

214 **Fig 7. Model performance is generally improved for each wild-type and variant amino acid.** Prediction  
215 Spearman's  $\rho$  when using (y-axis) or not using (x-axis) AS data on each wild-type (**A**) or variant (**B**) amino acid  
216 is shown in the scatter plots. The results are coloured according to the property of each amino acid type. Alanine  
217 (**A**) result is not applicable in the first figure since alanine scanning data are always missing when the wildtype is  
218 alanine itself. Absolute count for each amino acid can be found in Fig S14. (Neg.: negatively, Pos.: positively)

219

### 220 **3 Discussion**

221 In this study, we integrated alanine scanning (AS) data into deep mutational scanning (DMS)  
222 score prediction, leading to modest improvements in the accuracy of variant score prediction.  
223 We also explored the impact of the diversity of protein properties measured by DMS and AS.  
224 Filtering DMS and AS data based on our manual classification of assay type compatibility led  
225 to improved prediction performance.

226

227 A potential shortcoming of our current approach is that AS data were available for only a small  
228 proportion of the DMS data. Although most recent DMS studies can analyze variants of the  
229 whole protein, most AS experiments only cover a handful of residues in the target protein,  
230 leaving missing AS scores for the vast majority of residues. We explored this here and found  
231 that alternative methods for addressing the sparsity of AS data did not improve or degrade  
232 performance, but we anticipate further improved prediction accuracy if the low completeness  
233 and unevenness of AS data are appropriately handled before modelling.

234

235 In this study, we identified the importance of DMS/AS assay compatibility as a crucial factor  
236 for improving prediction accuracy. An issue with using this concept is that it further shrinks  
237 already sparse data. It also fails to take advantage of the fact that even for low compatible  
238 assays some fundamental information like protein abundance can still be mutually captured.

239 Instead of hard filtering, proper implementation of this underlying information may facilitate  
240 variant impact prediction in the future. Nonetheless, filtering on assay compatibility still leads  
241 to performance improvement. We also briefly explored whether the consistency of DMS and  
242 AS scores can be considered more directly by matching the best correlated AS data for each  
243 DMS dataset. Consistency is partially driven by assay compatibility but also reflects other fea-  
244 tures of the data, such as bias and noise.

245

246 The concepts of compatibility and data quality are also relevant to training any DMS-based  
247 predictors. DMS assays have been developed to measure variant impacts to distinct protein  
248 properties, and a variant can behave similarly to wildtype when measured by one assay yet  
249 show altered protein properties in other assay results, which are frequently found in regions  
250 with specific biochemical functions [25,52–56]. With more experimental assays to be applied,  
251 the diverse measurements may impede the progress of future DMS-based predictors unless this  
252 assay effect is properly addressed, for example, by building assay specific predictors. Meas-  
253 urement error is another source of DMS data heterogeneity that potentially affects the model  
254 performance. In our current study, DMS scores of protein variants are weighted equally while  
255 training. Adjustable weighting can be applied in future studies to adapt the distinct experi-  
256 mental error between individual variants and datasets, reducing the influence of low-confident  
257 data.

258

259 In summary, we conclude that the careful inclusion of low-throughput mutagenesis data im-  
260 proves the prediction of DMS scores, and the approaches described here can potentially be  
261 applied to other prediction methods.

262



263 **4 Availability of supporting source code and requirements**

264 **Project name:** DMS\_with\_Alanine\_scan

265 **Project home page:** [https://github.com/PapenfussLab/DMS\\_with\\_Alanine\\_scan](https://github.com/PapenfussLab/DMS_with_Alanine_scan)

266 **Operating system:** Platform independent

267 **Programming language:** Python

268 **Other requirements:** Python 3.10 or higher

269 **Licence:** MIT licence

270 **RRID:** SCR\_023949

271

272 **5 Data availability**

273 A copy of the data analysis code and a full set of data files required to reproduce this work  
274 are openly available in the GigaScience repository, GigaDB, under the record described in  
275 [57]. MaveDB accession numbers, UniProt accession numbers and other metadata describing  
276 the matched DMS-AS datasets are listed in Supplementary Table 1 (see Supporting infor-  
277 mation).

278

279 **6 List of abbreviations**

280 DMS: deep mutational scanning

281 AS: alanine scanning

282

283 **7 Supporting information**

284 **Supplementary Table 1:** All candidate DMS and alanine scanning data with detailed dataset  
285 information.

286 **Supplementary Table 2:** Normalized DMS dataset with protein property features.

287 **Supplementary Table 3:** Normalized alanine scanning dataset.

288

## 289 **8 Author contributions**

290 YF developed the software and wrote the initial draft of the manuscript. AFR conceived the  
291 study. JB, AFR, and ATP oversaw the project. All authors reviewed, contributed to, and ap-  
292 proved the manuscript.

293

## 294 **9 Funding**

295 YF is supported by Melbourne Research Scholarship. ATP was supported by an Australian Na-  
296 tional Health and Medical Research Council (NHMRC) Senior Research Fellowship (1116955).  
297 JB, AFR and ATP were supported by the Lorenzo and Pamela Galli Medical Research Trust.  
298 JB and ATP were supported by the Stafford Fox Medical Research Foundation. AFR was sup-  
299 ported by the National Human Genome Research Institute of the NIH under award numbers  
300 RM1HG010461 and UM1HG011969. The research benefitted from support from the Victorian  
301 State Government Operational Infrastructure Support and Australian Government NHMRC  
302 Independent Research Institute Infrastructure Support.

303

## 304 **10 Methods**

### 305 **10.1 DMS data collection**

306 DMS data were downloaded from MaveDB [40,41] which were then filtered and curated. DMS  
307 experiments targeting antibody and virus proteins were removed because of their potentially  
308 unique functionality. We retrieved the UniProt accession ID of target proteins by searching the  
309 protein names or sequences in UniProt [58], and proteins lacking available UniProt ID were  
310 also excluded. Datasets that are computationally processed or their wildtype-like and nonsense-

311 like scores (see Normalization) cannot be identified were also filtered out (Supplementary Ta-  
312 ble 1). All missense variants with only a single amino acid substitution were curated from the  
313 DMS studies for our analysis. A total of 130 DMS experiments from 53 studies [5,6,9–  
314 14,24,31,37–39,42,59–95] were collected for our analysis.

315

## 316 **10.2 Collection of AS data and other features**

317 The following process was used to search for candidate AS studies. Papers were identified by  
318 searching on PubMed and Google Scholar for the “alanine scan” or “alanine scanning” together  
319 with the name of candidate proteins. While searching in Google Scholar, we included the pro-  
320 tein’s UniProt ID rather than molecule name as the search term to reduce false positives. Ap-  
321 propriate AS data were collected from the search results. Western blot results were transformed  
322 to values by ImageJ if it was the only experimental data available in the study. A total 146 AS  
323 experiments were collected from 45 distinct studies [26–28,30,31,43–46,48,49,85,96–128].

324 Protein features of Shannon entropy and the logarithm of variant amino acid frequency were  
325 downloaded from the DeMaSk online toolkit [19]. The substitution score matrix feature was  
326 calculated from the mean of training DMS scores for each of the 380 possible amino acid sub-  
327 stitutions before each iteration of cross-validation.

328

## 329 **10.3 Normalization**

330 DMS and AS datasets were normalized to a common scale using the following approach  
331 adapted from previous studies [17,47]. Let  $D$  denotes a protein study measuring scores  $s_i^D$  for  
332 a single variant  $i$ ,  $s_{wt}^D$  denotes the scores for wildtype and  $s_{non}^D$  represents the score for non-  
333 sense-like variants. The normalized scores  $s_i'^D$  are given by:

$$334 \quad s_i'^D := \frac{s_i^D - s_{wt}^D}{s_{wt}^D - s_{non}^D} + 1$$

335 Wild-type scores were directly identified from the paper or the median score of synonymous  
336 variants. For DMS data, since not all DMS studies report score of nonsense variants, we defined  
337 the nonsense-like scores as the median DMS scores for the 1% missense variants with the  
338 strongest loss of function for each dataset. For AS data, nonsense-like scores were either de-  
339 fined according to the paper or using the extreme values (Supplementary Table 1).

340

#### 341 **10.4 AS data filtering and matching**

342 AS data subsets were filtered/matched according to either assay compatibility or score corre-  
343 lation. For assay compatibility filtering, we first categorized each DMS or AS assay by the  
344 protein property or function using the following assay types: binding affinity, enzyme activity,  
345 protein abundance, cell survival, pathogen infection, drug response, ability to perform a novel  
346 function, or other protein-specific activities (e.g., transcription activity for transcription factors)  
347 (Supplementary Table 1). The DMS/AS assay pairs were then classified into three levels of  
348 compatibility based on these categories (Fig S2). For each DMS dataset, we first tried to use  
349 only AS data with high assay compatibility for further modelling, removing AS data of medium  
350 and low assay compatibility. We then also tried to model with AS data of both high and medium  
351 assay compatibility.

352 For score correlation matching, Spearman's correlation ( $\rho$ ) is calculated between alanine sub-  
353 stitution scores in each pair of AS and DMS data. To avoid influence from the size of AS  
354 datasets, we estimated the  $\rho$  value with the empirical copula, which is related to the standard  
355 estimator by a factor of  $(n-1)/(n+1)$  [129,130]:

356 
$$\rho_r := \rho \times \frac{n-1}{n+1}$$

357 where  $\rho_r$  is the regularised correlation coefficient, and  $n$  is the number of alanine substitutions  
358 used for correlation calculation. For each DMS dataset, AS result with the highest  $\rho_r$  was  
359 picked for modelling.

360

### 361 **10.5 AS data pre-processing**

362 AS data were pre-processed prior to modelling. For variants without available (fil-  
363 tered/matched) AS data, their AS scores were imputed with the mean value of all available AS  
364 scores across all studies. Then the AS data were encoded by the wild-type and variant amino  
365 acid type with one-hot-encoding. For each variant, the AS feature is expanded with two one-  
366 hot vectors. Each of the vectors has 19 zeros and one non-zero value which was the AS score,  
367 with the location of the non-zero value indicating the wild-type or variant amino acid type.

