# GigaScience

## Integrating deep mutational scanning and lowthroughput mutagenesis data to predict the impact of amino acid variants
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-23-00040 |
| Full Title: | Integrating deep mutational scanning and lowthroughput mutagenesis data to predict the impact of amino acid variants |
| Article Type: | Research |

| Abstract: | 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. 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. 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. |
|---|---|

| | |
|---|---|
| Corresponding Author: | Alan F Rubin, PhD<br>Walter and Eliza Hall Institute of Medical Research<br>Parkville, VIC AUSTRALIA |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | Walter and Eliza Hall Institute of Medical Research |
| Corresponding Author's Secondary Institution: | |
| First Author: | Yunfan Fu |
| First Author Secondary Information: | |
| Order of Authors: | Yunfan Fu |
| | Justin Bedő |
| | Anthony Troy Papenfuss, BSc (Hons) PhD |
| | Alan F. Rubin |
| Order of Authors Secondary Information: | |
| Additional Information: | |

| Question | Response |
|---|---|
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our [Minimum Standards Reporting Checklist](#). Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>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 [Research Resource Identifiers](#) (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our [Minimum Standards Reporting Checklist](#)? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript. | Yes |

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

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

**Authors:**

Yunfan Fu[1,2], Justin Bedő[1,3,*], Anthony T. Papenfuss[1,2,4,*,**], Alan F. Rubin[1,2,*,**]


**Affiliations:**

[1]The Walter and Eliza Hall Institute of Medical Research, Parkville, 3052, Victoria, Australia.

[2]Department of Medical Biology, The University of Melbourne, Melbourne, VIC 3010, Australia.

[3]School of Computing and Information Systems, The University of Melbourne, Melbourne, VIC 3010, Australia.

[4]Peter MacCallum Cancer Centre, Melbourne, VIC 3000, Australia.


\* Contributed equally

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

## Abstract

**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.

**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.

**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.

# 1    Introduction

Deep mutational scanning (DMS) is a functional genomics method that can experimentally measure the impact of many thousands of protein variants by combining high-throughput sequencing with a functional assay [1]. In a typical DMS, a cDNA library of genetic variants of a target gene is generated, containing all possible single amino acid substitutions. This variant library is then expressed in a functional assay system where the variants can be selected based on their properties. The change in variant frequency in the pre- and post-selection populations is determined by high-throughput sequencing which is then used to calculate a multiplexed functional score that captures the variant's impact [2–4]. The versatility of DMS assays makes it possible to measure variant impact on a wide range of protein properties, including protein binding [5,6], protein abundance [7–9], catalytic activity [10,11] and cell growth rate [12–14].

48 Computational studies have used DMS data to build predictive models of variant impact. These

49 predictors use supervised or semi-supervised learning models trained on experimental DMS

50 data and various protein features to make predictions [15–21]. Envision is one such method

51 that used protein structural, physicochemical, and evolutionary features to predict variant effect

52 scores and was trained on DMS data from 8 proteins using gradient boosting [15]. Another

53 method, DeMaSk, predicted DMS scores by combining two evolutionary features (protein

54 positional conservation and variant homologous frequency) with a DMS substitution matrix

55 and was trained on data from 17 proteins using a linear model [17]. Deep learning algorithms

56 have also been applied to build protein fitness predictors [16,18], which are usually based only
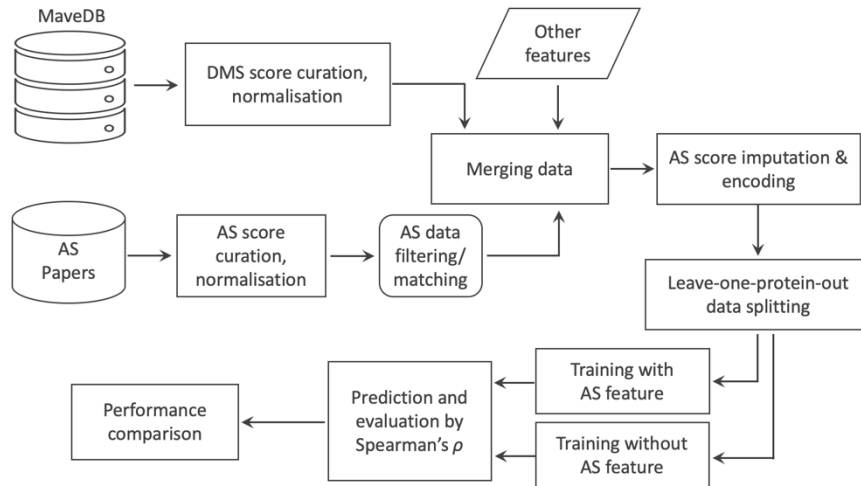
57 on variant sequences.

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 [22,23], protein stability [24,25], and protein activity [26,27]. Alanine scanning (AS) is a

62 widely-used low-throughput mutagenesis method [28,29], and AS data are available for many

63 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 [30]. AS experiments are typically

65 used to identify functional hot spots or critical residues in the target protein [31,32] and have

66 been used as a source of independent validation for DMS studies [27,33–35].

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

3

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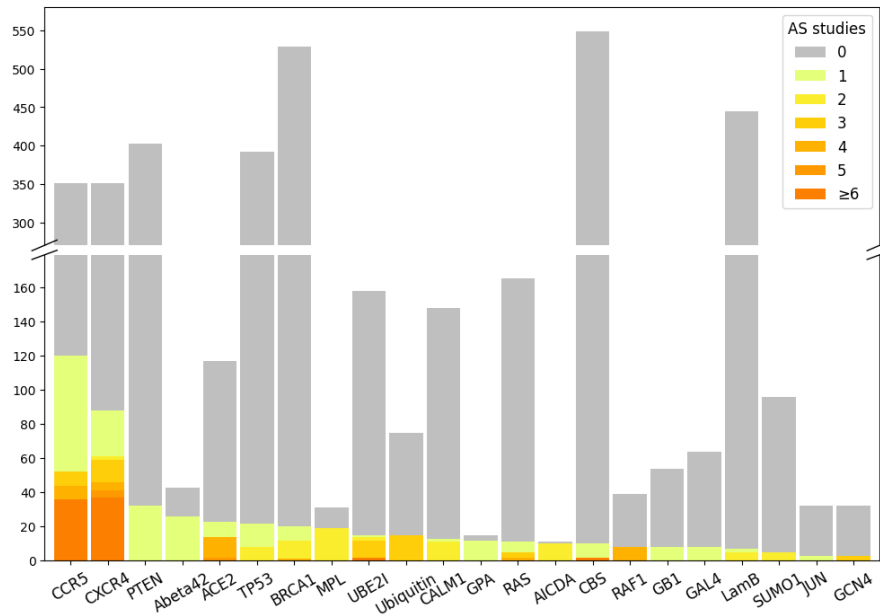
## 2    Results

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

To build the predictive model, 130 DMS datasets were collected from MaveDB [36,37] (Supplementary table 1). We searched the literature and found 146 AS datasets targeting the same proteins as 54 of the DMS datasets. In total, we obtained both DMS and AS data for 22 different proteins: 17 human proteins, three yeast proteins, and two bacterial proteins. Most DMS experiments were highly complete, with a mean coverage of 95.0% of all possible single amino acid substitutions assayed in the target region, comprising 373,219 total protein variant measurements. AS data were only available on a small number of protein residues (Fig 2), and we were able to curate 1,480 alanine substitution scores from the 146 studies. Variant scores from collected DMS and AS studies were linearly normalized to a common scale (see Methods) to make them comparable across datasets (Fig S1).

92



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

96

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

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

103

104

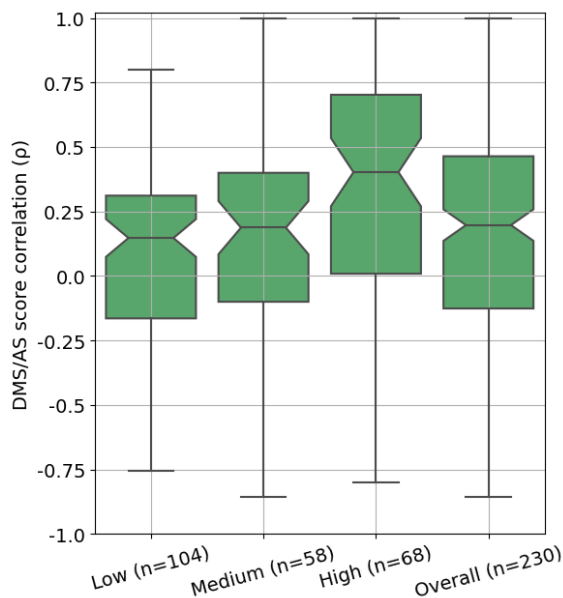**Fig 3. Correlation between DMS and AS data shows substantial variation.** We calculated Spearman's $\rho$

between alanine substitution scores in each pair of AS and DMS data. The results for pairs with less than three

alanine substitutions are removed. The red dashed line shows the median $\rho$.

108

We then considered if differences between AS and DMS assay designs might contribute to this

low agreement between scores. To explore this, we developed a decision tree (Fig S2) to

classify whether DMS/AS pairs had low, medium, or high assay compatibility, which we

defined as a similarity measurement of the functional assays performed. For example, the DMS

assay measuring the binding affinity of a cell surface protein, CXCR4, to its natural ligand [38]

has high compatibility with the AS experiment also measuring this ligand binding but has low

compatibility with the study on CXCR4's ability to facilitate virus infection [39]. A full assay

compatibility table can be found in Supplementary Table 1 with the compatibility

classifications and justification for each pair. We then compared DMS and AS score correlation

for each compatibility class and found that score correlations were closely related to assay

compatibility. Data from low compatibility assays had a median correlation of 0.15, rising to

0.19 for medium compatibility assays and 0.40 for high compatibility assays (Fig 4). This link

6

121 between assay compatibility and score correlation indicates that our decision tree approach was

122 able to capture the similarity between assay systems.

123



124

**Fig 4. DMS and AS data pairs with high assay compatibility show a higher score correlation.** Each box

represents Spearman's $\rho$ between DMS and AS data pairs of classified assay compatibility or the overall result.

