1    **Title:** Inverse folding of protein complexes with a structure-informed language model enables

2    unsupervised antibody evolution

3

4    **Authors:** Varun R. Shanker[1,2,3], Theodora U.J. Bruun[2,3,4], Brian L. Hie[3,4]*, Peter S. Kim[3,4,5]*

5

6    **Affiliations**:

7    [1]Stanford Biophysics Program, Stanford University School of Medicine, Stanford, CA 94305,

8    USA

9    [2]Stanford Medical Scientist Training Program, Stanford University School of Medicine, Stanford

10    CA 94305, USA

11    [3]Sarafan ChEM-H, Stanford University, Stanford, CA 94305, USA

12    [4]Department of Biochemistry, Stanford University School of Medicine, Stanford, CA 94305,

13    USA

14    [5]Chan Zuckerberg Biohub, San Francisco, CA 94158, USA.

15    *Correspondence: B.L.H. (brianhie@stanford.edu), P.S.K. (kimpeter@stanford.edu)

16

17    **Abstract**

18    Large language models trained on sequence information alone are capable of learning high level

19    principles of protein design. However, beyond sequence, the three-dimensional structures of

20    proteins determine their specific function, activity, and evolvability. Here we show that a general

21    protein language model augmented with protein structure backbone coordinates and trained on

22    the inverse folding problem can guide evolution for diverse proteins without needing to

23    explicitly model individual functional tasks. We demonstrate inverse folding to be an effective

24      unsupervised, structure-based sequence optimization strategy that also generalizes to multimeric

25      complexes by implicitly learning features of binding and amino acid epistasis. Using this

26      approach, we screened ~30 variants of two therapeutic clinical antibodies used to treat SARS-

27      CoV-2 infection and achieved up to 26-fold improvement in neutralization and 37-fold

28      improvement in affinity against antibody-escaped viral variants-of-concern BQ.1.1 and XBB.1.5,

29      respectively. In addition to substantial overall improvements in protein function, we find inverse

30      folding performs with leading experimental success rates among other reported machine

31      learning-guided directed evolution methods, without requiring any task-specific training data.

**Introduction**

Evolution generates diverse proteins at the level of biological sequences by exploring a

vast search space of potential mutations and acquiring those that improve fitness. However, it is

the three-dimensional structure encoded by these sequences that ultimately determines the

function and activity of a protein. Consequently, as proteins accumulate mutations, they undergo

corresponding structural changes, which in turn facilitate functional adaptations[1].

In the laboratory, this tendency for greater sequence change to cause structural

divergence poses a major challenge to engineering better proteins via a stepwise evolutionary

process. Mutations added in sequential rounds of artificial evolution are increasingly likely to

destabilize the structure and therefore diminish the protein's evolvability[2]. Identifying beneficial

mutations is further challenged by the fact that almost all mutations to a prototypical protein are

deleterious, or at best neutral, and only a rare subset are beneficial on its fitness landscape[3–8]. In

total, these phenomena can often reduce the evolutionarily accessible paths and make evolution

more susceptible to local fitness optima[9,10], further complicating attempts to increase fitness.

To address both the structural constraints of protein design and the high dimensionality of

the mutational search space, we utilized a general protein language model augmented with

structural information and trained across millions of non-redundant single sequence-structure

pairs on the inverse folding objective[11]. Most simply, the inverse folding problem considers the

task opposite of that performed by many of the recent powerful structure-prediction tools,

including AlphaFold and ESMFold[12,13]: recovery of a protein's native sequence, given its three-

dimensional backbone coordinates (**Figure 1a**). This is accomplished by predicting the identity

of an amino acid given both the preceding amino acid sequence (referred to as autoregressive

modeling) and the entire structure's backbone coordinates (**Methods**). Thus, sequences assigned

3

60    high likelihood scores by the inverse folding language model are expected to fold into the

61    backbone of the input structure with high confidence (**Figure 1b**).

62        Our inverse folding framework for protein design does not model an explicit protein

63    function or definition of protein fitness. Rather, using a structure-guided paradigm, we indirectly

64    explore the underlying fitness landscape by focusing exploration to regions where the backbone

65    fold of the protein is preserved. We hypothesize constraining evolution to regimes of high

66    inverse folding likelihood can serve as an effective prior for high-fitness variants, and thereby

67    improve the efficiency of evolution (**Figure 1c**).

68        We reasoned that this approach may be particularly valuable for the evolution of human

69    antibodies, which are used clinically to treat a broad range of diseases[14]. Antibodies are used

70    therapeutically to bind to a target antigen mediating pathogenesis, and modify or disrupt its

71    function[15] . A central concept of this study is to use the complete structure of the antibody-

72    antigen complex to guide evolution. By conditioning the inverse folding model on the entire

73    antibody-antigen complex, we sought to enable the discovery of mutations that preserve or

74    enhance the stability of the entire complex, and thus that improve antibody function.

75        Indeed, we show that as an unsupervised machine learning-guided evolution strategy,

76    inverse folding is capable of identifying high fitness mutations across several protein families

77    and tasks, performing better than sequence-only methods. We found that inverse folding

78    generalizes to protein complexes with improved antibody variant prediction when antigen

79    structural information is also included as input. To demonstrate the practical utility of this

80    method, we improved the potency of mature, clinical SARS-CoV-2 monoclonal antibody

81    therapies, in a low-throughput setting, against both their original viral target as well as viral

82    escape variants that reduced their efficacy, namely variants-of-concern (VOC) BQ.1.1 and

83   XBB.1.5. We achieved up to 26-fold improvement in the neutralization potency of Ly-1404

84   (Bebtelovimab) against BQ.1.1, and 11-fold for SA58, testing only a total of 31 and 25 antibody

85   variants, respectively. We also achieved 27-fold improvement in affinity against BQ.1.1 and 37-

86   fold improvement in affinity against XBB.1.5. Notably, all experimentally tested combinations

87   of inverse folding-recommended mutations showed improved activity, with many designs

88   comprising multiple synergistic mutations. With our approach, we report experimental success

89   rates that surpass all previous machine learning-guided protein evolution methods[8,16–28], including

90   those based on supervision with task-specific training data. These findings highlight the

91   advantage of an unsupervised, structure-based paradigm to identify efficient evolutionary

92   trajectories.

93

94   **Results**

95   *Inverse folding enriches sequence exploration for high function protein variants across diverse*

96   *tasks*

97          We evaluated whether inverse folding can be used to guide protein evolution, without

98   needing to explicitly model specific functional tasks, by assessing its ability to identify mutations

99   resulting in high levels of protein activity for a desired functional property, or fitness measure.

100  Accordingly, for 10 proteins from diverse families among four organisms, and with functions

101  ranging from enzyme catalysis (TPMT) to oncogenesis (HRAS) to transcriptional regulation

102  (GAL4), we used inverse folding likelihoods to score variants profiled in large datasets from

103  deep mutational scanning experiments[29–38]against a target backbone of the wild-type protein[39–48]

104  (**Methods**, **Supplementary Table 1**).

105        From the thousands of tested variants for each of the 10 proteins, we identified numerous

106        with experimentally determined protein activities ranking in the top percentiles of the entire

107        screen within just the set of top ten inverse-folding predictions (**Figure 1d**). Our analysis also

108        demonstrates that conditioning on structural information serves to improve predictive capabilities

109        of protein language models as we successfully identified mutations in the top fifth percentile for

110        9 out of the 10 proteins using inverse folding compared to just 2 proteins using a state-of-the-art

111        general protein language model trained only on sequence information and specifically for variant

112        prediction (ESM-1v)[49] (**Figure 1d**). This improvement in prediction also holds with increasingly

113        relaxed thresholds for classification as high-fitness variants.