368

### 369 **10.6 Training and evaluation of DMS score predictor**

370 To build the predictors, we performed linear regression using the function `sklearn.lin-`  
371 `ear_model.LinearRegression` from scikit-learn [131]. Training and validation data  
372 were separated with leave-one-protein-out cross-validation. In this process, data from one pro-  
373 tein were withheld for subsequent validation, and the rest were used for training. This process  
374 was iterated over all proteins in the data. Variants were inversely weighted during the training  
375 process by the number of measurements available, thus compensating for some regions having  
376 greater coverage with DMS and AS assays. Predictors were trained on protein features, DMS  
377 data and (optionally) AS data using four different filtering or matching strategies: i) all  
378 DMS/AS data, ii) compatibility-filtered DMS/AS data, iii) correlation-matched DMS/AS data,  
379 and iv) a control, constructed using DMS data only.

380 In the evaluation process, let  $V$  be protein variants assayed by both DMS study  $D$  and AS study  
381 A. Variant scores are predicted by the previously mentioned predictors either using AS data  
382 ( $\hat{s}_V^A$ ) or not ( $\hat{s}_V$ ). Spearman's correlation ( $\rho$ ) was calculated between the DMS scores  $s_V^D$  and  
383 each set of predicted scores. The difference of  $\rho$  was used to evaluate the performance change  
384 ( $\Delta\rho_V$ ).

$$385 \quad \rho_V^A = \text{Spearman's correlation}(\hat{s}_V^A, s_V^D)$$

$$386 \quad \rho_V = \text{Spearman's correlation}(\hat{s}_V, s_V^D)$$

$$387 \quad \Delta\rho_V = \rho_V^A - \rho_V$$

388 To evaluate, we iterated over variants from each pair of DMS/AS studies. Results were dropped  
389 for variants  $V$  with only one protein residue available during analysis and visualization. Model  
390 performance was compared using the following statistical tests. Results in Fig 5 & Fig S5 were  
391 tested with Welch's test, and results in Fig S4 & Fig S6 were tested with paired t-tests. The p-  
392 values were jointly corrected using the Holm-Šidák method. The 95% confidence interval of  
393 median values are calculated by Gaussian-based asymptotic approximation [132].

394

### 395 **10.7 Prediction with other variant effect predictors**

396 For PROVEAN [133] and SIFT [134], prediction results on target variants were directly down-  
397 loaded from the pre-calculated database for PROVEAN. For PolyPhen-2 [135] and GEMME  
398 [136], variant scores were computed through their online toolkits, using the default settings.  
399 ESM-1v [137] was set up locally and run according to its examples and documentations. EVE  
400 [138] results were collected from their pre-calculated database and a benchmarking study [139].

401

402 **11 References**

- 403 1. Fowler DM, Fields S. Deep mutational scanning: a new style of protein science. *Nature*  
404 *Methods*. 2014; doi: 10.1038/nmeth.3027.
- 405 2. Findlay GM. Linking genome variants to disease: scalable approaches to test the functional  
406 impact of human mutations. *Human Molecular Genetics*. 2021; doi: 10.1093/hmg/ddab219.
- 407 3. Geck RC, Boyle G, Amorosi CJ, Fowler DM, Dunham MJ. Measuring Pharmacogene Var-  
408 iant Function at Scale Using Multiplexed Assays. *Annual Review of Pharmacology and Toxi-*  
409 *cology*. 2022; doi: 10.1146/annurev-pharmtox-032221-085807.
- 410 4. Weile J, Roth FP. Multiplexed assays of variant effects contribute to a growing genotype-  
411 phenotype atlas. *Hum Genet*. 2018; doi: 10.1007/s00439-018-1916-x.
- 412 5. Diss G, Lehner B. The genetic landscape of a physical interaction. *eLife*. 2018; doi:  
413 10.7554/eLife.32472.
- 414 6. Fowler DM, Araya CL, Fleishman SJ, Kellogg EH, Stephany JJ, Baker D, et al.. High-reso-  
415 lution mapping of protein sequence-function relationships. *Nature Methods*. 2010; doi:  
416 10.1038/nmeth.1492.
- 417 7. Amorosi CJ, Chiasson MA, McDonald MG, Wong LH, Sitko KA, Boyle G, et al.. Massively  
418 parallel characterization of CYP2C9 variant enzyme activity and abundance. *The American*  
419 *Journal of Human Genetics*. 2021; doi: 10.1016/j.ajhg.2021.07.001.
- 420 8. Faure AJ, Domingo J, Schmiedel JM, Hidalgo-Carcedo C, Diss G, Lehner B. Mapping the  
421 energetic and allosteric landscapes of protein binding domains. *Nature*. 2022; doi:  
422 10.1038/s41586-022-04586-4.

**Commented [NN1]:** Authors: please go through the refer-  
ences and check all preprint citations - they are missing  
DOIs. Full citations are required including DOIs.

If the preprint has already been published - the full journal  
citation needs to be cited, instead of the preprint.

Note I have added in Ref #140 - the GigaDB DOI Citation.

**Commented [AR2R1]:** Thank you for this. We have up-  
dated two of the six preprints to their published versions and  
added DOIs to the remaining four preprint references (three  
biorxiv one arxiv).

- 423 9. Matreyek KA, Starita LM, Stephany JJ, Martin B, Chiasson MA, Gray VE, et al.. Multiplex  
424 assessment of protein variant abundance by massively parallel sequencing. *Nature Genetics*.  
425 2018; doi: 10.1038/s41588-018-0122-z.
- 426 10. Mighell TL, Evans-Dutson S, O’Roak BJ. A Saturation Mutagenesis Approach to Under-  
427 standing PTEN Lipid Phosphatase Activity and Genotype-Phenotype Relationships. *The Amer-  
428 ican Journal of Human Genetics*. 2018; doi: 10.1016/j.ajhg.2018.03.018.
- 429 11. Stiffler MA, Hekstra DR, Ranganathan R. Evolvability as a Function of Purifying Selection  
430 in TEM-1  $\beta$ -Lactamase. *Cell*. 2015; doi: 10.1016/j.cell.2015.01.035.
- 431 12. Ahler E, Register AC, Chakraborty S, Fang L, Dieter EM, Sitko KA, et al.. A Combined  
432 Approach Reveals a Regulatory Mechanism Coupling Src’s Kinase Activity, Localization, and  
433 Phosphotransferase-Independent Functions. *Molecular Cell*. 2019; doi: 10.1016/j.mol-  
434 cel.2019.02.003.
- 435 13. Giacomelli AO, Yang X, Lintner RE, McFarland JM, Duby M, Kim J, et al.. Mutational  
436 processes shape the landscape of TP53 mutations in human cancer. *Nature Genetics*. Nature  
437 Publishing Group; 2018; doi: 10.1038/s41588-018-0204-y.
- 438 14. Roscoe BP, Thayer KM, Zeldovich KB, Fushman D, Bolon DNA. Analyses of the Effects  
439 of All Ubiquitin Point Mutants on Yeast Growth Rate. *Journal of Molecular Biology*. 2013;  
440 doi: 10.1016/j.jmb.2013.01.032.
- 441 15. Tabet D, Parikh V, Mali P, Roth FP, Claussnitzer M. Scalable Functional Assays for the  
442 Interpretation of Human Genetic Variation. *Annu Rev Genet*. 2022; doi: 10.1146/annurev-  
443 genet-072920-032107.



- 444 16. Kuang D, Weile J, Kishore N, Nguyen M, Rubin AF, Fields S, et al.. MaveRegistry: a  
445 collaboration platform for multiplexed assays of variant effect. Lu Z, editor. *Bioinformatics*.  
446 2021; doi: 10.1093/bioinformatics/btab215.
- 447 17. Gray VE, Hause RJ, Luebeck J, Shendure J, Fowler DM. Quantitative Missense Variant  
448 Effect Prediction Using Large-Scale Mutagenesis Data. *Cell Systems*. 2018; doi:  
449 10.1016/j.cels.2017.11.003.
- 450 18. Alley EC, Khimulya G, Biswas S, AlQuraishi M, Church GM. Unified rational protein  
451 engineering with sequence-based deep representation learning. *Nat Methods*. 2019; doi:  
452 10.1038/s41592-019-0598-1.
- 453 19. Munro D, Singh M. DeMaSk: a deep mutational scanning substitution matrix and its use  
454 for variant impact prediction. Xu J, editor. *Bioinformatics*. 2020; doi: 10.1093/bioinformat-  
455 ics/btaa1030.
- 456 20. Biswas S, Khimulya G, Alley EC, Esvelt KM, Church GM. Low- N protein engineering  
457 with data-efficient deep learning. *Nature Methods*. Nature Publishing Group; 2021; doi:  
458 10.1038/s41592-021-01100-y.
- 459 21. Høie MH, Cagiada M, Beck Frederiksen AH, Stein A, Lindorff-Larsen K. Predicting and  
460 interpreting large-scale mutagenesis data using analyses of protein stability and conservation.  
461 *Cell Reports*. 2022; doi: 10.1016/j.celrep.2021.110207.
- 462 22. Wu Y, Li R, Sun S, Weile J, Roth FP. Improved pathogenicity prediction for rare human  
463 missense variants. *The American Journal of Human Genetics*. 2021; doi:  
464 10.1016/j.ajhg.2021.08.012.

- 465 23. Hsu C, Nisonoff H, Fannjiang C, Listgarten J. Learning protein fitness models from evolu-  
466 tionary and assay-labeled data. *Nat Biotechnol.* 2022; doi: 10.1038/s41587-021-01146-5.
- 467 24. Findlay GM, Daza RM, Martin B, Zhang MD, Leith AP, Gasperini M, et al.. Accurate  
468 classification of BRCA1 variants with saturation genome editing. *Nature.* 2018; doi:  
469 10.1038/s41586-018-0461-z.
- 470 25. Cagiada M, Bottaro S, Lindemose S, Schenstrøm SM, Stein A, Hartmann-Petersen R, et  
471 al.. Discovering functionally important sites in proteins. *Nat Commun.* Nature Publishing  
472 Group; 2023; doi: 10.1038/s41467-023-39909-0.
- 473 26. Block C, Janknecht R, Herrmann C, Nassar N, Wittinghofer A. Quantitative structure-ac-  
474 tivity analysis correlating Ras/Raf interaction in vitro to Raf activation in vivo. *Nature Struc-*  
475 *tural Biology.* Nature Publishing Group; 1996; doi: 10.1038/nsb0396-244.
- 476 27. Sloan DJ, Hellinga HW. Dissection of the protein G B1 domain binding site for human IgG  
477 Fc fragment. *Protein Science.* 1999; doi: 10.1110/ps.8.8.1643.
- 478 28. Fleming KG, Engelman DM. Specificity in transmembrane helix–helix interactions can  
479 define a hierarchy of stability for sequence variants. *PNAS.* National Academy of Sciences;  
480 2001; doi: 10.1073/pnas.251367498.
- 481 29. Shibata Y, White JF, Serrano-Vega MJ, Magnani F, Aloia AL, Grishammer R, et al..  
482 Thermostabilization of the Neurotensin Receptor NTS1. *Journal of Molecular Biology.* 2009;  
483 doi: 10.1016/j.jmb.2009.04.068.