The correlation coefficients are calculated between alanine substitution scores in each pair of AS and DMS data.

Results for data pairs with less than three alanine substitutions are removed.

129

## 2.3 Compatible AS data improve DMS score prediction accuracy

131 To test if incorporating AS data into DMS score models would improve prediction accuracy,

132 we decided to build a new model based on DeMaSk [17]. We chose DeMaSk because it showed

133 better performance compared to similar methods and was straightforward to modify. The

134 published DeMaSk model predicts DMS scores using protein positional conservation, variant

135 homologous frequency, and substitution score matrix, and we incorporated AS data as an

136 additional feature. Our new predictor was modelled with all 130 DMS we collected and we

137 applied a leave-one-protein-out cross-validation approach to training and testing [15].

138 Prediction performance was evaluated using the Spearman's correlation ($\rho$) between the

139    experimentally-derived DMS scores and the predicted scores for each pair of DMS and AS

140    studies. The performance of our DMS/AS model was compared with a model trained only on

141    DMS data, equivalent to retrained DeMaSk (Fig S3), by calculating the change of prediction $\rho$

142    (see Methods).

143

144    We trained our model with either all or a subset of AS data we collected (Fig 5, Table S1). We

145    first integrated all 146 AS data collected for training and evaluation but observed only a modest

146    improvement of prediction $\rho$ (Fig 5 left box, and Fig S4). We then retrained and evaluated our

147    model on filtered AS data with only high compatibility assays, and observed a median increase

148    in prediction Spearman's $\rho$ of 0.1 compared to the results with no AS data (Fig 5 middle box,

149    and Fig S4). However, training with both high and medium compatibility pairs reduced the

150    performance improvement (Fig S5). These results indicate that medium and low compatibility

151    pairs might provide inconsistent training data, degrading model performance. We also

152    evaluated the impact of including high compatibility AS data in an alternative model based on

153    Envison [15], and found similar results (Fig S6). To differentiate between high assay

154    compatibility and high DMS/AS score correlation, we trained the model using the most highly

155    correlated AS result for each DMS dataset (see Methods). Although the upper quartile was

156    high, the median performance change of this predictor was lower than the high assay

157    compatibility model, suggesting that matching with the highest score correlation alone is

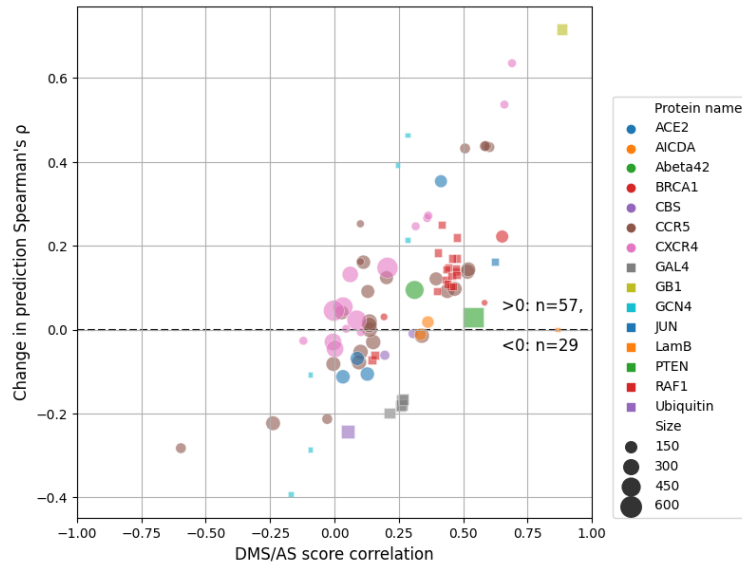158    insufficient (Fig 5 right box).

159

160

**Fig 5.  Performance of variant impact prediction is improved using AS data with high assay compatibility.**
The change of prediction $\rho$ for each DMS and AS data pair is shown as box plots. A higher value represents higher
prediction accuracy achieved for using AS data. Different approaches to filtering/matching the data are shown on
the x-axis: "All AS data" used all available data; "Compatibility filtered" used only data of high assay
compatibility; "Correlation matched" used only data with the highest regularised correlation for each DMS dataset.

166

To further explore the higher performance of compatibility-filtered predictor, we examined the
relationship between prediction $\rho$ change and score correlation for each high compatibility
DMS/AS pair (Fig 6). For most pairs, prediction performance was improved by using AS data,
and the scale of improvement was also related to the score correlation. This relationship could
also be observed for multiple DMS/AS pairs from an individual protein, such as CXCR4 and
CCR5. We saw the same trend in the predictor trained with all DMS/AS pairs but noted that
the performance even of highly correlated pairs was worse, likely due to the influence of low
compatibility training data on the model (Fig S7).

175

**Fig 6. Prediction performance change is related to DMS and AS score correlation.** Each dot represents a filtered DMS/AS data pair of high assay compatibility. The vertical axis shows the change of prediction $\rho$ by using AS data (larger means higher performance achieved by using AS data). The horizontal axis shows the DMS/AS score correlation for *all* variants on the matched residues rather than just alanine substitutions. The colours and shapes of the dots correspond to the target protein, and size indicates the number of variants in each data pair.

We also explored the consequences of the sparsity of AS data on our model in two ways: by using a boosting approach that focuses only on residues with AS data (Fig S8) and by using complete alanine substitution information from DMS as the AS feature (Fig S9). Both of these approaches performed very similarly to the primary model constructed using high-compatibility DMS/AS data and simple mean score imputation.

To test the influence of amino acids on our predictor, we grouped the prediction results by either wild-type or variant amino acid and calculated the prediction improvement when AS data were included (Fig 7). We found that 14 of 19 wild-type amino acids performed better with the addition of AS data, with cysteine showing the largest improvement and performing worst in the model lacking AS data. 18 of 20 variant amino acids benefited from the inclusion

194 of AS data, with marginal performance decrease on lysine and aspartic acid ($|\Delta\rho|<0.01$) (Fig

195 7).

196



197

**Fig 7. Model perfomance is generally improved for each wild-type and variant amino acid.** Prediction

199 Spearman's $\rho$ when using (y-axis) or not using (x-axis) AS data on each wild-type (left) or variant (right) amino

200 acid is shown in the scatter plots. The results are coloured according to the property of each amino acid type.

201 Alanine (A) result is not applicable in the first figure since alanine scanning data are always missing when the

202 wildtype is alanine itself. Absolute count for each amino acid can be found in Fig S10. (Neg.: negatively, Pos.:

203 positively)

204

## 3    Discussion

206 In this study, we integrated alanine scanning (AS) data into deep mutational scanning (DMS)

207 score prediction, leading to modest improvements in the accuracy of variant score prediction.

208 We also explored the impact of the diversity of protein properties measured by DMS and AS.

209 Filtering DMS and AS data based on our manual classification of assay type compatibility led

210 to improved prediction performance.

211

212  A potential shortcoming of our current approach is that AS data were available for only a small

213  proportion of the DMS data. Although most recent DMS studies can analyze variants of the

214  whole protein, most AS experiments only cover a handful of residues in the target protein,

215  leaving missing AS scores for the vast majority of residues. We explored this here and found

216  that alternative methods for addressing the sparsity of AS data did not improve or degrade

217  performance, but we anticipate further improved prediction accuracy if the low completeness

218  and unevenness of AS data are appropriately handled before modelling, such as by advanced

219  imputation methods [48,49].

220

221  In this study, we identified the importance of DMS/AS assay compatibility as a crucial factor

222  for improving prediction accuracy. An issue with using this concept is that it further shrinks

223  already sparse data. It also fails to take advantage of the fact that even for low compatible

224  assays some fundamental information like protein stability can still be mutually captured.

225  Instead of hard filtering, proper implementation of this underlying information may facilitate

226  variant impact prediction in the future. Nonetheless, filtering on assay compatibility still leads

227  to performance improvement. We also briefly explored whether the consistency of DMS and

228  AS scores can be considered more directly by matching the best correlated AS data for each

229  DMS dataset. Consistency is partially driven by assay compatibility but also reflects other

230  features of the data, such as bias and noise. While we picked the most correlated pair for each

231  DMS, we did not threshold the correlation, potentially including data pairs that were poor

232  matches.

233

234  The concepts of compatibility and data quality are also relevant to training any DMS-based

235  predictors. DMS assays have been developed to measure variant impacts to distinct protein

236  properties, and a variant can behave similarly to wildtype when measured by one assay yet

237  show altered protein properties in other assay results, which are frequently found in regions

238  with specific biochemical functions [50–55]. With more experimental assays to be applied, the

239  diverse measurements may impede the progress of future DMS-based predictors unless this

240  assay effect is properly addressed, for example, by building assay specific predictors.

241  Measurement error is another source of DMS data heterogeneity that potentially affects the

242  model performance. In our current study, DMS scores of protein variants are weighted equally

243  while training. Adjustable weighting can be applied in future studies to adapt the distinct

244  experimental error between individual variants and datasets, reducing the influence of low-

245  confident data.

246

247  In summary, we conclude that the careful inclusion of low-throughput mutagenesis data

248  improves the prediction of DMS scores, and the approaches described here can potentially be

249  applied to other prediction methods.

250

## 251  4    Availability of supporting source code and requirements

252  **Project name:** DMS_with_Alanine_scan

253  **Project home page:** https://github.com/PapenfussLab/DMS_with_Alanine_scan

254  **Operating system:** Platform independent

255  **Programming language:** Python

256  **Other requirements:** Python 3.10.6

257  **Licence:** MIT Licence

258

## 259  5    List of abbreviations

260  DMS: deep mutational scanning

261  AS: alanine scanning

262

## 6    Supporting information

263

**Supplementary Table 1:** All candidate DMS and alanine scanning data with detailed dataset

264

265 information.

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

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

268

## 7    Author contributions

269

270 YF developed the software and wrote the initial draft of the manuscript. AFR conceived the

271 study. JB, AFR, and ATP oversaw the project. All authors reviewed, contributed to, and

272 approved the manuscript.