114        These results suggest that inverse folding offers a promising alternative to brute force

115        experimental searches for beneficial mutations. Notably, some of the top mutations predicted by

116        inverse folding are also the same ones recovered from exhaustive experimental exploration. For

117        example, for restriction enzyme haeIIIM, variant Q18E is recommended within the top five

118        inverse folding predictions and experimentally ranks as the second-best substitution (and $> 5$

119        standard deviations above the mean) out of the nearly 2000 substitutions screened to the

120        endonuclease[38]. Another key advantage of our task-independent framework, in addition to being

121        broadly applicable across diverse proteins, is the ability to improve a single protein for multiple

122        desired properties without needing to develop specialized high-throughput assays to screen each

123        independently. From just the top 10 inverse folding predictions for MAPK1, we identify

124        substitutions Q105M and Y64D, which are experimentally shown to confer resistance to two

125        different oncogenic-targeting MAPK1 kinase inhibitors[32].

126

127        *Inverse folding is a state-of-the-art zero-shot mutational effect predictor for antibodies*

128    To analyze the effectiveness of augmenting a general protein language model with

129    structural information, specifically for antibody variant prediction, we compared the inverse

130    folding likelihoods of sequences across entire mutational landscapes against the corresponding

131    experimental fitness values from three existing mutagenesis datasets. The first two of the datasets

132    profile the scFv equilibrium dissociation constants ($K_D$) of all possible evolutionary intermediates

133    between the inferred germline and somatic sequence of naturally affinity-matured influenza

134    broadly neutralizing antibodies (bnAbs) CR9114 and CR6261, which bind the conserved stem

135    epitope of influenza surface protein hemagglutinin (HA)[50]. For both bnAbs, only mutations in the

136    heavy chain, which is responsible for antigen binding, were characterized. The profiled

137    mutational landscape of CR9114 includes all possible combinations of 16 substitutions while that

138    of CR6261 includes all possible combinations of 11 substitutions, totaling $2^{16} = 65,536$ and $2^{11} =$

139    2,048 variant antibody sequences respectively. Each of these libraries were screened for binding

140    against two distinct influenza HA subtypes (H1 and H3 for CR9114 and H1 and H9 for

141    CR6261). The third dataset assesses the effects of all possible single amino acid substitutions

142    with a deep mutational scan profiling 4,275 mutations in the variable regions for both heavy

143    chain (VH) and light chain (VL) of antibody G6.31 to binding with its ligand, vascular

144    endothelial growth factor A (VEGF-A)[51].

145    For each dataset, we computed the Spearman correlation between the log likelihood

146    estimated by the inverse folding model and the experimentally determined binding measure for a

147    given antigen, across all sequences in the mutational library. We scored the inverse folding

148    likelihood of each candidate sequence in the library using the backbone coordinates of a structure

149    with the mature antibody bound to its target antigen[52–54].

7

150        Across all five experimental binding datasets, we found that inverse folding performs

151    better than both a sequence-only language model, ESM-1v[49], and a site-independent model of

152    mutational frequency curated with extensive antibody sequence alignments, abYsis[55]. In nearly

153    all experimental scenarios, supplementing sequence information with the backbone coordinates

154    of the antibody alone, without providing antigen information, as input to inverse folding is

155    sufficient to outperform other sequence-only methods. A notable feature of the autoregressive

156    architecture is that it computes the joint likelihood over all positions in a sequence, making it

157    well-suited to score combinatorial sequence changes. We find that inverse folding can capture

158    complex epistatic interactions, or potential interdependence among individual amino acids, as it

159    performs well on the CR9114 and CR6261 libraries composed of sequences with multiple

160    mutations (**Figure 2a,b**).

161        We achieved the greatest improvement in performance on all five experimental screens

162    by incorporating the structure of both the antibody and antigen (**Figure 2a**), indicating that the

163    inverse folding model can implicitly learn features of binding (**Figure 2c**). This result is

164    particularly significant, given that the inverse folding model is only trained on single-chain

165    protein structures, while the antibody-antigen complexes we use as inputs are composed of either

166    three (G6.31) or four (CR9114, CR6261) protein chains. The most substantial contribution of

167    antigen information is observed in the case of CR9114-H1, for which the correlation increases

168    from 0.17 with only antibody information to 0.65 with sequence and backbone coordinates of the

169    entire complex.

170        Remarkably, we could still predict effects of mutations on binding for a cross-reactive

171    antibody while using a different antigen as input to the model. (**Figure 2a,b**). Despite using a

172    complex with HA from H5N1 influenza as input to score CR9114 variants, we obtain

8

173    correlations of 0.65 and 0.50 with experimental binding data for H1 and H3, respectively. This is

174    particularly striking since, for example, H5 and H1 only share 63% sequence identify across both

175    HA subunits (**Supplementary Figure 3**). This same cross-reactive predictive capability is

176    observed for CR6261, which is tested experimentally against H1 and H9 while we use an input

177    structure with HA from 1918 H1N1 influenza (**Figure 2a**). Although inverse folding cannot learn

178    explicit chemical rules of binding (e.g., hydrogen bonding or disulfide bridge formation) since it

179    does not have access to amino acid side chain atomic coordinates, these results suggest that

180    structural principles like interface packing or potential steric interference are not only implicitly

181    accessible from residue identities, but are also informative for binding prediction.

182         Our model's top recommended mutations are made independently of a specific definition

183    of fitness; they simply represent a set of variants with a high likelihood of folding into the input

184    backbone structure. Therefore, our model's recommendations may also help identify mutations

185    that improve other useful biochemical properties beyond affinity. Impressively, for example, the

186    top inverse folding-recommended mutation to the VL of G6.31 is F83A, which was identified in

187    the original screening study[51] to be particularly interesting as it confers a three-fold increase in

188    VEGF-A binding affinity and a 5°C improvement in melting temperature, despite being 25Å

189    from the antigen and in the antibody framework region. It was determined that the VL F83A

190    substitution induces more compact packing and the site serves as a conformational switch that

191    affects biological activity at the antibody-antigen interface by modulating both interdomain and

192    elbow angle dynamics[51].

193

194    *Engineering therapeutic antibodies for increased potency and resilience*

195  Finally, we aimed to assess if the structure-augmented language model's predictive

196 capabilities could not only resolve trends on large sets of experimental data, but also enable

197 efficient and successful directed evolution campaigns while testing only a small number (on the

198 order of tens) of variants. To do so, we considered the task of improving the potency and

199 resilience (effectiveness against a virus as it mutates over time) of two mature, clinical

200 monoclonal antibody therapies.

201

202 &bull; Ly-1404 (Bebtelovimab) was isolated from a COVID-19 convalescent donor and binds to

203  the receptor binding domain (RBD) of the SARS-CoV-2 Spike protein[56]. It was approved

204  by the U.S. F.D.A. on February 11, 2022 given its activity against both the original

205  Wuhan and Omicron SARS-CoV-2 variants and was the last remaining approved

206  monoclonal antibody therapy withstanding against viral evolution[57] until its

207  discontinuation on November 30, 2023 due to antibody evasion by VOC BQ.1.1.[58]

208 &bull; SA58 (BD55-5840) was isolated from a vaccinated individual and is one of two RBD-

209  targeting neutralizing antibodies (NAb) in a rationally developed antibody cocktail. SA58

210  alone retained efficacy against all Omicron subvariants, including *in vivo* protection

211  against BA.5[59,60] and was shown to be effective as a post-exposure prophylaxis in a

212  clinical study[61].