484 30. Brzovic PS, Heikaus CC, Kisselev L, Vernon R, Herbig E, Pacheco D, et al.. The Acidic  
485 Transcription Activator Gcn4 Binds the Mediator Subunit Gal11/Med15 Using a Simple Pro-  
486 tein Interface Forming a Fuzzy Complex. *Molecular Cell*. 2011; doi: 10.1016/j.mol-  
487 cel.2011.11.008.

488 31. Gajula KS, Huwe PJ, Mo CY, Crawford DJ, Stivers JT, Radhakrishnan R, et al.. High-  
489 throughput mutagenesis reveals functional determinants for DNA targeting by activation-in-  
490 duced deaminase. *Nucleic Acids Research*. 2014; doi: 10.1093/nar/gku689.

491 32. Kortemme T, Kim DE, Baker D. Computational Alanine Scanning of Protein-Protein In-  
492 terfaces. *Science's STKE*. American Association for the Advancement of Science; 2004; doi:  
493 10.1126/stke.2192004pl2.

494 33. Morrison KL, Weiss GA. Combinatorial alanine-scanning. *Current Opinion in Chemical*  
495 *Biology*. 2001; doi: 10.1016/S1367-5931(00)00206-4.

496 34. Cunningham BC, Wells JA. High-resolution epitope mapping of hGH-receptor interactions  
497 by alanine-scanning mutagenesis. *Science*. American Association for the Advancement of Sci-  
498 ence; 1989; doi: 10.1126/science.2471267.

499 35. DeLano WL. Unraveling hot spots in binding interfaces: progress and challenges. *Current*  
500 *Opinion in Structural Biology*. 2002; doi: 10.1016/S0959-440X(02)00283-X.

501 36. Eustache S, Leprince J, Tufféry P. Progress with peptide scanning to study structure-activ-  
502 ity relationships: the implications for drug discovery. *Expert Opinion on Drug Discovery*. 2016;  
503 doi: 10.1080/17460441.2016.1201058.

- 504 37. Olson CA, Wu NC, Sun R. A Comprehensive Biophysical Description of Pairwise Epistasis  
505 throughout an Entire Protein Domain. *Current Biology*. 2014; doi: 10.1016/j.cub.2014.09.072.
- 506 38. Staller MV, Holehouse AS, Swain-Lenz D, Das RK, Pappu RV, Cohen BA. A High-  
507 Throughput Mutational Scan of an Intrinsically Disordered Acidic Transcriptional Activation  
508 Domain. *Cell Systems*. 2018; doi: 10.1016/j.cels.2018.01.015.
- 509 39. Gray VE, Sitko K, Kameni FZN, Williamson M, Stephany JJ, Hasle N, et al.. Elucidating  
510 the Molecular Determinants of A $\beta$  Aggregation with Deep Mutational Scanning. *G3 (Be-*  
511 *thesda)*. 2019; doi: 10.1534/g3.119.400535.
- 512 40. Esposito D, Weile J, Shendure J, Starita LM, Papenfuss AT, Roth FP, et al.. MaveDB: an  
513 open-source platform to distribute and interpret data from multiplexed assays of variant effect.  
514 *Genome Biol*. 2019; doi: 10.1186/s13059-019-1845-6.
- 515 41. Rubin AF, Min JK, Rollins NJ, Da EY, Esposito D, Harrington M, et al.. MaveDB v2: a  
516 curated community database with over three million variant effects from multiplexed func-  
517 tional assays. bioRxiv; doi: 10.1101/2021.11.29.470445.
- 518 42. Heredia JD, Park J, Brubaker RJ, Szymanski SK, Gill KS, Procko E. Mapping Interaction  
519 Sites on Human Chemokine Receptors by Deep Mutational Scanning. *The Journal of Immu-*  
520 *nology*. American Association of Immunologists; 2018; doi: 10.4049/jimmunol.1800343.
- 521 43. Tian S, Choi W-T, Liu D, Pesavento J, Wang Y, An J, et al.. Distinct Functional Sites for  
522 Human Immunodeficiency Virus Type 1 and Stromal Cell-Derived Factor 1 $\alpha$  on CXCR4  
523 Transmembrane Helical Domains. *JVI*. 2005; doi: 10.1128/JVI.79.20.12667-12673.2005.

524 44. Chabot DJ, Zhang P-F, Quinnan GV, Broder CC. Mutagenesis of CXCR4 Identifies Im-  
525 portant Domains for Human Immunodeficiency Virus Type 1 X4 Isolate Envelope-Mediated  
526 Membrane Fusion and Virus Entry and Reveals Cryptic Coreceptor Activity for R5 Isolates. *J*  
527 *Virol.* 1999; doi: 10.1128/JVI.73.8.6598-6609.1999.

528 45. Han DP, Penn-Nicholson A, Cho MW. Identification of critical determinants on ACE2 for  
529 SARS-CoV entry and development of a potent entry inhibitor. *Virology.* 2006; doi:  
530 10.1016/j.virol.2006.01.029.

531 46. Fujita-Yoshigaki J, Shirouzu M, Ito Y, Hattori S, Furuyama S, Nishimura S, et al.. A Con-  
532 stitutive Effector Region on the C-terminal Side of Switch I of the Ras Protein. *J Biol Chem.*  
533 American Society for Biochemistry and Molecular Biology; 1995; doi: 10.1074/jbc.270.9.4661.

534 47. Gray VE, Hause RJ, Fowler DM. Analysis of Large-Scale Mutagenesis Data To Assess the  
535 Impact of Single Amino Acid Substitutions. *Genetics.* 2017; doi: 10.1534/genetics.117.300064.

536 48. Hidalgo P, Ansari AZ, Schmidt P, Hare B, Simkovich N, Farrell S, et al.. Recruitment of  
537 the transcriptional machinery through GAL11P: structure and interactions of the GAL4 dimer-  
538 ization domain. *Genes Dev.* 2001; doi: 10.1101/gad.873901.

539 49. Rodríguez-Escudero I, Oliver MD, Andrés-Pons A, Molina M, Cid VJ, Pulido R. A com-  
540 prehensive functional analysis of PTEN mutations: implications in tumor- and autism-related  
541 syndromes. *Human Molecular Genetics.* 2011; doi: 10.1093/hmg/ddr337.

542 50. Schröter C, Günther R, Rhiel L, Becker S, Toleikis L, Doerner A, et al.. A generic approach  
543 to engineer antibody pH-switches using combinatorial histidine scanning libraries and yeast  
544 display. *mAbs.* 2015; doi: 10.4161/19420862.2014.985993.

545 51. Starace DM, Bezanilla F. Histidine Scanning Mutagenesis of Basic Residues of the S4  
546 Segment of the Shaker K<sup>+</sup> Channel. *J Gen Physiol.* 117:469–902001;

547 52. Cagiada M, Johansson KE, Valanciute A, Nielsen SV, Hartmann-Petersen R, Yang JJ, et  
548 al.. Understanding the Origins of Loss of Protein Function by Analyzing the Effects of Thou-  
549 sands of Variants on Activity and Abundance. Ozkan B, editor. *Molecular Biology and Evolu-*  
550 *tion.* 2021; doi: 10.1093/molbev/msab095.

551 53. Jepsen MM, Fowler DM, Hartmann-Petersen R, Stein A, Lindorff-Larsen K. Chapter 5 -  
552 Classifying disease-associated variants using measures of protein activity and stability. In: Pey  
553 AL, editor. *Protein Homeostasis Diseases.* Academic Press;

554 54. Matreyek KA, Stephany JJ, Ahler E, Fowler DM. Integrating thousands of PTEN variant  
555 activity and abundance measurements reveals variant subgroups and new dominant negatives  
556 in cancers. *Genome Med.* 2021; doi: 10.1186/s13073-021-00984-x.

557 55. Mighell TL, Thacker S, Fombonne E, Eng C, O’Roak BJ. An Integrated Deep-Mutational-  
558 Scanning Approach Provides Clinical Insights on PTEN Genotype-Phenotype Relationships.  
559 *The American Journal of Human Genetics.* 2020; doi: 10.1016/j.ajhg.2020.04.014.

560 56. Nielsen SV, Hartmann-Petersen R, Stein A, Lindorff-Larsen K. Multiplexed assays reveal  
561 effects of missense variants in MSH2 and cancer predisposition. *PLOS Genetics.* Public Li-  
562 brary of Science; 2021; doi: 10.1371/journal.pgen.1009496.

563 57. Fu Y, Bedö J, Papenfuss AT, Rubin AF. Supporting data for “Integrating deep mutational  
564 scanning and low-throughput mutagenesis data to predict the impact of amino acid variants.”  
565 GigaScience Database. 2023. <http://dx.doi.org/10.5524/102429>.

566 58. The UniProt Consortium, Bateman A, Martin M-J, Orchard S, Magrane M, Agivetova R,  
567 et al.. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Research*. 2021;  
568 doi: 10.1093/nar/gkaa1100.

569 59. Andrews B, Fields S. Distinct patterns of mutational sensitivity for  $\lambda$  resistance and malto-  
570 dextrin transport in *Escherichia coli* LamB. *Microb Genom*. 2020; doi:  
571 10.1099/mgen.0.000364.

572 60. Bandaru P, Shah NH, Bhattacharyya M, Barton JP, Kondo Y, Cofsky JC, et al.. Decon-  
573 struction of the Ras switching cycle through saturation mutagenesis. *eLife*. 2017; doi:  
574 10.7554/eLife.27810.

575 61. Bolognesi B, Faure AJ, Seuma M, Schmiedel JM, Tartaglia GG, Lehner B. The mutational  
576 landscape of a prion-like domain. *Nat Commun*. 2019; doi: 10.1038/s41467-019-12101-z.