273

## 8    Funding

274

283

## 9 Methods

### 9.1 DMS data collection

DMS data were downloaded from MaveDB [36,37] which were then filtered and curated. DMS experiments targeting antibody and virus proteins were removed because of their potentially unique functionality. We retrieved the UniProt accession ID of target proteins by searching the protein names or sequences in UniProt [56], and proteins lacking available UniProt ID were also excluded. Datasets that are computationally processed or their wildtype-like and nonsense-like scores (see Normalization) cannot be identified were also filtered out (Supplementary Table 1). All missense variants with only a single amino acid substitution were curated from the DMS studies for our analysis. A total of 130 DMS experiments from 53 studies [5,6,9–14,27,33–35,38,57–94] were collected for our analysis.

### 9.2 Collection of AS data and other features

The following process was used to search for candidate AS studies. Papers were identified by searching on PubMed and Google Scholar for the "alanine scan" or "alanine scanning" together with the name of candidate proteins. While searching in Google Scholar, we included the protein's UniProt ID rather than molecule name as the search term to reduce false positives. Appropriate AS data were collected from the search results. Western blot results were transformed to values by ImageJ if it was the only experimental data available in the study. A total 146 AS experiments were collected from 45 distinct studies [22–24,26,27,39–42,44,45,84,95–127].

Protein features of Shannon entropy and the logarithm of variant amino acid frequency were downloaded from the DeMaSk online toolkit [17]. The substitution score matrix feature was calculated from the mean of training DMS scores for each of the 380 possible amino acid substitutions before each iteration of cross-validation.

15

309

## 9.3 Normalization

DMS and AS datasets were normalized to a common scale using the following approach adapted from previous studies [15,43]. Let $D$ denotes a protein study measuring scores $s_i^D$ for a single variant $i$, $s_{wt}^D$ denotes the scores for wildtype and $s_{non}^D$ represents the score for nonsense-like variants. The normalized scores $s_i'^D$ are given by:

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

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

## 9.4 AS data filtering and matching

AS data subsets were filtered/matched according to either assay compatibility or score correlation. For assay compatibility filtering, DMS and AS assay pairs were first classified into three levels of compatibility (Fig S2). For each DMS dataset, we first tried to use only AS data with high assay compatibility for further modelling, removing AS data of medium and low assay compatibility. We then also tried to model with AS data of both high and medium assay compatibility.

For score correlation matching, Spearman's correlation ($\rho$) is calculated between alanine substitution scores in each pair of AS and DMS data. To avoid influence from the size of AS datasets, we regularised the $\rho$ value by empirical copula [128]:

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

333  where $\rho_r$ is the regularised correlation coefficient, and $n$ is the number of alanine substitutions

334  used for correlation calculation. For each DMS dataset, AS result with the highest $\rho_r$ was

335  picked for modelling.

336

337  **9.5   AS data pre-processing**

338  AS data were pre-processed prior to modelling. For variants without available

339  (filtered/matched) AS data, their AS scores were imputed with the mean value of all available

340  AS scores. Then the AS data were encoded by the wild-type and variant amino acid type with

341  one-hot-encoding. For each variant, the AS feature is expanded with two one-hot vectors. Each

342  of the vectors has 19 zeros and one non-zero value which was the AS score, with the location

343  of the non-zero value indicating the wild-type or variant amino acid type.

344

345  **9.6   Training and evaluation of DMS score predictor**

346  To build the predictors, we performed linear regression using the function

347  `sklearn.linear_model.LinearRegression` from scikit-learn [129]. Training and

348  validation data were separated with leave-one-protein-out cross-validation. In this process, data

349  from one protein were withheld for subsequent validation, and the rest were used for training.

350  This process was iterated over all proteins in the data. Variants were inversely weighted during

351  the training process by the number of measurements available, thus compensating for some

352  regions having greater coverage with DMS and AS assays. Predictors were trained on protein

353  features, DMS data and (optionally) AS data using four different filtering or matching

354  strategies: i) all DMS/AS data, ii) compatibility-filtered DMS/AS data, iii) correlation-matched

355  DMS/AS data, and iv) a control, constructed using DMS data only.

356  In the evaluation process, let $V$ be protein variants assayed by both DMS study $D$ and AS study

357  $A$. Variant scores are predicted by the previously mentioned predictors either using AS data

17

358    ($\hat{s}_V^A$) or not ($\hat{s}_V$). Spearman's correlation ($\rho$) was calculated between the DMS scores $s_V^D$ and

359    each set of predicted scores. The difference of $\rho$ was used to evaluate the performance change

360    ($\Delta\rho_V$).

361 $$\rho_V^A = \text{Spearman's correlation}(\hat{s}_V^A, s_V^D)$$

362 $$\rho_V = \text{Spearman's correlation}(\hat{s}_V, s_V^D)$$

363 $$\Delta\rho_V = \rho_V^A - \rho_V$$

364    To evaluate, we iterated over variants from each pair of DMS/AS studies. Results were dropped

365    for variants *V* with only one protein residue available during analysis and visualization.

366

367    **10     References**

368    1. Fowler DM, Fields S. Deep mutational scanning: a new style of protein science. *Nature*

369    *Methods*. 2014; doi: 10.1038/nmeth.3027.

370    2. Findlay GM. Linking genome variants to disease: scalable approaches to test the functional

371    impact of human mutations. *Human Molecular Genetics*. 2021; doi: 10.1093/hmg/ddab219.

372    3. Geck RC, Boyle G, Amorosi CJ, Fowler DM, Dunham MJ. Measuring Pharmacogene

373    Variant Function at Scale Using Multiplexed Assays. *Annual Review of Pharmacology and*

374    *Toxicology*. 2022; doi: 10.1146/annurev-pharmtox-032221-085807.

375    4. Weile J, Roth FP. Multiplexed assays of variant effects contribute to a growing genotype–

376    phenotype atlas. *Hum Genet*. 2018; doi: 10.1007/s00439-018-1916-x.

377    5. Diss G, Lehner B. The genetic landscape of a physical interaction. *eLife*. 2018; doi:

378    10.7554/eLife.32472.

379  6. Fowler DM, Araya CL, Fleishman SJ, Kellogg EH, Stephany JJ, Baker D, et al.. High-

380  resolution mapping of protein sequence-function relationships. *Nature Methods*. 2010; doi:

381  10.1038/nmeth.1492.

382  7. Amorosi CJ, Chiasson MA, McDonald MG, Wong LH, Sitko KA, Boyle G, et al.. Massively

383  parallel characterization of CYP2C9 variant enzyme activity and abundance. *The American*

384  *Journal of Human Genetics*. 2021; doi: 10.1016/j.ajhg.2021.07.001.

385  8. Faure AJ, Domingo J, Schmiedel JM, Hidalgo-Carcedo C, Diss G, Lehner B. Mapping the

386  energetic and allosteric landscapes of protein binding domains. *Nature*. 2022; doi:

387  10.1038/s41586-022-04586-4.

388  9. Matreyek KA, Starita LM, Stephany JJ, Martin B, Chiasson MA, Gray VE, et al.. Multiplex

389  assessment of protein variant abundance by massively parallel sequencing. *Nature Genetics*.

390  2018; doi: 10.1038/s41588-018-0122-z.

391  10. Mighell TL, Evans-Dutson S, O'Roak BJ. A Saturation Mutagenesis Approach to

392  Understanding PTEN Lipid Phosphatase Activity and Genotype-Phenotype Relationships. *The*

393  *American Journal of Human Genetics*. 2018; doi: 10.1016/j.ajhg.2018.03.018.

394  11. Stiffler MA, Hekstra DR, Ranganathan R. Evolvability as a Function of Purifying Selection

395  in TEM-1 β-Lactamase. *Cell*. 2015; doi: 10.1016/j.cell.2015.01.035.

396  12. Ahler E, Register AC, Chakraborty S, Fang L, Dieter EM, Sitko KA, et al.. A Combined

397  Approach Reveals a Regulatory Mechanism Coupling Src's Kinase Activity, Localization, and

398  Phosphotransferase-Independent         Functions.       *Molecular       Cell*.       2019;       doi:

399  10.1016/j.molcel.2019.02.003.

400    13. Giacomelli AO, Yang X, Lintner RE, McFarland JM, Duby M, Kim J, et al.. Mutational

401    processes shape the landscape of TP53 mutations in human cancer. *Nature Genetics*. Nature

402    Publishing Group; 2018; doi: 10.1038/s41588-018-0204-y.

403    14. Roscoe BP, Thayer KM, Zeldovich KB, Fushman D, Bolon DNA. Analyses of the Effects

404    of All Ubiquitin Point Mutants on Yeast Growth Rate. *Journal of Molecular Biology*. 2013;

405    doi: 10.1016/j.jmb.2013.01.032.

406    15. Gray VE, Hause RJ, Luebeck J, Shendure J, Fowler DM. Quantitative Missense Variant

407    Effect Prediction Using Large-Scale Mutagenesis Data. *Cell Systems*. 2018; doi:

408    10.1016/j.cels.2017.11.003.

409    16. Alley EC, Khimulya G, Biswas S, AlQuraishi M, Church GM. Unified rational protein

410    engineering with sequence-based deep representation learning. *Nat Methods*. 2019; doi:

411    10.1038/s41592-019-0598-1.

412    17. Munro D, Singh M. DeMaSk: a deep mutational scanning substitution matrix and its use

413    for    variant    impact    prediction.    Xu    J,    editor.    *Bioinformatics*.    2020;    doi:

414    10.1093/bioinformatics/btaa1030.

415    18. Biswas S, Khimulya G, Alley EC, Esvelt KM, Church GM. Low- N protein engineering

416    with data-efficient deep learning. *Nature Methods*. Nature Publishing Group; 2021; doi:

417    10.1038/s41592-021-01100-y.

418    19. Høie MH, Cagiada M, Beck Frederiksen AH, Stein A, Lindorff-Larsen K. Predicting and

419    interpreting large-scale mutagenesis data using analyses of protein stability and conservation.

420    *Cell Reports*. 2022; doi: 10.1016/j.celrep.2021.110207.

421  20. Wu Y, Li R, Sun S, Weile J, Roth FP. Improved pathogenicity prediction for rare human

422  missense variants. *The American Journal of Human Genetics*. 2021; doi:

423  10.1016/j.ajhg.2021.08.012.