213

214  For both antibody engineering campaigns, we used the inverse folding language model to

215 compute likelihoods of all ~4,300 possible single-residue substitutions in the VH or VL regions

216 of the antibody. In the first round of evolution, we selected only the top ten predictions at unique

217 residues in each chain for experimental validation. An important practical benefit of our method

218   is the ability to optimize against measures of fitness most relevant to the protein's downstream

219   function, rather than being limited to indirect and less accurate surrogate measures that are more

220   amenable to high-throughput screening[4,16]. We leverage this advantage to directly evolve these

221   antibodies for their ability to more potently neutralize SARS-CoV-2 pseudotyped lentivirus.

222          Variants recommended by the inverse folding language model were assessed by

223   comparing the half-maximal inhibitory concentration ($IC_{50}$) relative to the wild-type antibody.

224   Remarkably, although we chose to only test 20 single-site substitutions for each of the two

225   clinical monoclonal antibody therapies, approximately one-third of them improved neutralizing

226   potency. Notably, several of these variants improve neutralization $IC_{50}$ by approximately 2-fold

227   with just a single amino acid change (**Figure 3a**, **Supplementary Data 1**).

228          Prompted by recent evidence showing that conservation of the overall RBD structure is

229   robust to SARS-CoV-2 evolution[62], we next sought to determine whether we could also evolve

230   the previously mature antibodies against SARS-CoV-2 BQ.1.1, the variant responsible for

231   diminished therapeutic efficacy. Although the antibodies were previously effective, a change in

232   antigen conceptually represents a fundamental shift in the underlying fitness landscape (**Figure**

233   **3b**). From the same set of 20 single amino acid substitutions to Ly-1404, we found that nearly

234   half improve neutralization of variant BQ.1.1. In addition to a high success rate, we also found

235   multiple of these mutations provided a large magnitude of improvement. Several single amino

236   acid substitutions to Ly-1404 individually result in over a 3-fold improvement while the most

237   beneficial mutation to SA58 results in a nearly 7-fold improvement (**Figure 3c**).

238          Taken together, approximately two-third and one-third of tested single amino acid

239   substitutions to Ly-1404 and SA58, respectively, were beneficial for neutralization of either the

240   original strain or BQ.1.1. These results reinforce that, despite all being predicted to have the

11

241 same backbone fold, inverse folding variants feature functional diversity and can be used for

242 distinct notions of protein fitness. Interestingly, for both antibodies, the most beneficial mutation,

243 is not shared by the each of the strains tested (**Supplementary Figure 4**).

244       A common challenge in directed evolution is contending with the combinatorial

245 explosion of possible sequences which emerges from trying to combine a set of individually

246 beneficial mutations. In the second round of evolution, we simply use the inverse folding model

247 again to acquire up to five top-scoring unique combinations of mutations to each antibody chain

248 (**Methods**). Notably, across both evolutionary trajectories, all 15 antibody designs with multiple

249 mutations have $IC_{50}$ values better than wild-type, with many designs showing synergistic effects

250 upon combination. For example, just a single amino acid mutation in each of the two chains of

251 SA58 leads to over an 11-fold improvement (**Figure 3c,d**). Similarly, the most potent evolved

252 design of Ly-1404 is a combination of seven of the eight beneficial single amino acid

253 substitution to the VH and improves neutralization 26-fold (**Figure 3d**). Critically, these

254 improvements to neutralizing potency against BQ.1.1 do not sacrifice potency against the

255 original strains. We found that the top SA58 design against BQ.1.1 after the second round of

256 evolution also improves BA.1 neutralization nearly 3-fold (**Supplementary Data 1**).

257

258 *Additional characterization of evolved antibodies*

259       To further characterize the basis for enhanced neutralization of SARS-CoV-2 VOC

260 BQ.1.1, we tested the binding affinity of all variant antibodies to RBD as bivalent IgG using

261 biolayer interferometry (BLI) to obtain the apparent dissociation constant ($K_{D,app}$). For Ly-1404,

262 all 23 variants with improved neutralization also have improved binding affinity up to ~27-fold.

263 Interestingly, we found four additional inverse folding-recommended mutations, which were

12

264    neutral or deleterious to neutralization, also improved binding affinity. Across all variants there

265    is a Spearman correlation of 0.47 between fold-change in $IC_{50}$ and fold-change in $K_{D,app}$ (**Figure**

266    **4a**).

267        We similarly screened the SA58 variants for binding to the RBD of BQ.1.1. However,

268    since the $K_D$ of the wildtype antibody as IgG was already sub-picomolar, further improvements to

269    binding were below the limit of quantitation and indistinguishable using this measure. Given this

270    strong binding affinity of wildtype SA58 to BQ.1.1 RBD, we also screened this same set of

271    variants against emerging VOC XBB.1.5 and observe improvements in $K_{D,app}$ up to 37-fold

272    (**Figure 4c**).

273        By testing several top affinity-matured designs in a polyspecificity assay, we also

274    confirmed that improvements in binding are not mediated by generalized enhancements of non-

275    specific interactions (**Supplementary Figure 5a**). In this assay, we observed no substantial

276    changes in off-target binding of the evolved antibodies to membrane soluble proteins,

277    particularly within a therapeutically viable range (as defined by controls of clinically approved

278    antibodies with recorded high and low polyspecificity). Furthermore, we found no correlation

279    between fold-change in polyspecificity and affinity fold-change (**Supplementary Figure 5b**).

280

281    *Analysis of evolutionary exploration*

282        Confronted by the large number of possible mutations, traditional experimental-based

283    methods for antibody affinity maturation often restrict the mutational search space to only a few

284    regions of the antibody. Specifically, binding optimization efforts are typically focused within

285    the complementarity determining regions (CDR), which are hotspots for natural somatic

286    hypermutation. However, using our unbiased approach to consider all regions of the variable

13

287    domain allows for many discoveries that may be less intuitive to a rational designer. For

288    example, the most beneficial substitutions to Ly-1404, VH F24Y and VH V90S, are located

289    within framework regions and positioned distally from the binding interface (**Supplementary**

290    **Figure 6, Supplementary Table 2**). Interestingly, they both improve neutralization of BQ.1.1

291    by over 3-fold and are not deleterious to Wuhan neutralization. In other cases, inverse folding

292    also successfully predicts beneficial substitutions using residues rarely observed among human

293    antibody sequences. Substitution VL N95V in SA58, which improves neutralization

294    approximately 7-fold against BQ.1.1, is mediated by the incorporation of a valine observed in

295    only 0.7% of human antibody sequences at that position and enhances antibody-antigen contact.

296    While inverse folding is capable of successfully making novel predictions, in some instances it

297    also does suggest reverting residues to ones frequently selected for in natural somatic

298    hypermutation. Mutation VL F51Y in Ly-1404 changes a phenylalanine observed in just 5% of

299    sequences to a tyrosine observed in 86% of sequences. However, this variant results in no change

300    to Wuhan neutralization. Overall, these results highlight the novelty and value in augmenting a

301    language model with structural information to evolve antibodies and proteins complexes.