577 62. Bridgford JL, Lee SM, Lee CMM, Guglielmelli P, Rumi E, Pietra D, et al.. Novel drivers  
578 and modifiers of MPL-dependent oncogenic transformation identified by deep mutational scan-  
579 ning. *Blood*. American Society of Hematology; 2020; doi: 10.1182/blood.2019002561.

580 63. Chan KK, Dorosky D, Sharma P, Abbasi SA, Dye JM, Kranz DM, et al.. Engineering hu-  
581 man ACE2 to optimize binding to the spike protein of SARS coronavirus 2. *Science*. American  
582 Association for the Advancement of Science; 2020; doi: 10.1126/science.abc0870.

583 64. Chiasson MA, Rollins NJ, Stephany JJ, Sitko KA, Matreyek KA, Verby M, et al.. Multi-  
584 plexed measurement of variant abundance and activity reveals VKOR topology, active site and  
585 human variant impact. *Elife*. 2020; doi: 10.7554/eLife.58026.

586 65. Elazar A, Weinstein J, Biran I, Fridman Y, Bibi E, Fleishman SJ. Mutational scanning  
587 reveals the determinants of protein insertion and association energetics in the plasma mem-  
588 brane. Shan Y, editor. *eLife*. eLife Sciences Publications, Ltd; 2016; doi: 10.7554/eLife.12125.

589 66. Firnberg E, Labonte JW, Gray JJ, Ostermeier M. A Comprehensive, High-Resolution Map  
590 of a Gene's Fitness Landscape. *Mol Biol Evol*. 2014; doi: 10.1093/molbev/msu081.

591 67. Hietpas RT, Jensen JD, Bolon DNA. Experimental illumination of a fitness landscape. *Pro-*  
592 *ceedings of the National Academy of Sciences*. 2011; doi: 10.1073/pnas.1016024108.

593 68. Hietpas RT, Bank C, Jensen JD, Bolon DNA. Shifting fitness landscapes in response to  
594 altered environments. *Evolution*. 2013; doi: 10.1111/evo.12207.

595 69. Jiang L, Mishra P, Hietpas RT, Zeldovich KB, Bolon DNA. Latent Effects of Hsp90 Mu-  
596 tants Revealed at Reduced Expression Levels. *PLOS Genetics*. Public Library of Science; 2013;  
597 doi: 10.1371/journal.pgen.1003600.

598 70. Jiang RJ. Exhaustive Mapping of Missense Variation in Coronary Heart Disease-related  
599 Genes [Thesis]. University of Toronto;

600 71. Keskin A, Akdoğan E, Dunn CD. Evidence for Amino Acid Snorkeling from a High-Res-  
601 olution, *In Vivo* Analysis of Fis1 Tail-Anchor Insertion at the Mitochondrial Outer Membrane.  
602 *Genetics*. 2017; doi: 10.1534/genetics.116.196428.

603 72. Kitzman JO, Starita LM, Lo RS, Fields S, Shendure J. Massively parallel single-amino-  
604 acid mutagenesis. *Nat Methods*. 2015; doi: 10.1038/nmeth.3223.

605 73. Kotler E, Shani O, Goldfeld G, Lotan-Pompan M, Tarcic O, Gershoni A, et al.. A System-  
606 atic p53 Mutation Library Links Differential Functional Impact to Cancer Mutation Pattern and



607 Evolutionary Conservation. *Molecular Cell*. Elsevier; 2018; doi: 10.1016/j.mol-  
608 cel.2018.06.012.

609 74. Kowalsky CA, Whitehead TA. Determination of binding affinity upon mutation for type I  
610 dockerin-cohesin complexes from *Clostridium thermocellum* and *Clostridium cellulolyticum*  
611 using deep sequencing. *Proteins: Structure, Function, and Bioinformatics*. 2016; doi:  
612 10.1002/prot.25175.

613 75. McLaughlin Jr RN, Poelwijk FJ, Raman A, Gosal WS, Ranganathan R. The spatial archi-  
614 tecture of protein function and adaptation. *Nature*. 2012; doi: 10.1038/nature11500.

615 76. Melamed D, Young DL, Gamble CE, Miller CR, Fields S. Deep mutational scanning of an  
616 RRM domain of the *Saccharomyces cerevisiae* poly(A)-binding protein. *RNA*. 2013; doi:  
617 10.1261/rna.040709.113.

618 77. Mishra P, Flynn JM, Starr TN, Bolon DNA. Systematic Mutant Analyses Elucidate General  
619 and Client-Specific Aspects of Hsp90 Function. *Cell Reports*. 2016; doi:  
620 10.1016/j.celrep.2016.03.046.

621 78. Nedrud D, Coyote-Maestas W, Schmidt D. A large-scale survey of pairwise epistasis re-  
622 veals a mechanism for evolutionary expansion and specialization of PDZ domains. *Proteins:  
623 Structure, Function, and Bioinformatics*. 2021; doi: 10.1002/prot.26067.

624 79. Newberry RW, Arhar T, Costello J, Hartoularos GC, Maxwell AM, Naing ZZC, et al..  
625 Robust Sequence Determinants of  $\alpha$ -Synuclein Toxicity in Yeast Implicate Membrane Binding.  
626 *ACS Chem Biol*. 2020; doi: 10.1021/acscembio.0c00339.

627 80. Newberry RW, Leong JT, Chow ED, Kampmann M, DeGrado WF. Deep mutational scan-  
628 ning reveals the structural basis for  $\alpha$ -synuclein activity. *Nat Chem Biol.* 2020; doi:  
629 10.1038/s41589-020-0480-6.

630 81. Roscoe BP, Bolon DNA. Systematic Exploration of Ubiquitin Sequence, E1 Activation  
631 Efficiency, and Experimental Fitness in Yeast. *Journal of Molecular Biology.* 2014; doi:  
632 10.1016/j.jmb.2014.05.019.

633 82. Sarkisyan KS, Bolotin DA, Meer MV, Usmanova DR, Mishin AS, Sharonov GV, et al..  
634 Local fitness landscape of the green fluorescent protein. *Nature.* Nature Publishing Group;  
635 2016; doi: 10.1038/nature17995.

636 83. Silverstein RA, Sun S, Verby M, Weile J, Wu Y, Gebbia M, et al.. A systematic genotype-  
637 phenotype map for missense variants in the human intellectual disability-associated gene GDI1.  
638 bioRxiv; doi: 10.1101/2021.10.06.463360.

639 84. Starita LM, Pruneda JN, Lo RS, Fowler DM, Kim HJ, Hiatt JB, et al.. Activity-enhancing  
640 mutations in an E3 ubiquitin ligase identified by high-throughput mutagenesis. *PNAS.* 2013;  
641 doi: 10.1073/pnas.1303309110.

642 85. Starita LM, Young DL, Islam M, Kitzman JO, Gullingsrud J, Hause RJ, et al.. Massively  
643 Parallel Functional Analysis of BRCA1 RING Domain Variants. *Genetics.* 2015; doi:  
644 10.1534/genetics.115.175802.

645 86. Starita LM, Islam MM, Banerjee T, Adamovich AI, Gullingsrud J, Fields S, et al.. A Mul-  
646 tiplex Homology-Directed DNA Repair Assay Reveals the Impact of More Than 1,000 BRCA1  
647 Missense Substitution Variants on Protein Function. *The American Journal of Human Genetics.*  
648 2018; doi: 10.1016/j.ajhg.2018.07.016.

649 87. Suiter CC, Moriyama T, Matreyek KA, Yang W, Scaletti ER, Nishii R, et al.. Massively  
650 parallel variant characterization identifies *NUDT15* alleles associated with thiopurine toxicity.  
651 *Proc Natl Acad Sci USA*. 2020; doi: 10.1073/pnas.1915680117.

652 88. Sun S, Weile J, Verby M, Wu Y, Wang Y, Cote AG, et al.. A proactive genotype-to-patient-  
653 phenotype map for cystathionine beta-synthase. *Genome Med*. 2020; doi: 10.1186/s13073-020-  
654 0711-1.

655 89. Thompson S, Zhang Y, Ingle C, Reynolds KA, Kortemme T. Altered expression of a quality  
656 control protease in *E. coli* reshapes the in vivo mutational landscape of a model enzyme. *eLife*.  
657 2020; doi: 10.7554/eLife.53476.

658 90. Trenker R, Wu X, Nguyen JV, Wilcox S, Rubin AF, Call ME, et al.. Human and viral  
659 membrane-associated E3 ubiquitin ligases MARCH1 and MIR2 recognize different features  
660 of CD86 to downregulate surface expression. *Journal of Biological Chemistry*. Elsevier; 2021;  
661 doi: 10.1016/j.jbc.2021.100900.

662 91. Weile J, Sun S, Cote AG, Knapp J, Verby M, Mellor JC, et al.. A framework for exhaust-  
663 ively mapping functional missense variants. *Mol Syst Biol*. 2017; doi: 10.15252/msb.20177908.

664 92. Weile J, Kishore N, Sun S, Maaieh R, Verby M, Li R, et al.. Shifting landscapes of human  
665 MTHFR missense-variant effects. *The American Journal of Human Genetics*. Elsevier; 2021;  
666 doi: 10.1016/j.ajhg.2021.05.009.

667 93. Wrenbeck EE, Bedewitz MA, Klesmith JR, Noshin S, Barry CS, Whitehead TA. An Auto-  
668 mated Data-Driven Pipeline for Improving Heterologous Enzyme Expression. *ACS Synth Biol*.  
669 American Chemical Society; 2019; doi: 10.1021/acssynbio.8b00486.

670 94. Zhang L, Sarangi V, Moon I, Yu J, Liu D, Devarajan S, et al.. CYP2C9 and CYP2C19:  
671 Deep Mutational Scanning and Functional Characterization of Genomic Missense Variants.  
672 *Clinical and Translational Science*. 2020; doi: <https://doi.org/10.1111/cts.12758>.

673 95. Zinkus-Boltz J, DeValk C, Dickinson BC. A Phage-Assisted Continuous Selection Ap-  
674 proach for Deep Mutational Scanning of Protein–Protein Interactions. *ACS Chem Biol*. Amer-  
675 ican Chemical Society; 2019; doi: 10.1021/acscchembio.9b00669.