424  21. Hsu C, Nisonoff H, Fannjiang C, Listgarten J. Learning protein fitness models from

425  evolutionary and assay-labeled data. *Nat Biotechnol*. 2022; doi: 10.1038/s41587-021-01146-5.

426  22. Block C, Janknecht R, Herrmann C, Nassar N, Wittinghofer A. Quantitative structure-

427  activity analysis correlating Ras/Raf interaction in vitro to Raf activation in vivo. *Nature*

428  *Structural Biology*. Nature Publishing Group; 1996; doi: 10.1038/nsb0396-244.

429  23. Sloan DJ, Hellinga HW. Dissection of the protein G B1 domain binding site for human IgG

430  Fc fragment. *Protein Science*. 1999; doi: 10.1110/ps.8.8.1643.

431  24. Fleming KG, Engelman DM. Specificity in transmembrane helix–helix interactions can

432  define a hierarchy of stability for sequence variants. *PNAS*. National Academy of Sciences;

433  2001; doi: 10.1073/pnas.251367498.

434  25. Shibata Y, White JF, Serrano-Vega MJ, Magnani F, Aloia AL, Grisshammer R, et al..

435  Thermostabilization of the Neurotensin Receptor NTS1. *Journal of Molecular Biology*. 2009;

436  doi: 10.1016/j.jmb.2009.04.068.

437  26. Brzovic PS, Heikaus CC, Kisselev L, Vernon R, Herbig E, Pacheco D, et al.. The Acidic

438  Transcription Activator Gcn4 Binds the Mediator Subunit Gal11/Med15 Using a Simple

439  Protein Interface Forming a Fuzzy Complex. *Molecular Cell*. 2011; doi:

440  10.1016/j.molcel.2011.11.008.

441  27. Gajula KS, Huwe PJ, Mo CY, Crawford DJ, Stivers JT, Radhakrishnan R, et al.. High-

442  throughput mutagenesis reveals functional determinants for DNA targeting by activation-

443  induced deaminase. *Nucleic Acids Research*. 2014; doi: 10.1093/nar/gku689.

444  28. Kortemme T, Kim DE, Baker D. Computational Alanine Scanning of Protein-Protein

445  Interfaces. *Science's STKE*. American Association for the Advancement of Science; 2004; doi:

446  10.1126/stke.2192004pl2.

447  29. Morrison KL, Weiss GA. Combinatorial alanine-scanning. *Current Opinion in Chemical*

448  *Biology*. 2001; doi: 10.1016/S1367-5931(00)00206-4.

449  30. Cunningham BC, Wells JA. High-resolution epitope mapping of hGH-receptor interactions

450  by alanine-scanning mutagenesis. *Science*. American Association for the Advancement of

451  Science; 1989; doi: 10.1126/science.2471267.

452  31. DeLano WL. Unraveling hot spots in binding interfaces: progress and challenges. *Current*

453  *Opinion in Structural Biology*. 2002; doi: 10.1016/S0959-440X(02)00283-X.

454  32. Eustache S, Leprince J, Tufféry P. Progress with peptide scanning to study structure-

455  activity relationships: the implications for drug discovery. *Expert Opinion on Drug Discovery*.

456  2016; doi: 10.1080/17460441.2016.1201058.

457  33. Olson CA, Wu NC, Sun R. A Comprehensive Biophysical Description of Pairwise Epistasis

458  throughout an Entire Protein Domain. *Current Biology*. 2014; doi: 10.1016/j.cub.2014.09.072.

459  34. Staller MV, Holehouse AS, Swain-Lenz D, Das RK, Pappu RV, Cohen BA. A High-

460  Throughput Mutational Scan of an Intrinsically Disordered Acidic Transcriptional Activation

461  Domain. *Cell Systems*. 2018; doi: 10.1016/j.cels.2018.01.015.

462  35. Gray VE, Sitko K, Kameni FZN, Williamson M, Stephany JJ, Hasle N, et al.. Elucidating

463  the Molecular Determinants of Aβ Aggregation with Deep Mutational Scanning. *G3*

464  *(Bethesda)*. 2019; doi: 10.1534/g3.119.400535.

465  36. Esposito D, Weile J, Shendure J, Starita LM, Papenfuss AT, Roth FP, et al.. MaveDB: an

466  open-source platform to distribute and interpret data from multiplexed assays of variant effect.

467  *Genome Biol*. 2019; doi: 10.1186/s13059-019-1845-6.

468  37. Rubin AF, Min JK, Rollins NJ, Da EY, Esposito D, Harrington M, et al.. MaveDB v2: a

469  curated community database with over three million variant effects from multiplexed

470  functional assays. bioRxiv;

471  38. Heredia JD, Park J, Brubaker RJ, Szymanski SK, Gill KS, Procko E. Mapping Interaction

472  Sites on Human Chemokine Receptors by Deep Mutational Scanning. *The Journal of*

473  *Immunology*. American Association of Immunologists; 2018; doi: 10.4049/jimmunol.1800343.

474  39. Tian S, Choi W-T, Liu D, Pesavento J, Wang Y, An J, et al.. Distinct Functional Sites for

475  Human Immunodeficiency Virus Type 1 and Stromal Cell-Derived Factor 1α on CXCR4

476  Transmembrane Helical Domains. *JVI*. 2005; doi: 10.1128/JVI.79.20.12667-12673.2005.

477  40. Chabot DJ, Zhang P-F, Quinnan GV, Broder CC. Mutagenesis of CXCR4 Identifies

478  Important Domains for Human Immunodeficiency Virus Type 1 X4 Isolate Envelope-

479  Mediated Membrane Fusion and Virus Entry and Reveals Cryptic Coreceptor Activity for R5

480  Isolates. *J Virol*. 1999; doi: 10.1128/JVI.73.8.6598-6609.1999.

481  41. Han DP, Penn-Nicholson A, Cho MW. Identification of critical determinants on ACE2 for

482  SARS-CoV entry and development of a potent entry inhibitor. *Virology*. 2006; doi:

483  10.1016/j.virol.2006.01.029.

484   42. Fujita–Yoshigaki J, Shirouzu M, Ito Y, Hattori S, Furuyama S, Nishimura S, et al.. A

485   Constitutive Effector Region on the C-terminal Side of Switch I of the Ras Protein. *J Biol Chem*.

486   American Society for Biochemistry and Molecular Biology; 1995; doi: 10.1074/jbc.270.9.4661.

487   43. Gray VE, Hause RJ, Fowler DM. Analysis of Large-Scale Mutagenesis Data To Assess the

488   Impact of Single Amino Acid Substitutions. *Genetics*. 2017; doi: 10.1534/genetics.117.300064.

489   44. Hidalgo P, Ansari AZ, Schmidt P, Hare B, Simkovich N, Farrell S, et al.. Recruitment of

490   the transcriptional machinery through GAL11P: structure and interactions of the GAL4

491   dimerization domain. *Genes Dev*. 2001; doi: 10.1101/gad.873901.

492   45. Rodríguez-Escudero I, Oliver MD, Andrés-Pons A, Molina M, Cid VJ, Pulido R. A

493   comprehensive functional analysis of PTEN mutations: implications in tumor- and autism-

494   related syndromes. *Human Molecular Genetics*. 2011; doi: 10.1093/hmg/ddr337.

495   46. Schröter C, Günther R, Rhiel L, Becker S, Toleikis L, Doerner A, et al.. A generic approach

496   to engineer antibody pH-switches using combinatorial histidine scanning libraries and yeast

497   display. *mAbs*. 2015; doi: 10.4161/19420862.2014.985993.

498   47. Starace DM, Bezanilla F. Histidine Scanning Mutagenesis of Basic Residues of the S4

499   Segment of the Shaker K+ Channel. *J Gen Physiol*. 117:469–902001;

500   48. Stekhoven DJ, Buhlmann P. MissForest--non-parametric missing value imputation for

501   mixed-type data. *Bioinformatics*. 2012; doi: 10.1093/bioinformatics/btr597.

502   49. Wu Y, Weile J, Cote AG, Sun S, Knapp J, Verby M, et al.. A web application and service

503   for imputing and visualizing missense variant effect maps. Schwartz R, editor. *Bioinformatics*.

504   2019; doi: 10.1093/bioinformatics/btz012.

505    50. Cagiada M, Johansson KE, Valanciute A, Nielsen SV, Hartmann-Petersen R, Yang JJ, et

506    al.. Understanding the Origins of Loss of Protein Function by Analyzing the Effects of

507    Thousands of Variants on Activity and Abundance. Ozkan B, editor. *Molecular Biology and*

508    *Evolution*. 2021; doi: 10.1093/molbev/msab095.

509    51. Cagiada M, Bottaro S, Lindemose S, Schenstrøm SM, Stein A, Hartmann-Petersen R, et

510    al.. Discovering functionally important sites in proteins. bioRxiv;

511    52. Jepsen MM, Fowler DM, Hartmann-Petersen R, Stein A, Lindorff-Larsen K. Chapter 5 -

512    Classifying disease-associated variants using measures of protein activity and stability. In: Pey

513    AL, editor. *Protein Homeostasis Diseases*. Academic Press;

514    53. Matreyek KA, Stephany JJ, Ahler E, Fowler DM. Integrating thousands of PTEN variant

515    activity and abundance measurements reveals variant subgroups and new dominant negatives

516    in cancers. *Genome Med*. 2021; doi: 10.1186/s13073-021-00984-x.

517    54. Mighell TL, Thacker S, Fombonne E, Eng C, O'Roak BJ. An Integrated Deep-Mutational-

518    Scanning Approach Provides Clinical Insights on PTEN Genotype-Phenotype Relationships.

519    *The American Journal of Human Genetics*. 2020; doi: 10.1016/j.ajhg.2020.04.014.

520    55. Nielsen SV, Hartmann-Petersen R, Stein A, Lindorff-Larsen K. Multiplexed assays reveal

521    effects of missense variants in MSH2 and cancer predisposition. *PLOS Genetics*. Public

522    Library of Science; 2021; doi: 10.1371/journal.pgen.1009496.

523    56. The UniProt Consortium, Bateman A, Martin M-J, Orchard S, Magrane M, Agivetova R,

524    et al.. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Research*. 2021;

525    doi: 10.1093/nar/gkaa1100.