302

303    **Discussion**

304          The discovery of mutations that improve protein function is inherently challenging due to

305    the large sequence search space and complex rules that govern the relationship between sequence

306    and function, such as stability or environmental selection pressures. We show that a general

307    inverse folding protein language model informed with the sequence and backbone structural

308    coordinates of a protein can considerably improve directed evolution efforts by serving as an

309    improved prior compared to sequence-only deep learning methods. Importantly, we highlight

14

310    that inverse folding can interrogate protein fitness landscapes indirectly, without needing to

311    explicitly model individual functional tasks or properties, making it broadly applicable to

312    proteins across diverse settings ranging from enzyme catalysis to antibiotic and chemotherapy

313    resistance (**Figure 1d**). We also demonstrate inverse folding generalizes to multimeric proteins,

314    despite being trained only on single-chain proteins, through its ability to implicitly learn features

315    of binding. This result is particularly remarkable considering inverse folding has no access to

316    amino acid side chain atoms, coordinates, or bond information.

317        Equipped with these capabilities, we use inverse folding to evolve clinical therapeutic

318    antibodies and identify several mutations which act synergistically to improve antibody potency

319    and resilience against emerging variants of concern. In the context of pandemics and emergency-

320    use situations, where monoclonal antibody therapies are limited in supply and vulnerable to

321    resistance from viral evolution, the ability to rapidly make improvements in potency with a

322    general method could have major clinical and economic implications.

323        In comparison to fourteen other promising machine learning-guided protein design

324    methods[8,16–28], we find that inverse folding has the strongest performance to date, even without

325    requiring any assay-labeled fitness data to use as training data for task-specific model supervision

326    (**Figure 5**, **Supplementary Data 5**). By eliminating the reliance on any initial data collection,

327    inverse folding has the potential to accelerate entire evolutionary campaigns.

328        Computational methods like the one we propose have the opportunity to democratize

329    protein engineering efforts. Not only is our approach more efficient than conventional resource-

330    intensive techniques that experimentally test the effects of all single-residue changes on

331    biochemical functions like binding affinity, but consequently it enables directed evolution based

332    on properties that are not easily measured at scale or that are incompatible with high-throughput

15

333    screening. Overcoming these limitations, we anticipate our structure-based paradigm will be

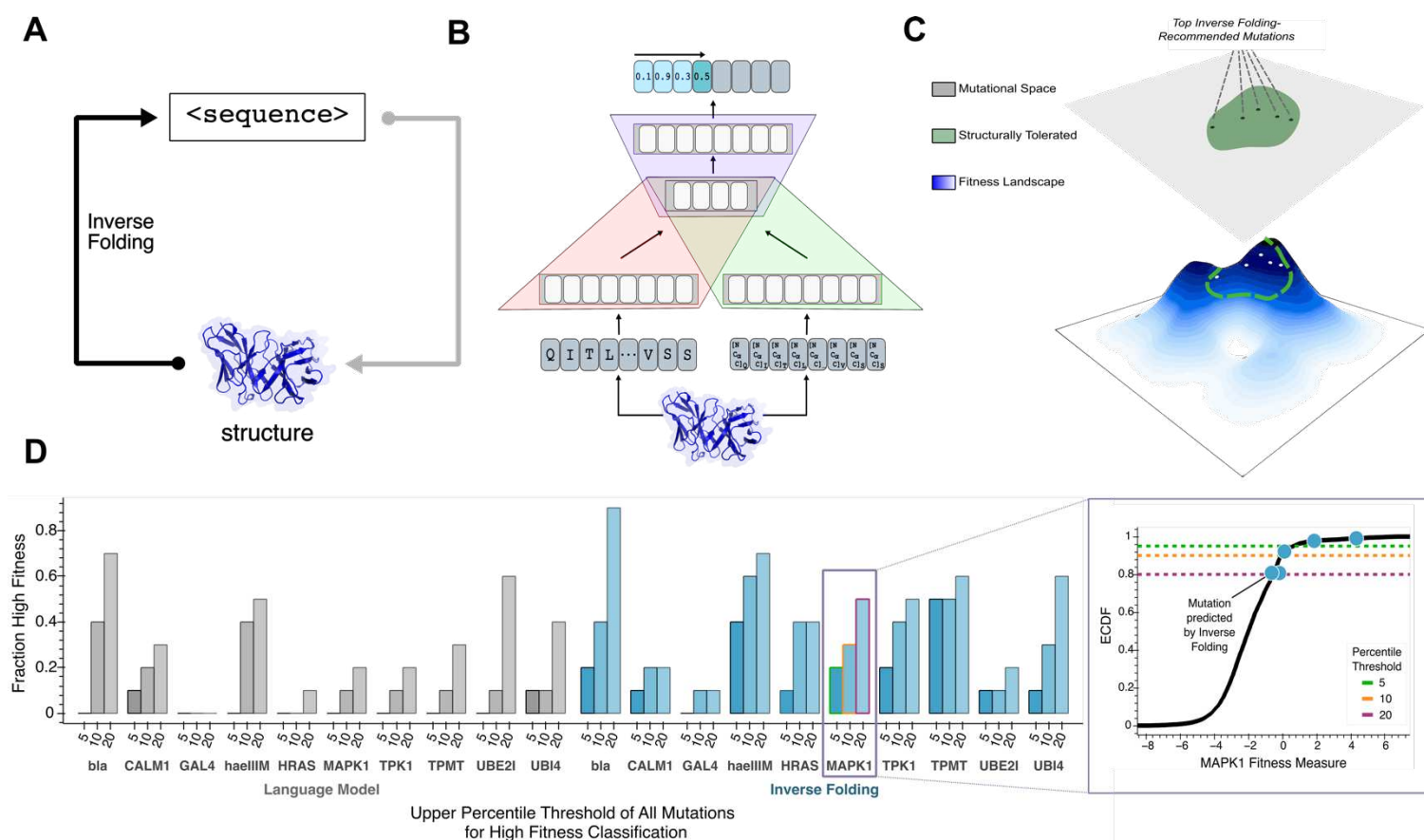334    useful for evolving proteins across many domains.

335

336

**Figure 1: Guiding evolution of diverse proteins via inverse folding**

**(A)** The inverse folding problem refers to the prediction of a protein's native amino acid

sequence, given its three-dimensional backbone structure, which is conceptually analogous to the

opposite problem solved by structure prediction tools like AlphaFold[12]. **(B)** A hybrid

autoregressive model[11] integrates amino acid values and backbone structural information to

evaluate the joint likelihood over all positions in a sequence. Amino acids from the protein

sequence are tokenized (red), combined with geometric features extracted from a structural

encoder (green), and modeled with an encoder-decoder transformer (purple). Sequences assigned

high likelihoods by the model represent high confidence in folding into the input backbone

structure. **(C)** Our structure-guided framework for protein design indirectly explores the

underlying fitness landscape, without modeling a specific definition of fitness or requiring any

348    task-specific training data, by constraining the search space to regions where the backbone fold

349    preserved. **(D)** High fitness sensitivity analysis reveals that multimodal input improves language

350    model performance compared to sequence-only input across 10 proteins from diverse protein

351    families (left). 'Fraction High fitness' is the fraction of the top ten single amino acid substitutions

352    recommended by each model that are ranked in the top indicated percentile of all experimentally

353    screened variants. A representative plot (right) demonstrates this metric for assessing enrichment

354    of high-fitness MAPK1 mutations, with successfully predicted mutations highlighted (blue) on

355    the empirical cumulative density function (ECDF) of the experimental data (black). The three

356    different thresholds, as defined by percentiles, are also shown as dashed lines. Inverse folding