676 96. Bernier-Villamor V, Sampson DA, Matunis MJ, Lima CD. Structural Basis for E2-Medi-  
677 ated SUMO Conjugation Revealed by a Complex between Ubiquitin-Conjugating Enzyme  
678 Ubc9 and RanGAP. *Cell*. 108:122002;

679 97. Blanpain C, Doranz BJ, Vakili J, Rucker J, Govaerts C, Baik SSW, et al.. Multiple Charged  
680 and Aromatic Residues in CCR5 Amino-terminal Domain Are Involved in High Affinity Bind-  
681 ing of Both Chemokines and HIV-1 Env Protein. *J Biol Chem*. 1999; doi:  
682 10.1074/jbc.274.49.34719.

683 98. Brzovic PS, Keefe JR, Nishikawa H, Miyamoto K, Fox D, Fukuda M, et al.. Binding and  
684 recognition in the assembly of an active BRCA1/BARD1 ubiquitin-ligase complex. *Proceed-*  
685 *ings of the National Academy of Sciences*. 2003; doi: 10.1073/pnas.0836054100.

686 99. Chen S, Wu J, Zhong S, Li Y, Zhang P, Ma J, et al.. iASPP mediates p53 selectivity through  
687 a modular mechanism fine-tuning DNA recognition. *Proc Natl Acad Sci USA*. 2019; doi:  
688 10.1073/pnas.1909393116.

689 100. Chupreta S, Holmstrom S, Subramanian L, Iñiguez-Lluhí JA. A Small Conserved Surface  
690 in SUMO Is the Critical Structural Determinant of Its Transcriptional Inhibitory Properties.  
691 *MCB*. 2005; doi: 10.1128/MCB.25.10.4272-4282.2005.

692 101. Cobb JA, Roberts DM. Structural Requirements for N-Trimethylation of Lysine 115 of  
693 Calmodulin. *Journal of Biological Chemistry*. 2000; doi: 10.1074/jbc.M002332200.

694 102. Coyne RS, McDonald HB, Edgemon K, Brody LC. Functional Characterization of  
695 BRCA1 Sequence Variants using a Yeast Small Colony Phenotype Assay. *Cancer Biology &*  
696 *Therapy*. 2004; doi: 10.4161/cbt.3.5.809.

697 103. Denker K, Orlik F, Schiffler B, Benz R. Site-directed Mutagenesis of the Greasy Slide  
698 Aromatic Residues Within the LamB (Maltoporin) Channel of Escherichia coli: Effect on Ion  
699 and Maltopentaose Transport. *Journal of Molecular Biology*. 2005; doi:  
700 10.1016/j.jmb.2005.07.025.

701 104. Dragic T, Trkola A, Lin SW, Nagashima KA, Kajumo F, Zhao L, et al.. Amino-Terminal  
702 Substitutions in the CCR5 Coreceptor Impair gp120 Binding and Human Immunodeficiency  
703 Virus Type 1 Entry. *J Virol*. 1998; doi: 10.1128/JVI.72.1.279-285.1998.

704 105. Dragic T, Trkola A, Thompson DAD, Cormier EG, Kajumo FA, Maxwell E, et al.. A  
705 binding pocket for a small molecule inhibitor of HIV-1 entry within the transmembrane helices  
706 of CCR5. *Proceedings of the National Academy of Sciences*. 2000; doi:  
707 10.1073/pnas.090576697.

708 106. Ecsédi P, Gógl G, Hóf H, Kiss B, Harmat V, Nyitray L. Structure Determination of the  
709 Transactivation Domain of p53 in Complex with S100A4 Using Annexin A2 as a Crystalliza-  
710 tion Chaperone. *Structure*. 2020; doi: 10.1016/j.str.2020.05.001.

711 107. Kopecká J, Krijt J, Raková K, Kožich V. Restoring assembly and activity of cystathionine  
712  $\beta$ -synthase mutants by ligands and chemical chaperones. *Journal of Inherited Metabolic Dis-*  
713 *ease*. 2011; doi: 10.1007/s10545-010-9087-5.

714 108. Kožich V, Sokolová J, Klatovská V, Krijt J, Janošik M, Jelínek K, et al.. Cystathionine  $\beta$ -  
715 synthase mutations: effect of mutation topology on folding and activity. *Hum Mutat.* 2010; doi:  
716 10.1002/humu.21273.

717 109. Kruger W d., Wang L, Jhee K h., Singh R h., Elsas II L j.. Cystathionine  $\beta$ -synthase defi-  
718 ciency in Georgia (USA): Correlation of clinical and biochemical phenotype with genotype.  
719 *Human Mutation.* 2003; doi: 10.1002/humu.10290.

720 110. Lee SY, Pullen L, Virgil DJ, Castañeda CA, Abeykoon D, Bolon DNA, et al.. Alanine  
721 Scan of Core Positions in Ubiquitin Reveals Links between Dynamics, Stability, and Function.  
722 *Journal of Molecular Biology.* 2014; doi: 10.1016/j.jmb.2013.10.042.

723 111. Li W, Zhang C, Sui J, Kuhn JH, Moore MJ, Luo S, et al.. Receptor and viral determinants  
724 of SARS-coronavirus adaptation to human ACE2. *EMBO J.* 2005; doi: 10.1038/sj.em-  
725 boj.7600640.

726 112. Lin G, Baribaud F, Romano J, Doms RW, Hoxie JA. Identification of gp120 Binding Sites  
727 on CXCR4 by Using CD4-Independent Human Immunodeficiency Virus Type 2 Env Proteins.  
728 *JVI.* 2003; doi: 10.1128/JVI.77.2.931-942.2003.

729 113. Mascle XH, Lussier-Price M, Cappadocia L, Estephan P, Raiola L, Omichinski JG, et al..  
730 Identification of a Non-covalent Ternary Complex Formed by PIAS1, SUMO1, and UBC9  
731 Proteins Involved in Transcriptional Regulation. *Journal of Biological Chemistry.* 2013; doi:  
732 10.1074/jbc.M113.486845.

733 114. Matthews EE, Thévenin D, Rogers JM, Gotow L, Lira PD, Reiter LA, et al.. Thrombo-  
734 poietin receptor activation: transmembrane helix dimerization, rotation, and allosteric modula-  
735 tion. *The FASEB Journal.* 2011; doi: <https://doi.org/10.1096/fj.10-178673>.

736 115. Mayfield JA, Davies MW, Dimster-Denk D, Pleskac N, McCarthy S, Boydston EA, et al..  
737 Surrogate Genetics and Metabolic Profiling for Characterization of Human Disease Alleles.  
738 *Genetics*. 2012; doi: 10.1534/genetics.111.137471.

739 116. Navenot J-M, Wang Z, Trent JO, Murray JL, Hu Q, DeLeeuw L, et al.. Molecular anatomy  
740 of CCR5 engagement by physiologic and viral chemokines and HIV-1 envelope glycoproteins:  
741 differences in primary structural requirements for RANTES, MIP-1 $\alpha$ , and vMIP-II bind-  
742 ing11Edited by P. E. Wright. *Journal of Molecular Biology*. 2001; doi:  
743 10.1006/jmbi.2001.5086.

744 117. Peng L, Damschroder MM, Cook KE, Wu H, Dall'Acqua WF. Molecular basis for the  
745 antagonistic activity of an anti-CXCR4 antibody. *mAbs*. 2016; doi:  
746 10.1080/19420862.2015.1113359.

747 118. Peterson BR, Sun LJ, Verdine GL. A critical arginine residue mediates cooperativity in  
748 the contact interface between transcription factors NFAT and AP-1. *Proceedings of the Na-  
749 tional Academy of Sciences*. 1996; doi: 10.1073/pnas.93.24.13671.

750 119. Rabut GEE, Konner JA, Kajumo F, Moore JP, Dragic T. Alanine Substitutions of Polar  
751 and Nonpolar Residues in the Amino-Terminal Domain of CCR5 Differently Impair Entry of  
752 Macrophage- and Dualtropic Isolates of Human Immunodeficiency Virus Type 1. *J Virol*. 1998;  
753 doi: 10.1128/JVI.72.4.3464-3468.1998.

754 120. Ransburgh DJR, Chiba N, Ishioka C, Toland AE, Parvin JD. Identification of Breast Tu-  
755 mor Mutations in *BRCA1* That Abolish Its Function in Homologous DNA Recombination.  
756 *Cancer Res*. 2010; doi: 10.1158/0008-5472.CAN-09-2850.

757 121. Tan Y, Tong P, Wang J, Zhao L, Li J, Yu Y, et al.. The Membrane-Proximal Region of  
758 C–C Chemokine Receptor Type 5 Participates in the Infection of HIV-1. *Front Immunol.* 2017;  
759 doi: 10.3389/fimmu.2017.00478.

760 122. Towler WI, Zhang J, Ransburgh DJR, Toland AE, Ishioka C, Chiba N, et al.. Analysis of  
761 BRCA1 Variants in Double-Strand Break Repair by Homologous Recombination and Single-  
762 Strand Annealing. *Human Mutation.* 2013; doi: 10.1002/humu.22251.

763 123. Trent JO, Wang Z, Murray JL, Shao W, Tamamura H, Fujii N, et al.. Lipid Bilayer Sim-  
764 ulations of CXCR4 with Inverse Agonists and Weak Partial Agonists. *J Biol Chem.* 2003; doi:  
765 10.1074/jbc.M307850200.

766 124. Van Gelder P, Dumas F, Bartoldus I, Saint N, Prilipov A, Winterhalter M, et al.. Sugar  
767 Transport through Maltoporin of *Escherichia coli* : Role of the Greasy Slide. *J Bacteriol.* 2002;  
768 doi: 10.1128/JB.184.11.2994-2999.2002.

769 125. VanBerkum MF, Means AR. Three amino acid substitutions in domain I of calmodulin  
770 prevent the activation of chicken smooth muscle myosin light chain kinase. *J Biol Chem.* Amer-  
771 ican Society for Biochemistry and Molecular Biology; 266:21488–951991;

772 126. Wei Q, Wang L, Wang Q, Kruger WD, Dunbrack RL. Testing computational prediction  
773 of missense mutation phenotypes: Functional characterization of 204 mutations of human  
774 cystathionine beta synthase. *Proteins: Structure, Function, and Bioinformatics.* 2010; doi:  
775 10.1002/prot.22722.

776 127. Williams AD, Shivaprasad S, Wetzel R. Alanine Scanning Mutagenesis of A $\beta$ (1-40) Am-  
777 yloid Fibril Stability. *Journal of Molecular Biology.* 2006; doi: 10.1016/j.jmb.2006.01.041.