526 57. Andrews B, Fields S. Distinct patterns of mutational sensitivity for λ resistance and
527 maltodextrin transport in Escherichia coli LamB. *Microb Genom*. 2020; doi:
528 10.1099/mgen.0.000364.

529 58. Bandaru P, Shah NH, Bhattacharyya M, Barton JP, Kondo Y, Cofsky JC, et al..
530 Deconstruction of the Ras switching cycle through saturation mutagenesis. *eLife*. 2017; doi:
531 10.7554/eLife.27810.

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

534 60. Bridgford JL, Lee SM, Lee CMM, Guglielmelli P, Rumi E, Pietra D, et al.. Novel drivers
535 and modifiers of MPL-dependent oncogenic transformation identified by deep mutational
536 scanning. *Blood*. American Society of Hematology; 2020; doi: 10.1182/blood.2019002561.

537 61. Chan KK, Dorosky D, Sharma P, Abbasi SA, Dye JM, Kranz DM, et al.. Engineering
538 human ACE2 to optimize binding to the spike protein of SARS coronavirus 2. *Science*.
539 American Association for the Advancement of Science; 2020; doi: 10.1126/science.abc0870.

540 62. Chiasson MA, Rollins NJ, Stephany JJ, Sitko KA, Matreyek KA, Verby M, et al..
541 Multiplexed measurement of variant abundance and activity reveals VKOR topology, active
542 site and human variant impact. *Elife*. 2020; doi: 10.7554/eLife.58026.

543 63. Elazar A, Weinstein J, Biran I, Fridman Y, Bibi E, Fleishman SJ. Mutational scanning
544 reveals the determinants of protein insertion and association energetics in the plasma
545 membrane. Shan Y, editor. *eLife*. eLife Sciences Publications, Ltd; 2016; doi:
546 10.7554/eLife.12125.

547   64. Findlay GM, Daza RM, Martin B, Zhang MD, Leith AP, Gasperini M, et al.. Accurate

548   classification of BRCA1 variants with saturation genome editing. *Nature*. 2018; doi:

549   10.1038/s41586-018-0461-z.

550   65. Firnberg E, Labonte JW, Gray JJ, Ostermeier M. A Comprehensive, High-Resolution Map

551   of a Gene's Fitness Landscape. *Mol Biol Evol*. 2014; doi: 10.1093/molbev/msu081.

552   66. Hietpas RT, Jensen JD, Bolon DNA. Experimental illumination of a fitness landscape.

553   *Proceedings of the National Academy of Sciences*. 2011; doi: 10.1073/pnas.1016024108.

554   67. Hietpas RT, Bank C, Jensen JD, Bolon DNA. Shifting fitness landscapes in response to

555   altered environments. *Evolution*. 2013; doi: 10.1111/evo.12207.

556   68. Jiang L, Mishra P, Hietpas RT, Zeldovich KB, Bolon DNA. Latent Effects of Hsp90

557   Mutants Revealed at Reduced Expression Levels. *PLOS Genetics*. Public Library of Science;

558   2013; doi: 10.1371/journal.pgen.1003600.

559   69. Jiang RJ. Exhaustive Mapping of Missense Variation in Coronary Heart Disease-related

560   Genes [Thesis]. University of Toronto;

561   70. Keskin A, Akdoğan E, Dunn CD. Evidence for Amino Acid Snorkeling from a High-

562   Resolution, *In Vivo* Analysis of Fis1 Tail-Anchor Insertion at the Mitochondrial Outer

563   Membrane. *Genetics*. 2017; doi: 10.1534/genetics.116.196428.

564   71. Kitzman JO, Starita LM, Lo RS, Fields S, Shendure J. Massively parallel single-amino-

565   acid mutagenesis. *Nat Methods*. 2015; doi: 10.1038/nmeth.3223.

566   72. Kotler E, Shani O, Goldfeld G, Lotan-Pompan M, Tarcic O, Gershoni A, et al.. A

567   Systematic p53 Mutation Library Links Differential Functional Impact to Cancer Mutation

568    Pattern and Evolutionary Conservation. *Molecular Cell*. Elsevier; 2018; doi:

569    10.1016/j.molcel.2018.06.012.

570    73. Kowalsky CA, Whitehead TA. Determination of binding affinity upon mutation for type I

571    dockerin–cohesin complexes from Clostridium thermocellum and Clostridium cellulolyticum

572    using deep sequencing. *Proteins: Structure, Function, and Bioinformatics*. 2016; doi:

573    10.1002/prot.25175.

574    74. McLaughlin Jr RN, Poelwijk FJ, Raman A, Gosal WS, Ranganathan R. The spatial

575    architecture of protein function and adaptation. *Nature*. 2012; doi: 10.1038/nature11500.

576    75. Melamed D, Young DL, Gamble CE, Miller CR, Fields S. Deep mutational scanning of an

577    RRM domain of the Saccharomyces cerevisiae poly(A)-binding protein. *RNA*. 2013; doi:

578    10.1261/rna.040709.113.

579    76. Mishra P, Flynn JM, Starr TN, Bolon DNA. Systematic Mutant Analyses Elucidate General

580    and    Client-Specific    Aspects    of    Hsp90    Function.    *Cell    Reports*.    2016;    doi:

581    10.1016/j.celrep.2016.03.046.

582    77. Nedrud D, Coyote-Maestas W, Schmidt D. A large-scale survey of pairwise epistasis

583    reveals a mechanism for evolutionary expansion and specialization of PDZ domains. *Proteins:*

584    *Structure, Function, and Bioinformatics*. 2021; doi: 10.1002/prot.26067.

585    78. Newberry RW, Arhar T, Costello J, Hartoularos GC, Maxwell AM, Naing ZZC, et al..

586    Robust Sequence Determinants of α-Synuclein Toxicity in Yeast Implicate Membrane Binding.

587    *ACS Chem Biol*. 2020; doi: 10.1021/acschembio.0c00339.

588  79. Newberry RW, Leong JT, Chow ED, Kampmann M, DeGrado WF. Deep mutational

589  scanning reveals the structural basis for α-synuclein activity. *Nat Chem Biol*. 2020; doi:

590  10.1038/s41589-020-0480-6.

591  80. Roscoe BP, Bolon DNA. Systematic Exploration of Ubiquitin Sequence, E1 Activation

592  Efficiency, and Experimental Fitness in Yeast. *Journal of Molecular Biology*. 2014; doi:

593  10.1016/j.jmb.2014.05.019.

594  81. Sarkisyan KS, Bolotin DA, Meer MV, Usmanova DR, Mishin AS, Sharonov GV, et al..

595  Local fitness landscape of the green fluorescent protein. *Nature*. Nature Publishing Group;

596  2016; doi: 10.1038/nature17995.

597  82. Silverstein RA, Sun S, Verby M, Weile J, Wu Y, Roth FP. A systematic genotype-

598  phenotype map for missense variants in the human intellectual disability-associated gene GDI1.

599  bioRxiv;

600  83. Starita LM, Pruneda JN, Lo RS, Fowler DM, Kim HJ, Hiatt JB, et al.. Activity-enhancing

601  mutations in an E3 ubiquitin ligase identified by high-throughput mutagenesis. *PNAS*. 2013;

602  doi: 10.1073/pnas.1303309110.

603  84. Starita LM, Young DL, Islam M, Kitzman JO, Gullingsrud J, Hause RJ, et al.. Massively

604  Parallel Functional Analysis of BRCA1 RING Domain Variants. *Genetics*. 2015; doi:

605  10.1534/genetics.115.175802.

606  85. Starita LM, Islam MM, Banerjee T, Adamovich AI, Gullingsrud J, Fields S, et al.. A

607  Multiplex Homology-Directed DNA Repair Assay Reveals the Impact of More Than 1,000

608  BRCA1 Missense Substitution Variants on Protein Function. *The American Journal of Human*

609  *Genetics*. 2018; doi: 10.1016/j.ajhg.2018.07.016.

610 86. Suiter CC, Moriyama T, Matreyek KA, Yang W, Scaletti ER, Nishii R, et al.. Massively

611 parallel variant characterization identifies *NUDT15* alleles associated with thiopurine toxicity.

612 *Proc Natl Acad Sci USA*. 2020; doi: 10.1073/pnas.1915680117.

613 87. Sun S, Weile J, Verby M, Wu Y, Wang Y, Cote AG, et al.. A proactive genotype-to-patient-

614 phenotype map for cystathionine beta-synthase. *Genome Med*. 2020; doi: 10.1186/s13073-020-

615 0711-1.

616 88. Thompson S, Zhang Y, Ingle C, Reynolds KA, Kortemme T. Altered expression of a quality

617 control protease in E. coli reshapes the in vivo mutational landscape of a model enzyme. *eLife*.

618 2020; doi: 10.7554/eLife.53476.

619 89. Trenker R, Wu X, Nguyen JV, Wilcox S, Rubin AF, Call ME, et al.. Human and viral

620 membrane–associated E3 ubiquitin ligases MARCH1 and MIR2 recognize different features

621 of CD86 to downregulate surface expression. *Journal of Biological Chemistry*. Elsevier; 2021;

622 doi: 10.1016/j.jbc.2021.100900.

623 90. Weile J, Sun S, Cote AG, Knapp J, Verby M, Mellor JC, et al.. A framework for

624 exhaustively mapping functional missense variants. *Mol Syst Biol*. 2017; doi:

625 10.15252/msb.20177908.

626 91. Weile J, Kishore N, Sun S, Maaieh R, Verby M, Li R, et al.. Shifting landscapes of human

627 MTHFR missense-variant effects. *The American Journal of Human Genetics*. Elsevier; 2021;

628 doi: 10.1016/j.ajhg.2021.05.009.

629 92. Wrenbeck EE, Bedewitz MA, Klesmith JR, Noshin S, Barry CS, Whitehead TA. An

630 Automated Data-Driven Pipeline for Improving Heterologous Enzyme Expression. *ACS Synth

631 Biol*. American Chemical Society; 2019; doi: 10.1021/acssynbio.8b00486.

632  93. Zhang L, Sarangi V, Moon I, Yu J, Liu D, Devarajan S, et al.. CYP2C9 and CYP2C19:

633  Deep Mutational Scanning and Functional Characterization of Genomic Missense Variants.

634  *Clinical and Translational Science*. 2020; doi: https://doi.org/10.1111/cts.12758.