357    predictions are more enriched, on average, for high fitness variants across various tested

358    thresholds for high fitness classification. bla, Beta-lactamase TEM; CALM1, Calmodulin-1;

359    haeIIIM, Type II methyltransferase M.HaeIII; HRAS, GTPase HRas; MAPK1, Mitogen-

360    activated protein kinase; TMPT, Thiopurine S-methyltransferase; TPK1, Thiamin

361    pyrophosphokinase 1; UBI4, Polyubiquitin; UBE2I, SUMO-conjugating enzyme UBC9
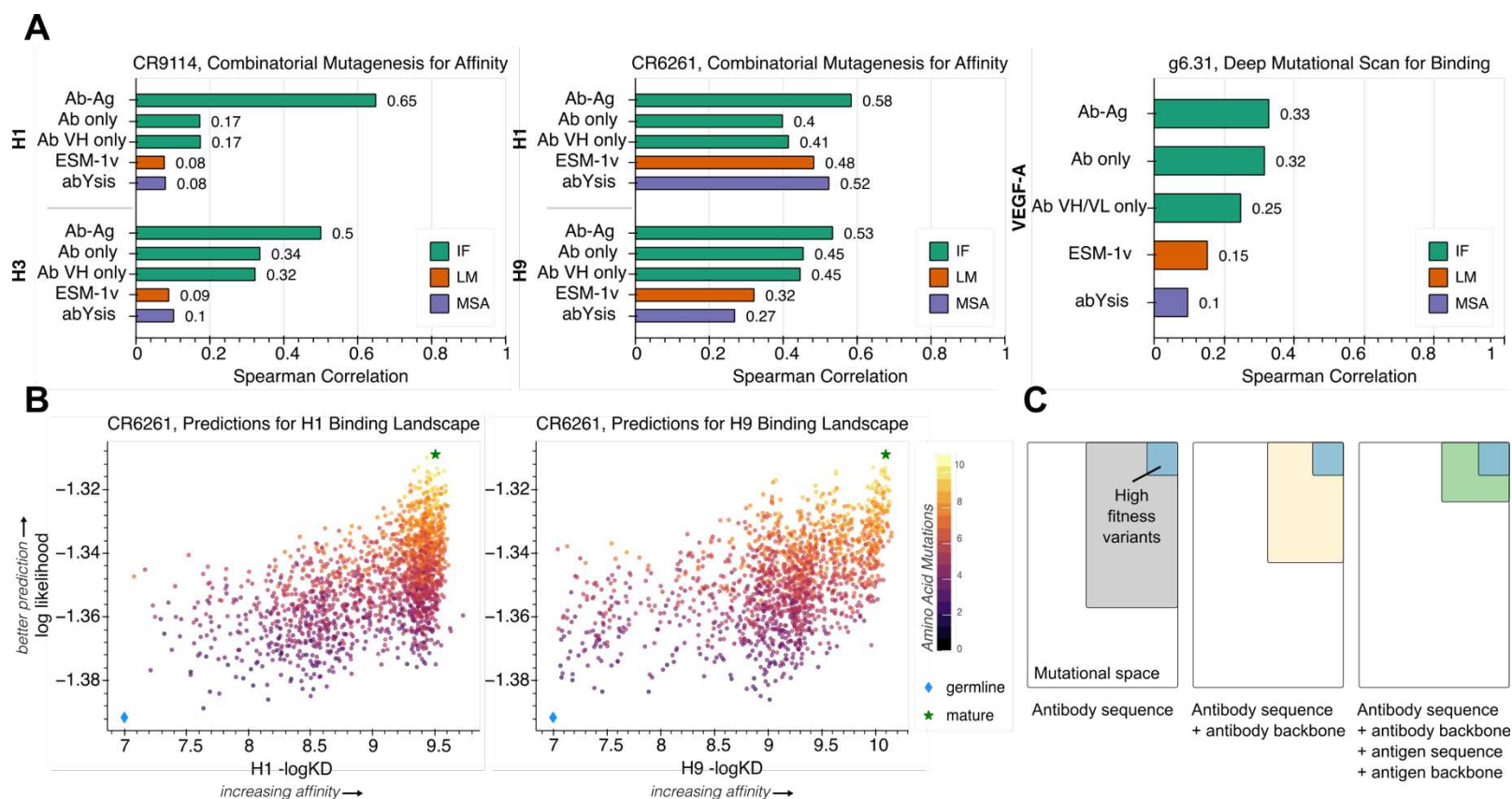
362

**A**



**B**



363

**Figure 2: Inverse folding of antibody-antigen complexes resolves mutational landscapes by**

**implicitly learning features of binding and protein epistasis**

**(A)** Spearman correlation using inverse folding as well as sequence-based modeling approaches

ESM-v[49] and abYsis[55] reported for three antibodies screened with corresponding influenza A HA

subtypes H1, H3, and H9. Bars are colored by the type of model used: IF, Inverse Folding

(green); LM, Language Model (orange); and MSA, Multiple Sequence Alignment (purple).

Inverse folding was evaluated in three different settings: i) providing the entire antibody variable

region and antigen complex (Ab-Ag) ii) providing only the antibody variable region (Ab only),

and iii) providing only the single antibody variable region of the chain responsible for binding or

being mutated (Ab VH only or Ab VH/VL only). Inverse folding implicitly learns features of

binding and protein epistasis. For example, when scoring combinatorial mutations to CR9114

19

375     against H1, we find that the model has much higher performance (Spearman $\rho = 0.65$ for H1, 0.5

376     for H3) than a masked language model ESM-1v (Spearman $\rho = 0.08$ for H1, 0.09 for H3) and a

377     site-independent, alignment-based model abYsis (Spearman $\rho = 0.08$ for H1, 0.1 for H3). This

378     performance improvement is also consistent across the other combinatorial landscapes tested. **(B)**

379     Scatter plots showing inverse folding predictions against experimentally determined dissociation

380     constants of CR6261 against HA-H1(left) and HA-H9 (right). The germline and mature

381     sequences are highlighted on all plots as indicated in the legend. For visualization, all scatter

382     plots omit points on the lower limit of quantitation. Further analysis of assay limit on predictive

383     performance is shown in **Supplementary Figure 2**. **(C)** Conceptual schematic representation of

384     protein language performance improvements with improved priors. Providing sequence and

385     structural information of both the antibody and antigen enables inverse folding to most