778 128. Zhang J, Rao E, Dioszegi M, Kondru R, DeRosier A, Chan E, et al.. The Second Extra-  
779 cellular Loop of CCR5 Contains the Dominant Epitopes for Highly Potent Anti-Human Immu-  
780 nodeficiency Virus Monoclonal Antibodies. *AAC*. 2007; doi: 10.1128/AAC.01302-06.

781 129. Nelsen RB. An introduction to copulas. 2nd ed. New York: Springer;

782 130. Bedó J, Ong CS. Multivariate Spearman's rho for aggregating ranks using copulas. *Jour-*  
783 *nal of Machine Learning Research*. arXiv; 2016; doi: 10.48550/ARXIV.1410.4391.

784 131. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al.. Scikit-  
785 learn: Machine Learning in Python. *Journal of machine Learning research*. :2825–30 2011;

786 132. Hunter JD. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineer-*  
787 *ing*. 2007; doi: 10.1109/MCSE.2007.55.

788 133. Choi Y, Sims GE, Murphy S, Miller JR, Chan AP. Predicting the Functional Effect of  
789 Amino Acid Substitutions and Indels. de Brevern AG, editor. *PLoS ONE*. 2012; doi:  
790 10.1371/journal.pone.0046688.

791 134. Vaser et al.. SIFT missense predictions for genomes. *Nature Protocols*. 2016;

792 135. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al.. A  
793 method and server for predicting damaging missense mutations. *Nature Methods*. 2010; doi:  
794 10.1038/nmeth0410-248.

795 136. Laine E, Karami Y, Carbone A. GEMME: A Simple and Fast Global Epistatic Model  
796 Predicting Mutational Effects. *Molecular Biology and Evolution*. 2019; doi: 10.1093/mol-  
797 bev/msz179.

798 137. Meier J, Rao R, Verkuil R, Liu J, Sercu T, Rives A. Language models enable zero-shot  
799 prediction of the effects of mutations on protein function. *bioRxiv*; doi:  
800 10.1101/2021.07.09.450648.

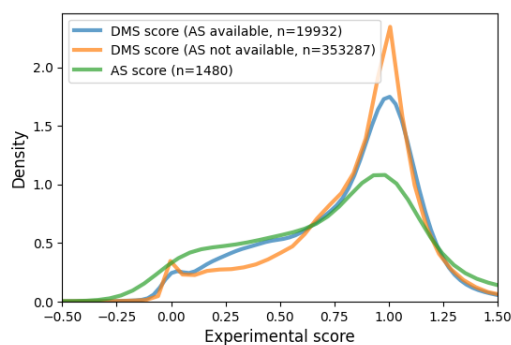
801 138. Frazer J, Notin P, Dias M, Gomez A, Min JK, Brock K, et al.. Disease variant prediction  
802 with deep generative models of evolutionary data. *Nature*. 2021; doi: 10.1038/s41586-021-  
803 04043-8.

804 139. Livesey BJ, Marsh JA. Updated benchmarking of variant effect predictors using deep  
805 mutational scanning. *Molecular Systems Biology*. John Wiley & Sons, Ltd; 2023; doi:  
806 10.15252/msb.202211474.

807 140. González J, Dai Z, Hennig P, Lawrence ND. Batch Bayesian Optimization via Local Pe-  
808 nalization. *arXiv*; doi: 10.48550/arXiv.1505.08052.

809

## 810 **Supplementary material**

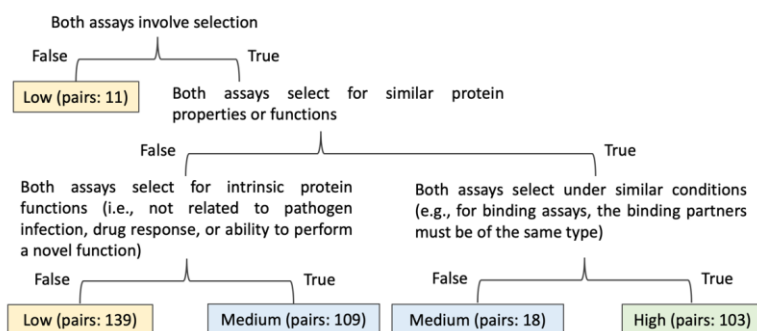


811

812 **Fig S1. DMS and AS score distribution.** The figure shows the kernel estimated density of normalized AS  
813 scores and DMS scores for variants with or without available AS data.

814

For each pair of DMS and AS experiments:



815

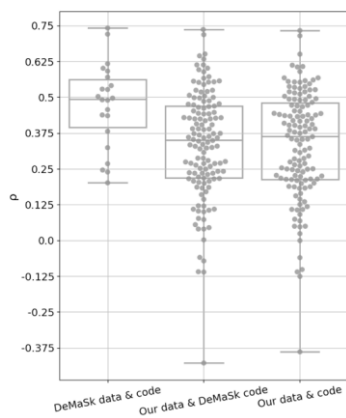
816 **Fig S2. Decision tree for classifying DMS and AS assay compatibility.** The similarity of DMS and AS assays

817 are compared (Methods) and the DMS/AS assay pairs are classified using three levels of compatibility (low,

818 medium, high). The leaf-node text and color show the classified assay compatibility. The number indicates the

819 count of assay pairs for each compatibility level.

820



821

822 **Fig S3. Comparison between published and re-implemented predictors.** The plot shows leave-one-protein-

823 out cross-validation performance on predictors built from the published DeMaSk code or our code. The predictors

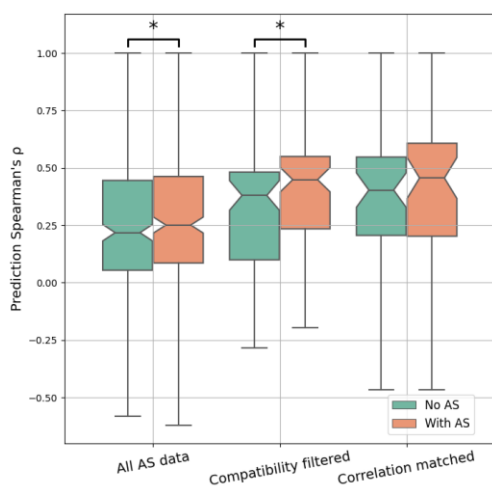
824 were trained and evaluated on DMS data either provided by the DeMaSk study or curated by our own. The

825 “DeMaSk data & code” result is similar to the published result. For the “Our data & DeMaSk code” result, we

826 used our own data and published code which shows a median performance around 0.35. This is probably because

827 many more DMS results are included in our data. The similarity of results achieved using “Our data & code”  
828 demonstrates the correctness of our re-implementation. (Whiskers show the full value range)

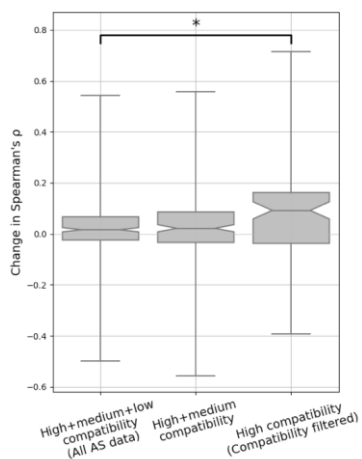
829



830

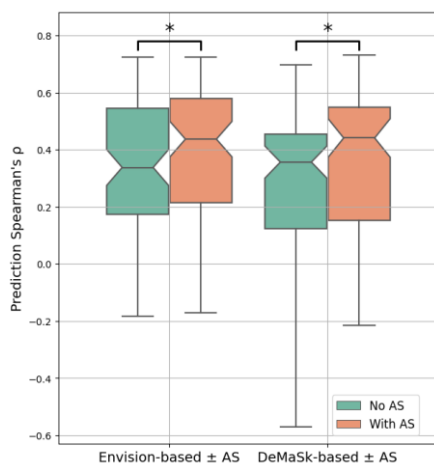
831 **Fig S4. Performance comparison between predictors with or without AS data.** The Spearman's  $\rho$  between  
832 DMS scores and predicted scores for each DMS and AS data pair are shown as box plots. Different approaches  
833 to filtering the data are shown on the x-axis: “All AS data” used all available data; “Compatibility filtered” used  
834 only data of high assay compatibility; “Correlation matched” used only data with the highest regularised correla-  
835 tion for each DMS dataset. The figure does not include data without available AS scores. This means that the  
836 different results are not directly comparable since they are computed for different subsets of DMS/AS data pairs  
837 (for example, “All AS data” contains all DMS/AS data pairs, but “Compatibility filtered” contains only data pairs  
838 of high assay compatibility). Control results are shown as green boxes for predictions on the same residues without  
839 AS data as a feature. The underlying  $\rho$  for each data pair in the control results is the same, but the boxes are shifted  
840 due to data filtering. Results for data pairs with only one residue are not shown. P-values were calculated using  
841 paired t-test and jointly corrected using Holm-Šidák (Methods), \*:  $p < 0.05$ . Notches show the 95% confidence  
842 interval around the median, and whiskers show the full value range.

843



844

845 **Fig S5. The change in prediction performance for using data of different assay compatibility levels.** The  
 846 change of prediction Spearman's  $\rho$  for each DMS and AS data pair is shown as box plots. A higher value represents  
 847 higher prediction accuracy achieved for using AS data. Different data filtering methods are shown on the x-axis.  
 848 Results for data pairs with only one residue are not shown. P-values were calculated using Welch's test and jointly  
 849 corrected using Holm-Šidák (Methods), \*:  $p < 0.05$ . Notches show the 95% confidence interval around the median,  
 850 and whiskers show the full value range.