635  94. Zinkus-Boltz J, DeValk C, Dickinson BC. A Phage-Assisted Continuous Selection

636  Approach for Deep Mutational Scanning of Protein–Protein Interactions. *ACS Chem Biol*.

637  American Chemical Society; 2019; doi: 10.1021/acschembio.9b00669.

638  95. Bernier-Villamor V, Sampson DA, Matunis MJ, Lima CD. Structural Basis for E2-

639  Mediated SUMO Conjugation Revealed by a Complex between Ubiquitin-Conjugating

640  Enzyme Ubc9 and RanGAP. *Cell*. 108:122002;

641  96. Blanpain C, Doranz BJ, Vakili J, Rucker J, Govaerts C, Baik SSW, et al.. Multiple Charged

642  and Aromatic Residues in CCR5 Amino-terminal Domain Are Involved in High Affinity

643  Binding of Both Chemokines and HIV-1 Env Protein. *J Biol Chem*. 1999; doi:

644  10.1074/jbc.274.49.34719.

645  97. Brzovic PS, Keeffe JR, Nishikawa H, Miyamoto K, Fox D, Fukuda M, et al.. Binding and

646  recognition in the assembly of an active BRCA1/BARD1 ubiquitin-ligase complex.

647  *Proceedings of the National Academy of Sciences*. 2003; doi: 10.1073/pnas.0836054100.

648  98. Chen S, Wu J, Zhong S, Li Y, Zhang P, Ma J, et al.. iASPP mediates p53 selectivity through

649  a modular mechanism fine-tuning DNA recognition. *Proc Natl Acad Sci USA*. 2019; doi:

650  10.1073/pnas.1909393116.

651  99. Chupreta S, Holmstrom S, Subramanian L, Iñiguez-Lluhí JA. A Small Conserved Surface

652  in SUMO Is the Critical Structural Determinant of Its Transcriptional Inhibitory Properties.

653  *MCB*. 2005; doi: 10.1128/MCB.25.10.4272-4282.2005.

654     100. Cobb JA, Roberts DM. Structural Requirements for N-Trimethylation of Lysine 115 of

655     Calmodulin. *Journal of Biological Chemistry*. 2000; doi: 10.1074/jbc.M002332200.

656     101. Coyne RS, McDonald HB, Edgemon K, Brody LC. Functional Characterization of

657     BRCA1 Sequence Variants using a Yeast Small Colony Phenotype Assay. *Cancer Biology &*

658     *Therapy*. 2004; doi: 10.4161/cbt.3.5.809.

659     102. Denker K, Orlik F, Schiffler B, Benz R. Site-directed Mutagenesis of the Greasy Slide

660     Aromatic Residues Within the LamB (Maltoporin) Channel of Escherichia coli: Effect on Ion

661     and     Maltopentaose     Transport.     *Journal     of     Molecular     Biology*.     2005;     doi:

662     10.1016/j.jmb.2005.07.025.

663     103. Dragic T, Trkola A, Lin SW, Nagashima KA, Kajumo F, Zhao L, et al.. Amino-Terminal

664     Substitutions in the CCR5 Coreceptor Impair gp120 Binding and Human Immunodeficiency

665     Virus Type 1 Entry. *J Virol*. 1998; doi: 10.1128/JVI.72.1.279-285.1998.

666     104. Dragic T, Trkola A, Thompson DAD, Cormier EG, Kajumo FA, Maxwell E, et al.. A

667     binding pocket for a small molecule inhibitor of HIV-1 entry within the transmembrane helices

668     of     CCR5.     *Proceedings     of     the     National     Academy     of     Sciences*.     2000;     doi:

669     10.1073/pnas.090576697.

670     105. Ecsédi P, Gógl G, Hóf H, Kiss B, Harmat V, Nyitray L. Structure Determination of the

671     Transactivation Domain of p53 in Complex with S100A4 Using Annexin A2 as a

672     Crystallization Chaperone. *Structure*. 2020; doi: 10.1016/j.str.2020.05.001.

673     106. Kopecká J, Krijt J, Raková K, Kožich V. Restoring assembly and activity of cystathionine

674     β-synthase mutants by ligands and chemical chaperones. *Journal of Inherited Metabolic*

675     *Disease*. 2011; doi: 10.1007/s10545-010-9087-5.

676    107. Kožich V, Sokolová J, Klatovská V, Krijt J, Janošík M, Jelínek K, et al.. Cystathionine β-

677    synthase mutations: effect of mutation topology on folding and activity. *Hum Mutat*. 2010; doi:

678    10.1002/humu.21273.

679    108. Kruger W d., Wang L, Jhee K h., Singh R h., Elsas II L j.. Cystathionine β-synthase

680    deficiency in Georgia (USA): Correlation of clinical and biochemical phenotype with genotype.

681    *Human Mutation*. 2003; doi: 10.1002/humu.10290.

682    109. Lee SY, Pullen L, Virgil DJ, Castañeda CA, Abeykoon D, Bolon DNA, et al.. Alanine

683    Scan of Core Positions in Ubiquitin Reveals Links between Dynamics, Stability, and Function.

684    *Journal of Molecular Biology*. 2014; doi: 10.1016/j.jmb.2013.10.042.

685    110. Li W, Zhang C, Sui J, Kuhn JH, Moore MJ, Luo S, et al.. Receptor and viral determinants

686    of    SARS-coronavirus    adaptation    to    human    ACE2.    *EMBO    J*.    2005;    doi:

687    10.1038/sj.emboj.7600640.

688    111. Lin G, Baribaud F, Romano J, Doms RW, Hoxie JA. Identification of gp120 Binding Sites

689    on CXCR4 by Using CD4-Independent Human Immunodeficiency Virus Type 2 Env Proteins.

690    *JVI*. 2003; doi: 10.1128/JVI.77.2.931-942.2003.

691    112. Mascle XH, Lussier-Price M, Cappadocia L, Estephan P, Raiola L, Omichinski JG, et al..

692    Identification of a Non-covalent Ternary Complex Formed by PIAS1, SUMO1, and UBC9

693    Proteins Involved in Transcriptional Regulation. *Journal of Biological Chemistry*. 2013; doi:

694    10.1074/jbc.M113.486845.

695    113. Matthews EE, Thévenin D, Rogers JM, Gotow L, Lira PD, Reiter LA, et al..

696    Thrombopoietin receptor activation: transmembrane helix dimerization, rotation, and allosteric

697    modulation. *The FASEB Journal*. 2011; doi: https://doi.org/10.1096/fj.10-178673.

698    114. Mayfield JA, Davies MW, Dimster-Denk D, Pleskac N, McCarthy S, Boydston EA, et al.. 

699    Surrogate Genetics and Metabolic Profiling for Characterization of Human Disease Alleles. 

700    *Genetics*. 2012; doi: 10.1534/genetics.111.137471.

701    115. Navenot J-M, Wang Z, Trent JO, Murray JL, Hu Q, DeLeeuw L, et al.. Molecular anatomy 

702    of CCR5 engagement by physiologic and viral chemokines and HIV-1 envelope glycoproteins: 

703    differences in primary structural requirements for RANTES, MIP-1α, and vMIP-II 

704    binding11Edited by P. E. Wright. *Journal of Molecular Biology*. 2001; doi: 

705    10.1006/jmbi.2001.5086.

706    116. Peng L, Damschroder MM, Cook KE, Wu H, Dall'Acqua WF. Molecular basis for the 

707    antagonistic activity of an anti-CXCR4 antibody. *mAbs*. 2016; doi: 

708    10.1080/19420862.2015.1113359.

709    117. Peterson BR, Sun LJ, Verdine GL. A critical arginine residue mediates cooperativity in 

710    the contact interface between transcription factors NFAT and AP-1. *Proceedings of the 

711    National Academy of Sciences*. 1996; doi: 10.1073/pnas.93.24.13671.

712    118. Rabut GEE, Konner JA, Kajumo F, Moore JP, Dragic T. Alanine Substitutions of Polar 

713    and Nonpolar Residues in the Amino-Terminal Domain of CCR5 Differently Impair Entry of 

714    Macrophage- and Dualtropic Isolates of Human Immunodeficiency Virus Type 1. *J Virol*. 1998; 

715    doi: 10.1128/JVI.72.4.3464-3468.1998.

716    119. Ransburgh DJR, Chiba N, Ishioka C, Toland AE, Parvin JD. Identification of Breast 

717    Tumor Mutations in *BRCA1* That Abolish Its Function in Homologous DNA Recombination. 

718    *Cancer Res*. 2010; doi: 10.1158/0008-5472.CAN-09-2850.

719    120. Tan Y, Tong P, Wang J, Zhao L, Li J, Yu Y, et al.. The Membrane-Proximal Region of

720    C–C Chemokine Receptor Type 5 Participates in the Infection of HIV-1. *Front Immunol*. 2017;

721    doi: 10.3389/fimmu.2017.00478.

722    121. Towler WI, Zhang J, Ransburgh DJR, Toland AE, Ishioka C, Chiba N, et al.. Analysis of

723    BRCA1 Variants in Double-Strand Break Repair by Homologous Recombination and Single-

724    Strand Annealing. *Human Mutation*. 2013; doi: 10.1002/humu.22251.

725    122. Trent JO, Wang Z, Murray JL, Shao W, Tamamura H, Fujii N, et al.. Lipid Bilayer

726    Simulations of CXCR4 with Inverse Agonists and Weak Partial Agonists. *J Biol Chem*. 2003;

727    doi: 10.1074/jbc.M307850200.

728    123. Van Gelder P, Dumas F, Bartoldus I, Saint N, Prilipov A, Winterhalter M, et al.. Sugar

729    Transport through Maltoporin of *Escherichia coli* : Role of the Greasy Slide. *J Bacteriol*. 2002;

730    doi: 10.1128/JB.184.11.2994-2999.2002.