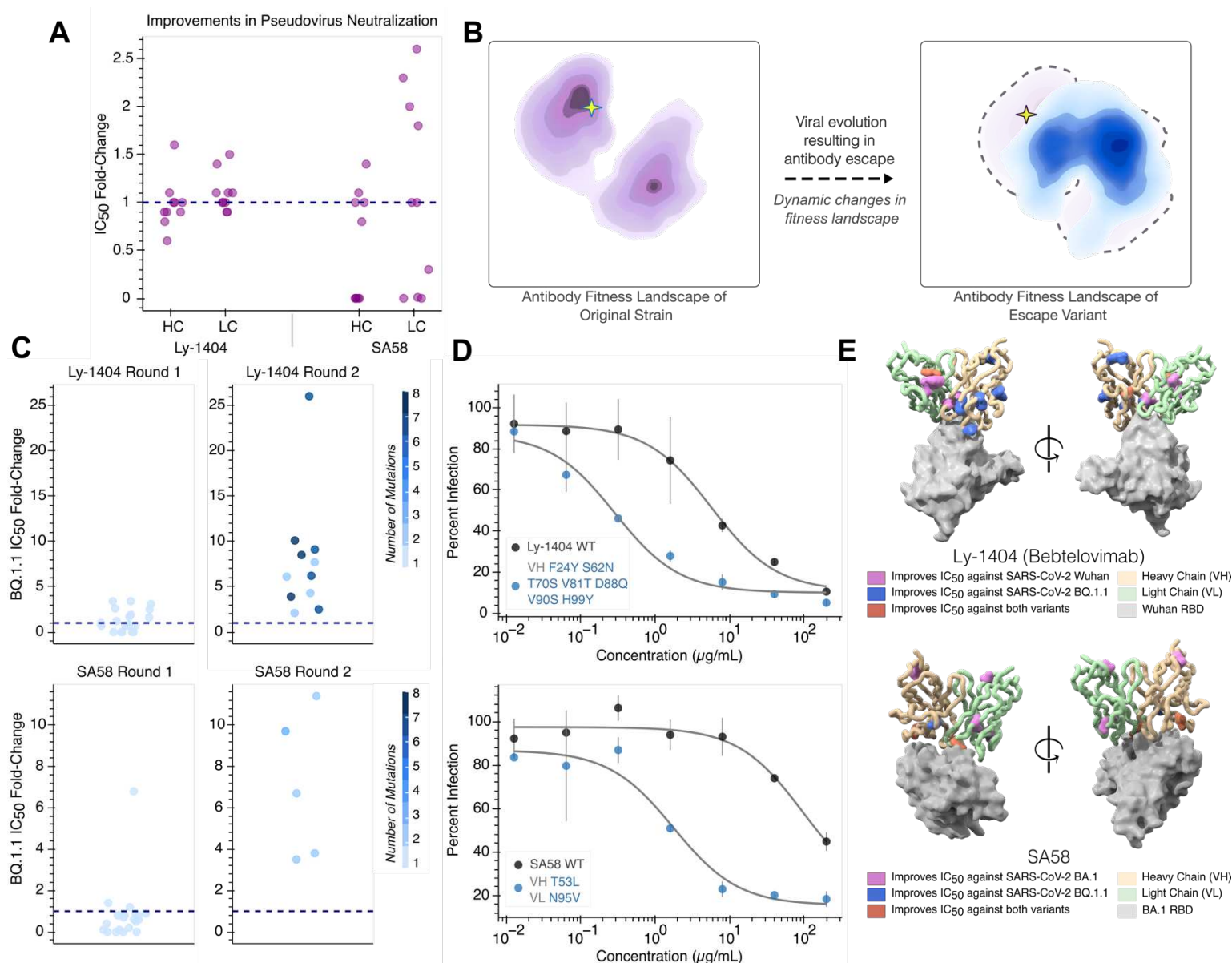386     efficiently identify complex destabilizing mutations and enrich for high fitness antibody variants.

387

**Figure 3. Inverse folding-guided evolution of antibodies improves neutralization potency and resilience**

(**A**) Each point represents the fold-change in IC50 of pseudovirus neutralization for antibody variants with single amino acid mutations. Antibodies are tested against the viral strain represented in the input structure (Ly1404- Wuhan, SA58-BA.1 Omicron). A dashed line is shown at fold-change of 1 corresponding to no change. 35% of Ly-1404 variants and 30% of SA58 variants improved antibody potency (defined as 1.1-fold or higher improvement in IC50

395    compared to wild-type). Among this subset of beneficial mutations, we identify single amino

396    acid mutations that provide a 1.6-fold improvement in Ly-1404 IC50 and a 2.6-fold

397    improvement in SA58 IC50. **(B)** Conceptual representation of viral evolution. Selection for

398    immune evasion drives antibody escape, which fundamentally represents a dynamic change in

399    the underlying fitness landscape for the antibody. This antigenic drift displaces a potent antibody

400    from a peak on the previous fitness landscape (left) to a new starting point at lower activity

401    (right). **(C)** Strip plots visualizing antibody evolution across two rounds. Each point shows the

402    corresponding fold-change in $IC_{50}$ of pseudovirus neutralization for a designed variant and is

403    colored according to the number of mutations it has (1-8). Consistent with preserving backbone

404    fold, all 55 designed variants across both antibody evolutionary campaigns could be expressed.

405    All round 1 variants are only composed of only single amino acid changes while beneficial

406    mutations are combined in round 2. All round 2 variants have improved neutralization activity

407    compared to their respective wild-type antibody (dotted line). **(D)** Pseudovirus neutralization

408    curves are shown for the most potent evolved antibody variant, consisting of mutations annotated

409    to the left.  The top Ly-1404 variant, bearing seven amino acid substitutions in VH, achieves a

410    26-fold improvement in neutralization against BQ.1.1 (top). The top SA58 variant, bearing single

411    amino acid mutations in both VH and VL, achieves an 11-fold improvement in neutralization

412    against BQ.1.1 (bottom). **(E)** Residues at which mutations improve neutralization against either

413    the structure-encoded strain, BQ.1.1, or both viral strains are highlighted with spheres for

414    antibodies Ly-1404 (PDB 7MMO) and SA58 (PDB 7Y0W). Notably, beneficial mutations are

415    identified both within the binding interface as well distal to the antigen. Neutralization enhancing

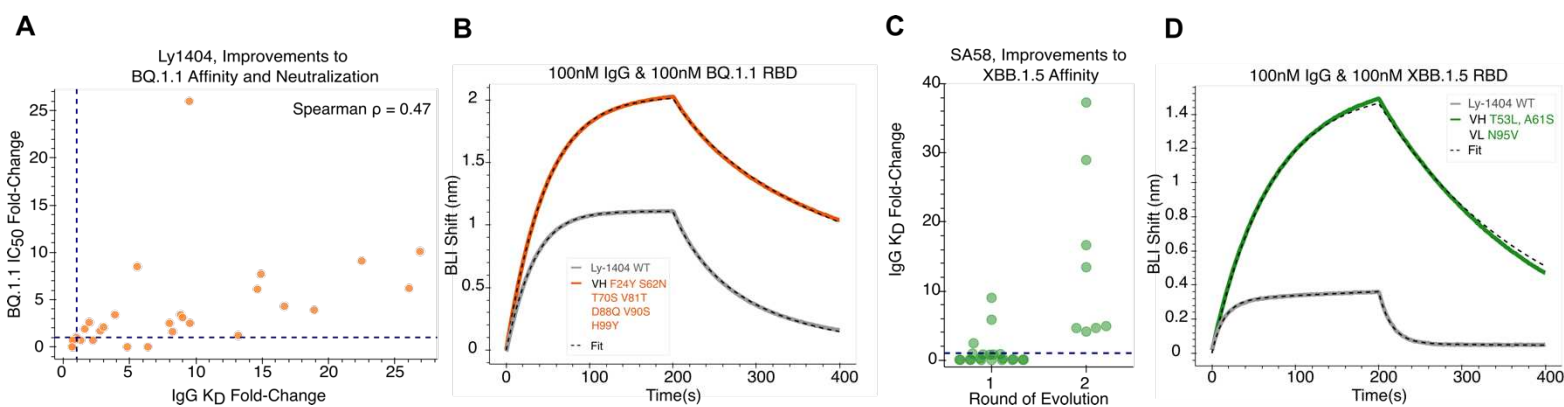416    mutations are labeled in **Supplementary Figure 6**.