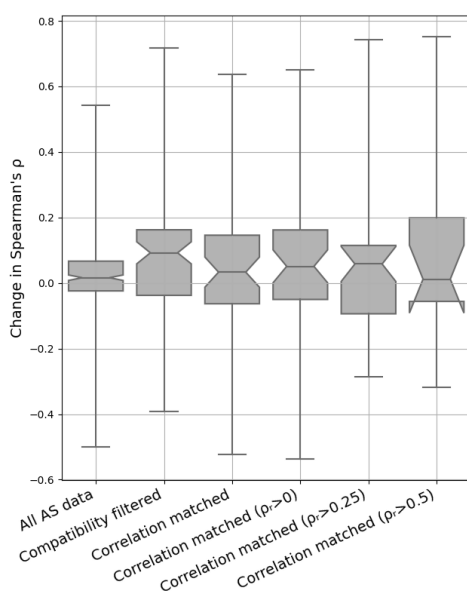


851

852 **Fig S6. Prediction performance is improved while incorporating high compatibility AS data into the En-**  
 853 **vision model.** The Spearman's  $\rho$  between experiment DMS scores and predicted scores for each DMS/AS assay

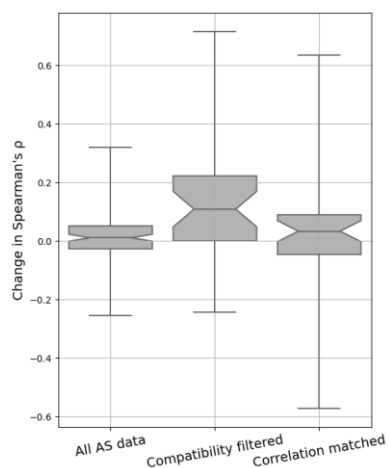
854 pair with high compatibility are shown as box plots. The x-axis shows the predictor used, either Envision or  
 855 DeMaSk. Control results are shown as green boxes for predictions on the same residues without AS data as a  
 856 feature. Results for data pairs with only one residue are not shown. P-values were calculated using paired t-test  
 857 and jointly corrected using Holm-Šidák (Methods), \*:  $p < 0.05$ . Notches show the 95% confidence interval around  
 858 the median, and whiskers show the full value range.

859



860 **Fig S7. Performance improvement on thresholded correlation matching.** The change of prediction  $\rho$  for  
 861 each DMS and AS data pair is shown as box plots. Different approaches to filtering/matching the data are shown  
 862 on the x-axis: “All AS data”, “Compatibility filtered” and “Correlation matched” are the same results as previously  
 863 discussed; while doing correlation matching, a further thresholding (0, 0.25 or 0.5) on the regularized DMS/AS  
 864 correlation values ( $\rho_r$ ) was applied. Notches show the 95% confidence interval around the median, and whiskers  
 865 show the full value range.

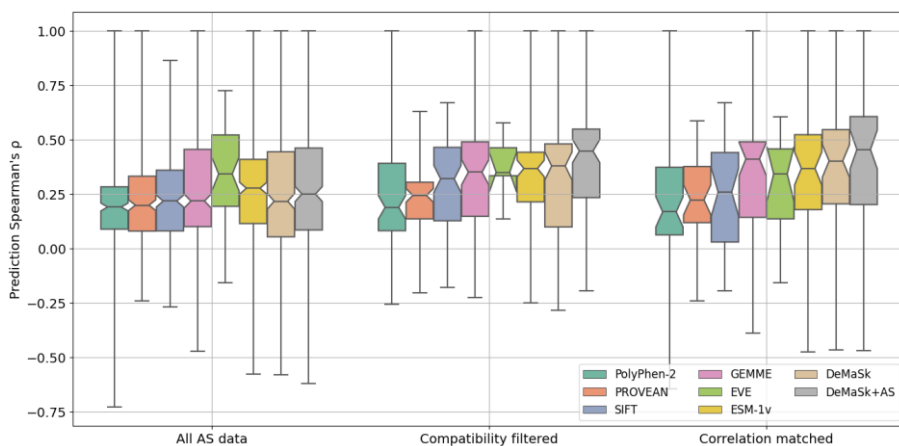
867



868

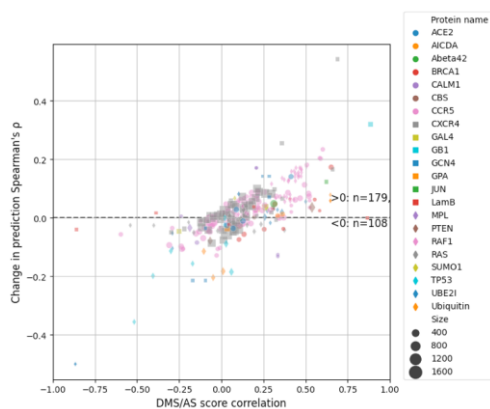
869 **Fig S8. Performance improvement on averaged DMS/AS testing data.** This figure shows model perfor-  
 870 mance when we averaged variant scores for DMS or AS data that are: i) published in the same paper; ii) targeting  
 871 the same protein region; iii) measured by the same type of assays (Supplementary Table 1). The change of pre-  
 872 diction  $\rho$  for each averaged DMS and AS data pair is shown. A higher value represents higher prediction accuracy  
 873 achieved when using AS data. Different approaches to filtering/matching the data are shown on the x-axis: “All  
 874 AS data” used all available data; “Compatibility filtered” used only data of high assay compatibility; “Correlation  
 875 matched” used only data with the highest regularised correlation for each DMS dataset. Results for data pairs  
 876 with only one residue are not shown. Notches show the 95% confidence interval around the median, and whiskers  
 877 show the full value range.

878



879  
 880 **Fig S9. Model performance on various variant effect predictors.** The Spearman's  $\rho$  between DMS scores  
 881 and predicted scores from different variant effect predictors for each DMS and AS pair are shown as box plots.  
 882 Results are evaluated on different sets of variant data shown on the x-axis: "All AS data" used all available data;  
 883 "Compatibility filtered" used only data of high assay compatibility; "Correlation matched" used only AS data  
 884 with the highest regularised correlation for each DMS dataset. The figure does not include residues without avail-  
 885 able AS scores. Results for data pairs with only one residue are not shown. Notches show the 95% confidence  
 886 interval around the median, and whiskers show the full value range.

887

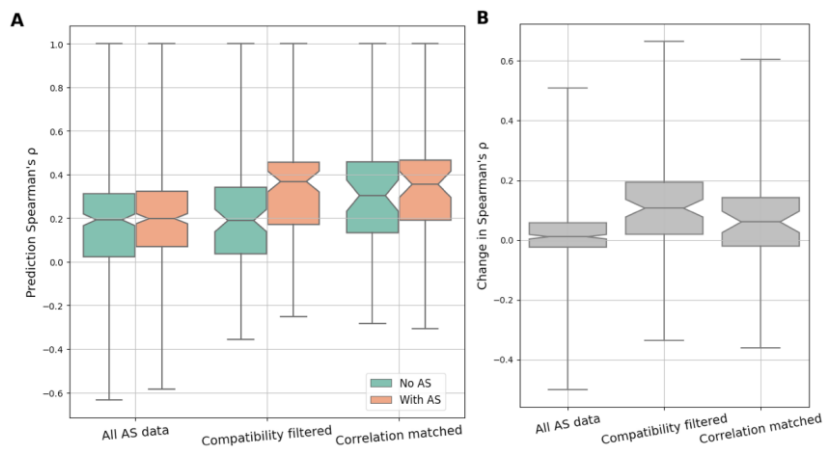


888  
 889 **Fig S10. Prediction performance change for using all AS data.** Each dot represents a DMS/AS data pair. The  
 890 vertical axis shows the change of prediction  $\rho$  by using AS data (larger means higher performance achieved by



891 using AS data). The horizontal axis shows the DMS/AS score correlation for *all* variants on the matched residues  
892 rather than just alanine substitutions. The colours and shapes of the dots correspond to the target protein, and size  
893 indicates the number of variants in each data pair. Results for data pairs with only one residue are not shown.

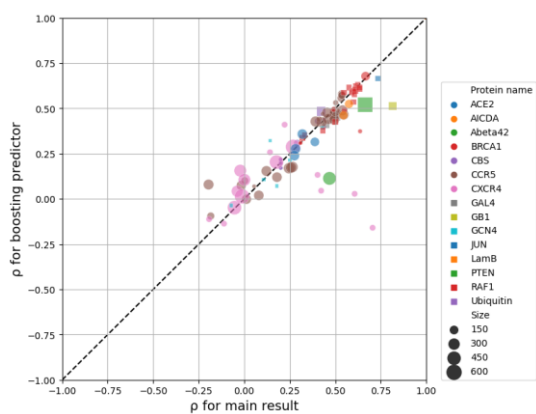
894



895

896 **Fig S11. Model performance for training with AS-data-available-residues.** The predictors were trained only  
897 on variants that have AS data available. Panel A shows the performance visualized by prediction Spearman's  $\rho$   
898 for DMS scores and predicted scores for each DMS and AS data pair. Different approaches to filtering the data  
899 are shown on the x-axis: "All AS data" used all available data; "Compatibility filtered" used only data of high  
900 assay compatibility; "Correlation matched" used only AS data with the highest regularised correlation for each  
901 DMS dataset. Control results are shown as green boxes for predictions on the same residues without AS data as a  
902 feature. Panel B shows change of prediction  $\rho$  for each DMS and AS data pair. A higher value indicates higher  
903 prediction accuracy achieved when using AS data. Different approaches to filtering the data are also shown on  
904 the x-axis as described. Notches show the 95% confidence interval around the median, and whiskers show the full  
905 value range.

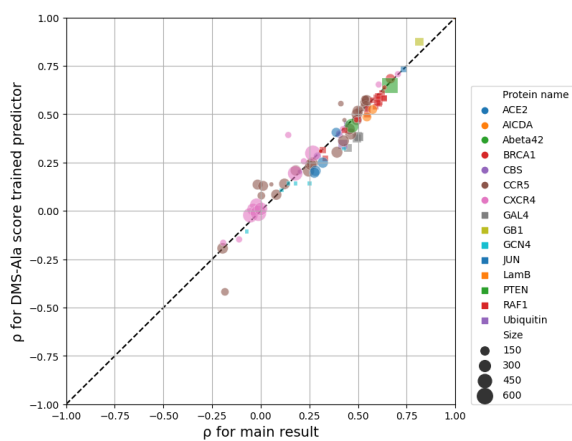
906



907

908 **Fig S12. Boosting setup shows similar performance as the main result.** Each dot represents a filtered  
 909 DMS/AS data pair of high assay compatibility. The vertical and horizontal axes show the prediction Spearman's  
 910  $\rho$  for either modelled with boosting or the one-step (main result) setup. The colours and shapes of the dots corre-  
 911 spond to the target protein, and size indicates the number of variants in each data pair.