731    124. VanBerkum MF, Means AR. Three amino acid substitutions in domain I of calmodulin

732    prevent the activation of chicken smooth muscle myosin light chain kinase. *J Biol Chem*.

733    American Society for Biochemistry and Molecular Biology; 266:21488–951991;

734    125. Wei Q, Wang L, Wang Q, Kruger WD, Dunbrack RL. Testing computational prediction

735    of missense mutation phenotypes: Functional characterization of 204 mutations of human

736    cystathionine beta synthase. *Proteins: Structure, Function, and Bioinformatics*. 2010; doi:

737    10.1002/prot.22722.

738    126. Williams AD, Shivaprasad S, Wetzel R. Alanine Scanning Mutagenesis of Aβ(1-40)

739    Amyloid Fibril Stability. *Journal of Molecular Biology*. 2006; doi: 10.1016/j.jmb.2006.01.041.

740  127. Zhang J, Rao E, Dioszegi M, Kondru R, DeRosier A, Chan E, et al.. The Second

741  Extracellular Loop of CCR5 Contains the Dominant Epitopes for Highly Potent Anti-Human

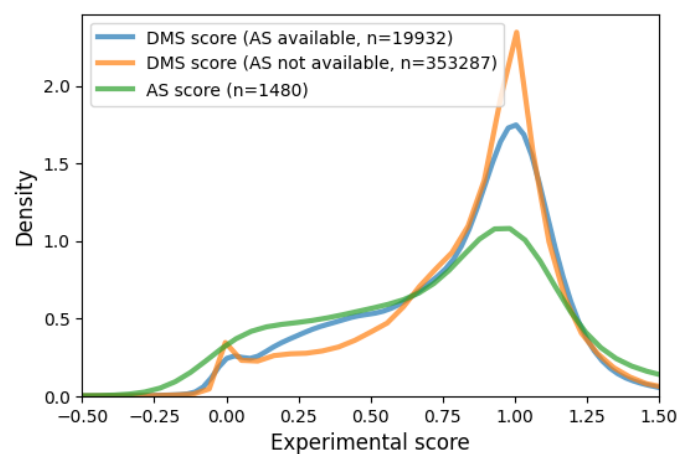742  Immunodeficiency Virus Monoclonal Antibodies. *AAC*. 2007; doi: 10.1128/AAC.01302-06.

743  128. Nelsen RB. An introduction to copulas. 2nd ed. New York: Springer;

744  129. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al.. Scikit-

745  learn: Machine Learning in Python. *Journal of machine Learning research*. :2825–30 2011;

746  130. González J, Dai Z, Hennig P, Lawrence ND. Batch Bayesian Optimization via Local

747  Penalization. arXiv;

748

749  **Supplementary material**



750

751  **Fig S1.  DMS and AS score distribution.** The figure shows the kernel estimated density of normalized AS scores

752  and DMS scores for variants with or without available AS data.

753

For each **pair** of DMS and AS experiments:

754

**Fig S2.  Decision tree for classifying the DMS and AS assay compatibility.** The end-nodes show the classified

assay compatibility. The number indicates the count of assay pairs for each compatibility level (low, medium,
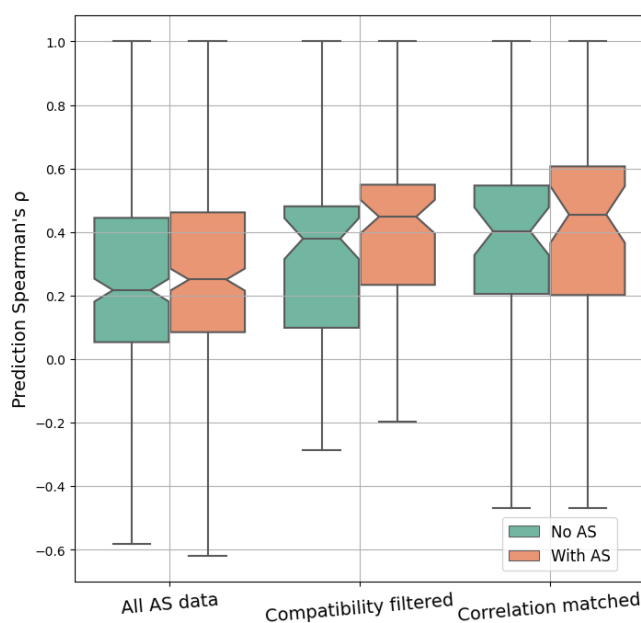
high).

758



759

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

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

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

"DeMaSk data & code" result is similar to the published result. For the "Our data & DeMaSk code" result, we

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

765 many more DMS results are included in our data. The similarity of results achieved using "Our data & code"

766 demonstrates the correctness of our re-implementation.
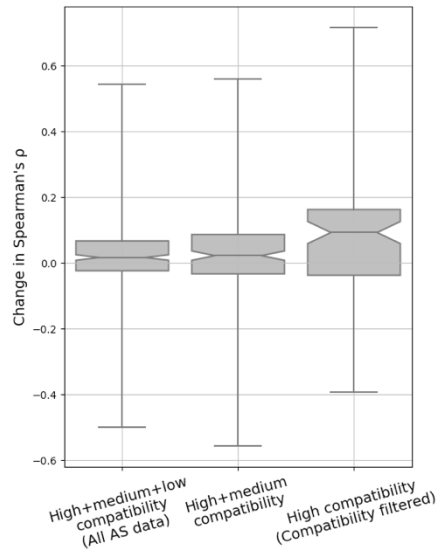
767

768



769 **Fig S4. Performance comparison between predictors using AS data or not.** The Spearman's $\rho$ between

770 experiment DMS scores and predicted scores for each DMS and AS data pair are shown as box plots. Different

771 approaches to filtering/matching the data are shown on the x-axis: "All AS data" used all available data;

772 "Compatibility filtered" used only data of high assay compatibility; "Correlation matched" used only data with

773 the highest regularised correlation for each DMS dataset. The figure does not include data without available

774 (filtered/matched) AS scores. This means that the different results are not directly comparable since they are

775 visualized on different subsets of DMS/AS data pairs (for example, "All AS data" contains all DMS/AS data pairs,

776 but "Compatibility filtered" contains only data pairs of high assay compatibility). Control results are shown as

777 green boxes for predicting without AS data as a feature. The underlying $\rho$ for each data pair in the control results

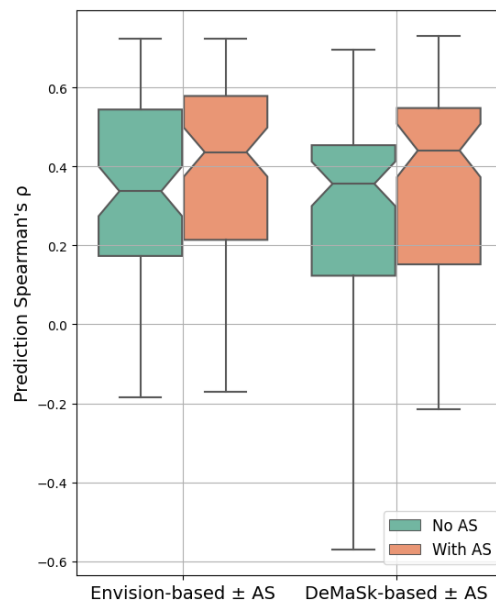778 is the same, but the boxes are shifted due to data filtering/matching.

779

780

**Fig S5. The performance of variant impact prediction for using data of different assay compatibility levels.** The change of prediction Spearman's $\rho$ for each DMS and AS data pair is shown as box plots. A higher value represents higher prediction accuracy achieved for using AS data. Different data filtering methods are shown on the x-axis.



785

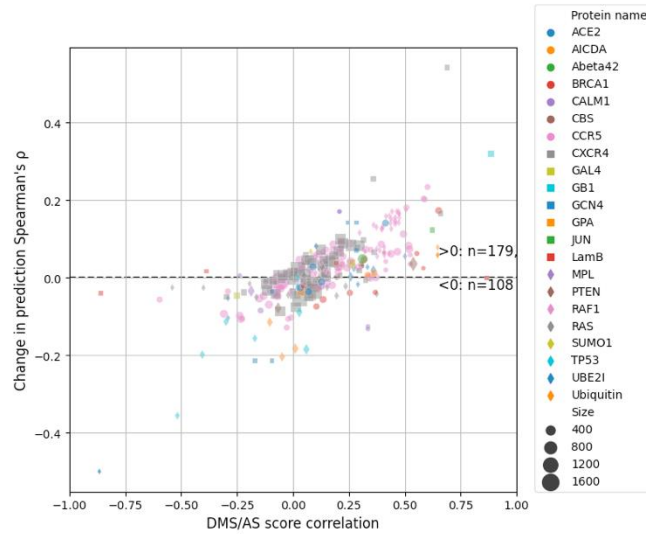**Fig S6. Prediction performance is improved while incorporating high compatibility AS data into the Envision model.** The Spearman's $\rho$ between experiment DMS scores and predicted scores for each high compatible DMS/AS assay pair are shown as box plots. The x-axis shows the predictor used, either Envision or DeMaSk. Control results are shown as green boxes for predicting without AS data as a feature.
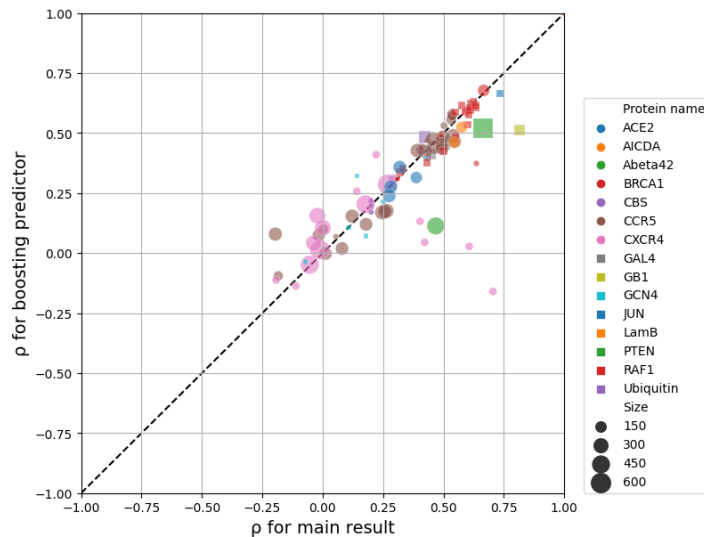
790