417

22

**Figure 4: Antibodies evolved for high potency also exhibit improved affinity**

**(A)** Ly-1404 antibody variants show a Spearman correlation of 0.47 between apparent affinity fold-change and potency fold-change. Improved affinity is observed to be necessary but not sufficient for improved neutralization activity. Four variants exhibit improved affinity but do not enhance neutralization. All variants with improved neutralization also display improved affinity. The top inverse folding Ly-1404 design with a 27-fold improvement in neutralization has a 9.5-fold improvement in affinity to BQ.1.1 RBD, as measured using BLI. **(C)** SA58 antibodies evolved for improved potency against BQ.1.1 also exhibit improved affinity against VOC XBB.1.5, up to 37-fold. **(B, D)** Representative traces of BLI binding assays for Ly-1404 and SA58 variants with improved affinity.
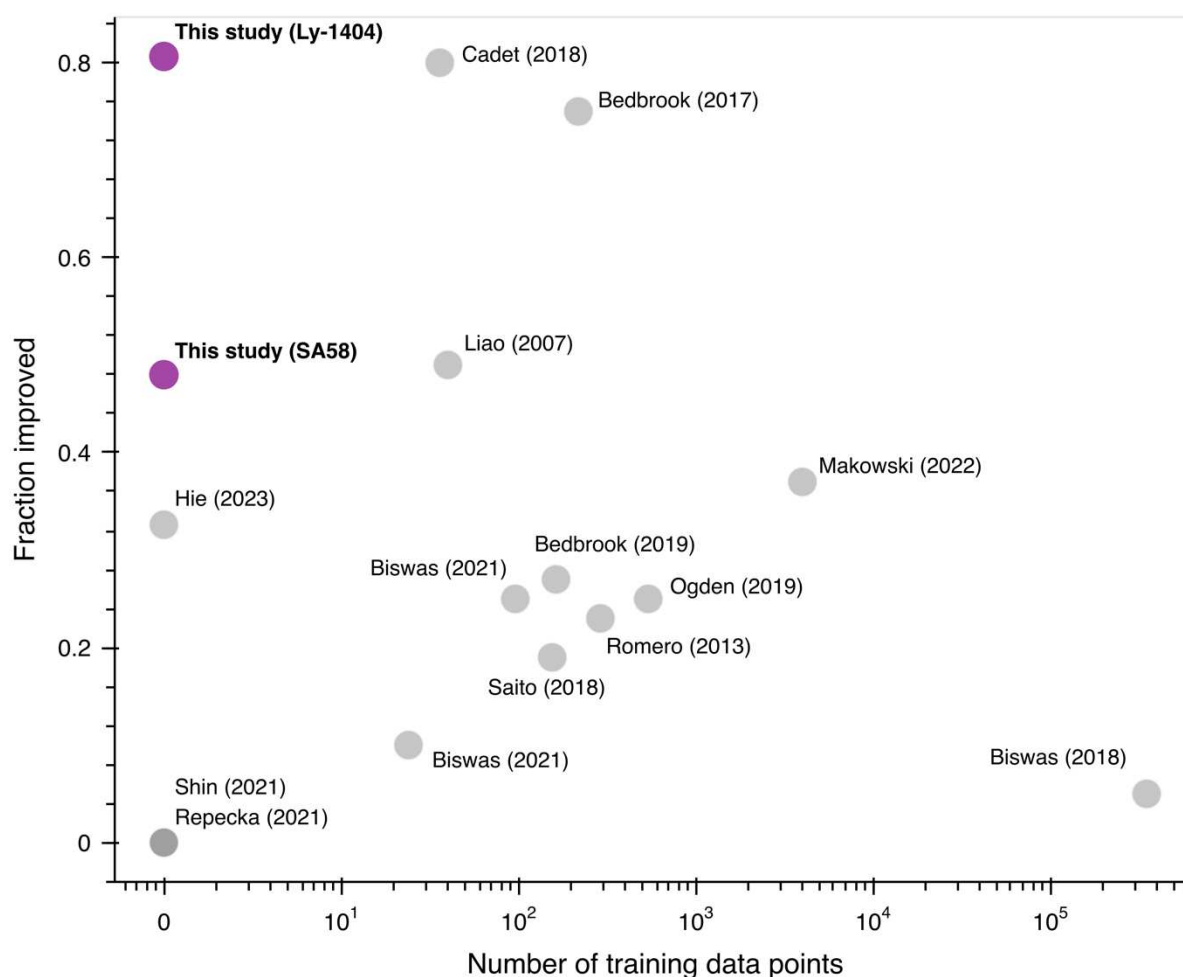
**Figure 5: Comparison to other machine learning-guided directed evolution methods**

'Fraction improved' refers to the hit rate of variants tested which are improved relative to a wildtype protein used as a starting point for directed evolution or a reference protein used as a design template. Higher hit rates indicate more efficient experimental exploration. Inverse folding achieves the highest hit rate with the lowest number of assay-labeled training data points to-date[8,16–28].

439 **Methods**

440 *Inverse folding model description and scoring of sequences*

441     As input to the inverse folding model, we provide a protein structure $\mathbf{Y} \in \mathbb{R}^{N \times 3 \times 3}$, where

442 $N$ is the number of amino acids, and each amino acid is featurized by the three-dimensional

443 physical coordinates of all three atoms in the protein backbone: the α-carbon, β-carbon, and

444 nitrogen atoms in the protein backbone (hence the dimensionality $N \times 3 \times 3$). The inverse

445 folding model learns the probability distribution $p$ of a protein sequence $\mathbf{x} = (x_1, \dots, x_N) \in \mathcal{X}^N$

446 (where $\mathcal{X}$ is the alphabet of amino acids) given a structure $\mathbf{Y}$ via the chain rule of probability

447     $$p(\mathbf{x}|\mathbf{Y}) = p(x_1|\mathbf{Y})p(x_2|x_1, \mathbf{Y}) \dots p(x_N|x_1, \dots, x_{N-1}, \mathbf{Y}).$$

448     The probability distribution at each position is defined over $\mathcal{X}$, such that it is a 20-

449 dimensional vector with all constituent entries summing to 1.

450     Thus, for a given sequence $\hat{\mathbf{x}} = (\hat{x}_1, \dots, \hat{x}_N)$ and its corresponding given structure $\hat{\mathbf{Y}}$, we

451 can score the probability of $\hat{\mathbf{x}}$ folding into $\mathbf{Y}$ under the inverse folding model by computing the

452 value of $p(\mathbf{x} = \hat{\mathbf{x}}|\mathbf{Y})$, which we can do autoregressively as

453     $$p(\mathbf{x} = \hat{\mathbf{x}}|\hat{\mathbf{Y}}) = p(x_1 = \hat{x}_1|\hat{\mathbf{Y}}) \dots p(x_N = \hat{x}_N|\hat{x}_1, \dots, \hat{x}_{N-1}, \hat{\mathbf{Y}}).$$

454     This is evaluated output is a likelihood between 0 and 1, inclusive. The computed score

455 $p(\mathbf{x} = \hat{\mathbf{x}}|\hat{\mathbf{Y}})$ is used as prediction for "fitness" (e.g., binding affinity or enzymatic activity).

456 Importantly, the inverse folding model does not have any explicit access to "fitness" during

457 either training or evaluation, which we refer to as "zero shot" fitness prediction.

458     We use the inverse folding model checkpoint of ESM-IF1 GVP-Transformer as of April

459 10, 2022[11].