912

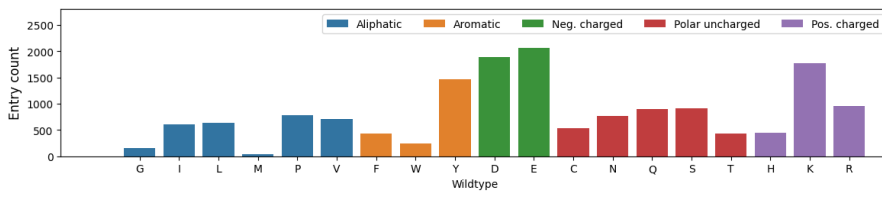


913

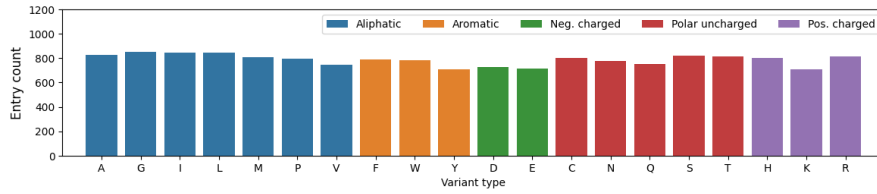
914 **Fig S13. Training with DMS scores of alanine substitutions shows similar performance as the main result.**  
 915 The vertical and horizontal axes show the prediction Spearman's  $\rho$  for predictors either trained with DMS score  
 916 of alanine substitutions (DMS-Ala) or AS data of high assay compatibility (main result), yet all evaluated on high

917 compatibility AS data. The colours and shapes of the dots correspond to the target protein, and size indicates the  
918 number of variants in each data pair.

919



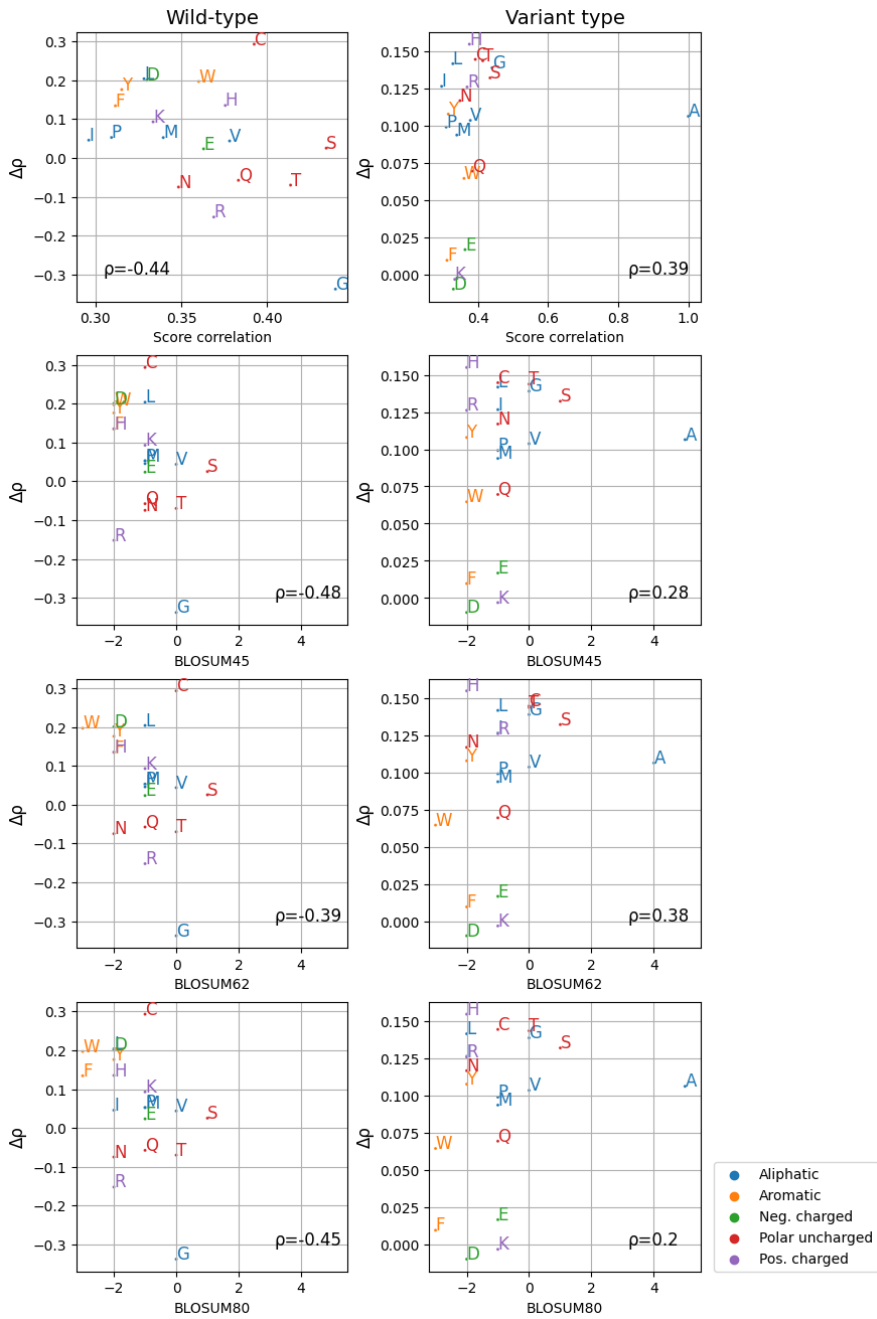
920



921

922 **Fig S14. Count of variant entries for each wild-type or variant amino acid of high assay compatibility data.**  
923 (Neg.: negatively, Pos.: positively)

924



926 **Fig S15. Relationship between amino acid similarity and model performance.** For each amino acid, its sim-  
 927 ilarity to alanine was computed by their DMS score correlation or using BLOSUM scores as shown on the x-axis.  
 928 The performance improvement ( $\Delta\rho$ ) for each wild-type (left) or variant (right) amino acid while using AS data  
 929 were computed as previously mentioned (Fig 7), with their Spearman's correlation against the similarity meas-  
 930 urements shown on the figure. The label for each amino acid is coloured by the amino acid physicochemical  
 931 property. (Neg.: negatively; Pos.: positively)

932

933 **Table S1. DMS/AS correlation on each secondary structural region.** The secondary structure of each variant  
 934 is determined by UniProt annotations. The Spearman's correlation between DMS and all or high compatibility  
 935 AS data on each structural region is computed, with the number of protein residues involved shown in parenthesis.

$\rho$ (n_residues)	HELIX	STRAND	TURN
All AS	0.13 (233)	0.13 (83)	0.17 (22)
AS of high com- patibility	0.28 (115)	0.26 (56)	0.41 (15)

936

937 **Table S2. Amount of data with AS scores available**

Data composition	Protein	DMS dataset	AS dataset <sup>1</sup>	Variant entries <sup>2</sup>
All AS	22	54	146	70446
Compatibility filtered	15	35	60	15739
High+medium assay com- patibility	21	51	105	28380
Correlation matched	22	54	32	7940

938 1. This column shows how many unique AS datasets are included.

939 2. Include duplicated variants caused by multiple experiments targeting the same protein variant.

940

## 941 **Supplementary information**

### 942 **Applying AS data to Envision method**

943 We re-implemented a predictor based on Envision [17] to incorporate AS data. Features used  
944 in Envision were downloaded from its online toolkit. All Envision features are used for mod-  
945 elling except for substitution type (wt\_mut) which has low importance according to the pub-  
946 lished result and our pilot studies yet is computationally expensive in our setup. Protein data  
947 were excluded if their features were not available online. DMS and AS data pairs with high  
948 assay compatibility were used for modelling. Missing feature values were imputed by the mean  
949 values for numerical features or the most frequent values for categorical features. Categorical  
950 features are encoded with the one-hot encoder. We used `sklearn.ensemble.Gradi-  
951 entBoostingRegressor` from scikit-learn package [131] to build the predictor, and hy-  
952 perparameters were tuned by Bayesian Optimization [140] with Group K-Fold (protein-30-fold)  
953 cross-validation. The training and evaluation process were similar to that previously described.  
954 For comparison, we repeated the DeMaSk-based analysis on the same subset of data.

955

### 956 **Boosting with AS data**

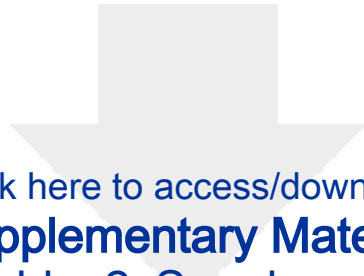
957 To deal with the sparsity of AS data, we tested a variant impact predictor based on boosting. A  
958 first linear regression predictor was trained with all training DMS data using the three DeMaSk  
959 features without AS data, which was the same as the control predictor mentioned previously.  
960 We then calculated the prediction error by subtracting the predicted scores from DMS scores,  
961 and a second linear regression predictor was trained to predict the error. The second predictor  
962 was trained only on DMS/AS data of high assay compatibility and used both protein features  
963 and the encoded AS scores. The final prediction result was the sum of the outputs from these  
964 two predictors.

965

966 **Replacing AS data with DMS scores of alanine substitutions**

967 We investigated another potential approach to overcome the sparsity of AS data by replacing  
968 the AS feature with the DMS scores of alanine substitutions (DMS-Ala). The intention of this  
969 study is to model the scenario of ideal AS data, which perfectly matches the DMS-Ala data  
970 during training. To do this, for all DMS datasets we collected, their AS feature values, regard-  
971 less of availability, were replaced by the DMS-Ala scores on the same residue. Missing scores  
972 were imputed by the mean value of all DMS-Ala scores. A regression model was trained and  
973 evaluated as previously described, using the three DeMaSk features as well as the DMS-Ala  
974 scores. The AS data of high assay compatibility are still used for the testing process.

975



Click here to access/download

**Supplementary Material**

[Supplementary\\_Table\\_2\\_Supplementary\\_Material.csv](#)





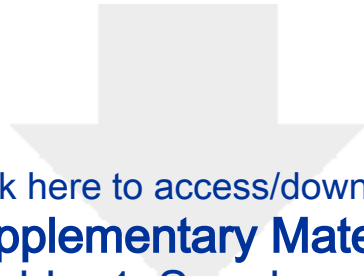


Click here to access/download

**Supplementary Material**

Supplementary\_Table\_3\_Supplementary Material.csv





Click here to access/download

**Supplementary Material**

Supplementary\_Table\_1\_Supplementary Material.xlsx