791



792 **Fig S7. Prediction performance change for using all AS data.** Each dot represents a DMS/AS data pair. The

793 vertical axis shows the change of prediction $\rho$ by using AS data (larger means higher performance achieved by

794 using AS data). The horizontal axis shows the DMS/AS score correlation for *all* variants on the matched residues

795 rather than just alanine substitutions. The colours and shapes of the dots correspond to the target protein, and size

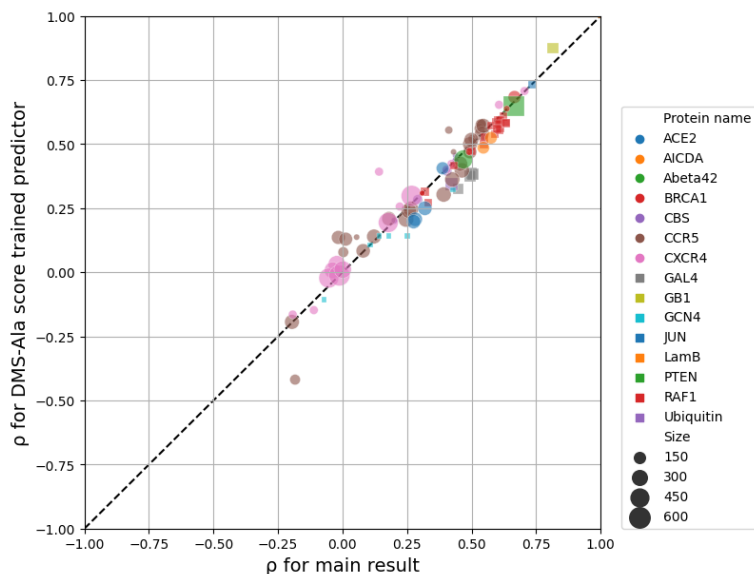796 indicates the number of variants in each data pair.

797



798

799 **Fig S8. Boosting setup shows similar performance as the main result.** Each dot represents a filtered DMS/AS

800 data pair of high assay compatibility. The vertical and horizontal axes show the prediction Spearman's $\rho$ for either

801 modelled with boosting or the one-step (main result) setup. The colours and shapes of the dots correspond to the

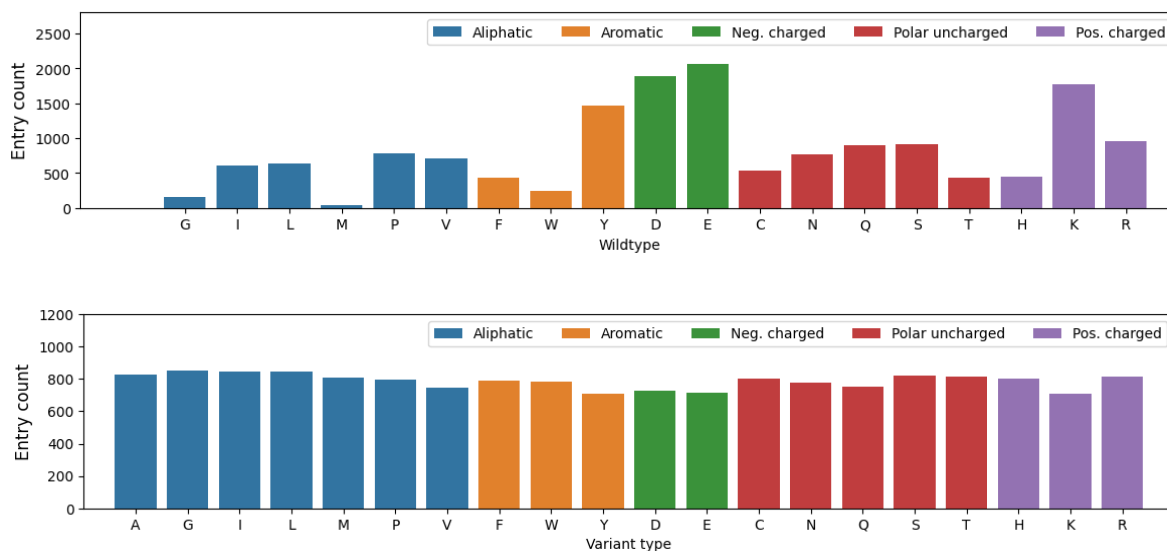802 target protein, and size indicates the number of variants in each data pair.

803



804

**Fig S9.  Training with DMS scores of alanine substitutions shows similar performance as the main result.**
The vertical and horizontal axes show the prediction Spearman's $\rho$ for predictors either trained with DMS score
of alanine substitutions (DMS-Ala) or AS data of high assay compatibility (main result), yet all evaluated on high
compatibility AS data. The colours and shapes of the dots correspond to the target protein, and size indicates the
number of variants in each data pair.

810



811



812

**Fig S10. Count of variant entries for each wild-type or variant amino acid of high assay compatibility data.**

(Neg.: negatively, Pos.: positively)

815

**Table S1. Amount of data with AS scores available**

| Data composition | Protein | DMS dataset | AS dataset[1] | Variant entries[2] |
|---|---|---|---|---|
| All AS | 22 | 54 | 146 | 70446 |
| Compatibility filtered | 15 | 35 | 60 | 15739 |
| High+medium assay compatibility | 21 | 51 | 105 | 28380 |
| Correlation matched | 22 | 54 | 32 | 7940 |

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

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

819

820 **Supplementary information**

821 **Applying AS data to Envision method**

822 We re-implemented a predictor based on Envision [15] to incorporate AS data. Features used

823 in Envision were downloaded from its online toolkit. All Envision features are used for

824 modelling except for substitution type (wt_mut) which has low importance according to the

825 published result and our pilot studies yet is computationally expensive in our setup. Protein

826 data were excluded if their features were not available online. DMS and AS data pairs with

827 high assay compatibility were used for modelling. Missing feature values were imputed by the

828 mean values for numerical features or the most frequent values for categorical features.

829 Categorical features are encoded with the one-hot encoder. We used

830 `sklearn.ensemble.GradientBoostingRegressor` from scikit-learn package [129]

831 to build the predictor, and hyperparameters were tuned by Bayesian Optimization [130] with

832 Group K-Fold (protein-30-fold) cross-validation. The training and evaluation process were

833 similar to that previously described. For comparison, we repeated the DeMaSk-based analysis

834 on the same subset of data.

835

**Boosting with AS data**

837   To deal with the sparsity of AS data, we tested a variant impact predictor based on boosting. A

838   first linear regression predictor was trained with all training DMS data using the three DeMaSk

839   features without AS data, which was the same as the control predictor mentioned previously.

840   We then calculated the prediction error by subtracting the predicted scores from DMS scores,

841   and a second linear regression predictor was trained to predict the error. The second predictor

842   was trained only on DMS/AS data of high assay compatibility and used both protein features

843   and the encoded AS scores. The final prediction result was the sum of the outputs from these

844   two predictors.

845

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

847   We investigated another potential approach to overcome the sparsity of AS data by replacing

848   the AS feature with the DMS scores of alanine substitutions (DMS-Ala). For all DMS datasets

849   we collected, their AS feature values, regardless of availability, were replaced by the DMS-

850   Ala scores on the same residue. Missing scores were imputed by the mean value of all DMS-

851   Ala scores. A regression model was trained and evaluated as previously described, using the

852   three DeMaSk features as well as the DMS-Ala scores. The AS data of high assay compatibility

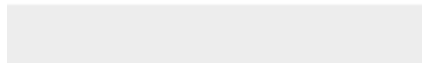853   are still used for the testing process.

854

Supplementary Table 1

Click here to access/download
**Supplementary Material**
Supplementary_Table_1_Supplementary Matrerial.xlsx
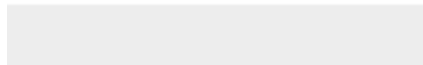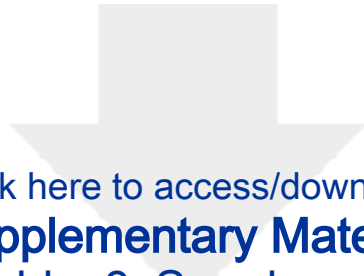
Click here to access/download

**Supplementary Material**

Supplementary_Table_2_Supplementary Matrerial.csv

Click here to access/download
**Supplementary Material**
Supplementary_Table_3_Supplementary Matrerial.csv

**WEHI**
brighter together

The Walter and Eliza Hall Institute of Medical Research
ABN 12 004 251 423

1G Royal Parade Parkville Victoria 3052 Australia
T +61 3 9345 2555  F +61 3 9347 0852
www.wehi.edu.au

Dr Scott Edmunds
Editor-in-Chief
Gigascience

Dear Dr Edmunds,

Please find our enclosed manuscript entitled "**Integrating deep mutational scanning and low-throughput mutagenesis data to predict the impact of amino acid variants**" for your consideration for publication in *Gigascience*.

The key contributions of our work are:
- We developed the first predictor of protein variant impact integrating high-throughput and low-throughput mutagenesis data, in our case, deep mutational scanning and alanine scan data.
- We demonstrate that integrative variant impact predictors improve model performance only when the high and low throughput data are generated by related assay types.

In this work, we collected high-throughput deep mutational scanning (DMS) data from an online database, with 370,000 protein variants and low throughput alanine scanning data of matched proteins from published papers. We defined a decision tree to classify low- and high-throughput assays to distinct levels of similarity across multiple categories, which we call assay compatibility. We then explored models of variant impact trained on these data. Our results showed the connection between experiment assay compatibility and the predictor's performance built from these data.

This is an original research article, and we have no conflicts of interest to disclose. All authors have participated in the preparation of this manuscript and approved the submission of it. We confirm that this work has not been published nor is currently under consideration for publication elsewhere.

Thank you for your consideration of this manuscript.

Yours sincerely,

Mr Yunfan Fu          Prof Anthony T. Papenfuss     Dr Alan F. Rubin