460

461 *Diverse proteins benchmarking experiment with scanning mutagenesis data*

25

462    We analyzed the effectiveness of using the inverse folding language model, ESM-IF1

463    model to identify high fitness variants from protein mutational scans as a proxy for the ability to

464    guide evolution without explicitly modeling a protein's function. We also compared its

465    performance to ESM-1v, a sequence-only general protein language model. To do so, we used all

466    deep mutational scanning (DMS) datasets from the benchmarking study by Livesey and Marsh[29]

467    profiling over 100 variants and reported to have 90% or higher coverage of DMS results across

468    the corresponding curated PDB structure (**Supplementary Table 1**). From this set of 12

469    proteins, Cas9 was excluded because its sequence length was larger than the maximum allowable

470    length of 1024 amino acids by ESM-1v and ccdB was excluded because the experimental values

471    were discretized within a small range. For each of the 10 mutagenesis datasets, all the sequence

472    likelihood of all variants with coverage in the structure were determined using inverse folding.

473    For ESM-1v, the average masked marginals likelihood score across all five models in the ESM-

474    1v group was used. The experimental data distribution was binarized for high-fitness

475    classification using a percentile-based threshold. The enrichment of high fitness variants was

476    then determined by using the metric of fraction high fitness as defined by the fraction of the top

477    10 model-predicted variants with experimental values above the high fitness threshold. The

478    analysis was performed at three different percentile thresholds, top 5th percentile (95th percentile),

479    top 10th percentile (90th percentile), and top 20th percentile (80th percentile), to determine

480    sensitivity of the result based on the stringency of the selected cutoff parameter.

481

482    *Benchmarking of antibody mutagenesis*

483    We use three antibody mutagenesis datasets[50,51] to benchmark the performance of

484    modeling variant effects on antibody binding using inverse folding against two sequence-only

485    methods, ESM-1v[49] and abYsis[55]. Variant sequences were scored using the inverse folding model

486    with three different forms of structure input: i) variable region of mutated antibody chain only ii)

487    variable regions of both antibody chains iii) variable regions of both antibody chains in complex

488    with antigen. The autoregressive scoring of sequences with inverse folding enables evaluation of

489    sequences with multiple mutations. The Spearman correlation was determined between the log

490    likelihood scores across all sequences and corresponding reported experimental binding

491    measurements: -log($K_D$) for CR9114 and CR6261; log(binding enrichment) g6.31. The following

492    structures were used for input backbone coordinates of the VH, VL, and antigen: PDB 4FQI[52],

493    CR9114-H5; PDB 3GBN[53], CR6261-H1; PDB 2FJG, g6.31-VEGF.

494            ESM-1v and abYsis were scored using the variant sequence of the antibody variable

495    region. For variants with multiple mutations, the average effect of all mutant amino acids in the

496    sequence was considered, namely

497    $$p(\mathbf{x}) = \frac{1}{|\mathcal{M}|} \sum_{i \in M} [\log\ p(\mathbf{x}_i = \mathbf{x}_i^{\mathrm{mt}}) - \log p(\mathbf{x}_i = \mathbf{x}_i^{\mathrm{wt}})]$$

498    where $\mathcal{M}$ is defined as the set of all mutations in the input sequence $\mathbf{x}$. For abYsis, individual

499    mutation likelihoods were determined using the frequency of amino acids at each position based

500    on multiple sequence alignment provided by the webtool version 3.4.1

501    (http://www.abysis.org/abysis/index.html). We aligned VH and VL protein sequences using the

502    default settings provided in the 'Annotate' tool, with the database of 'Homo sapiens' sequences

503    as of April 1, 2023.

504

505    *Acquisition of antibody amino acid substitutions using inverse folding*

506            We select amino acid substitutions recommended by the inverse folding model to test in

507    our directed evolution campaigns for Ly-1404 and SA58. For a given wild-type antibody

27

508    variable region sequence, $\mathbf{x} = (x_1, \ldots, x_N) \in \mathcal{X}^N$, where $\mathcal{X}$ is the set of amino acids and $N$ is the

509    sequence length, we score all possible single amino acid substitutions against a corresponding

510    structure of the variable regions of both antibody chains in complex with the RBD of SARS-

511    CoV-2 Spike protein, $\hat{\mathbf{Y}}$ by computing $p(\mathbf{x} = \hat{\mathbf{x}}|\hat{\mathbf{Y}})$. Protein structures used are reported in

512    Supplementary Table 1. We then select the set of top ten predicted single amino acid

513    substitutions at unique residues in each antibody variable region to test in the first round of

514    evolution.

515        After testing individual amino acid mutations in a pseudovirus neutralization screen, in

516    Round 2, beneficial mutations (defined as $IC_{50}$ fold-change > 1.1) were combined to assess the

517    combinatorial effects and potential for further neutralization improvement. We tested up to four

518    combinations of single amino acid mutations on each chain (two total mutations to the antibody).

519    We also used the inverse folding model to score a library of all possible combinations of the

520    beneficial mutations to an antibody chain (For example, VH Ly-1404 has 8 beneficial mutations

521    resulting in 255 total candidate sequences), and selected the top five scoring designs (or less if

522    there were a fewer number of total possible combinations). Lastly, we tested a maximum of two

523    variants consisting of the best single-chain designs together. In total, 31 variants were tested for

524    Ly-1404 and 25 variants were tested for SA58.

525

526    *Antibody cloning*

527        We cloned the antibody sequences into the CMV/R plasmid backbone for expression

528    under a CMV promoter. The heavy chain or light chain sequence was cloned between the CMV

529    promoter and the bGH poly(A) signal sequence of the CMV/R plasmid to facilitate improved

530    protein expression. Variable regions were cloned into the human IgG1 backbone; Ly-1404

28

531  variants were cloned with a lambda light chain, whereas variants of SA58 were cloned with a

532  kappa light chain. The vector for both heavy and light chain sequences also contained the

533  HVM06_Mouse (UniProt: P01750) Ig heavy chain V region 102 signal peptide

534  (MGWSCIILFLVATATGVHS) to allow for protein secretion and purification from the

535  supernatant. VH and VL segments were ordered as gene blocks from Integrated DNA

536  Technologies and were cloned into linearized CMV/R backbones with 5× In-Fusion HD Enzyme

537  Premix (Takara Bio).

538

539  *Antigen cloning*

540      RBD sequences were cloned into a pADD2 vector between the rBeta-globin intron and β-

541  globin poly(A). All RBD constructs contain an AviTag and 6×His tag. RBD sequences were

542  based off wild-type Wuhan-Hu-1 (GenBank: BCN86353.1), Omicron BA.1

543  (GenBank: UFO69279.1), BQ.1.1 (GenBank: OP412163.1 ), XBB.1.5 (GenBank: OP790748.1 ).

544

545  *DNA preparation*

546      Plasmids were transformed into Stellar competent cells (Takara Bio), and transformed

547  cells were plated and grown at 37 °C overnight. Colonies were mini-prepped per the

548  manufacturer's recommendations (GeneJET, K0502, Thermo Fisher Scientific) and sequence

549  confirmed (Sequetech) and then maxi-prepped per the manufacturer's protocols (ZymoPure II

550  Plasmid Maxiprep Kit, Zymo Research). Plasmids were sterile filtered using a 0.22-μm syringe

551  filter and stored at 4 °C.

552

553  *Protein expression*

29

554       All proteins were expressed in Expi293F cells (Thermo Fisher Scientific, A14527).

555    Proteins containing a biotinylation tag (AviTag) were also expressed in the presence of a BirA

556    enzyme, resulting in spontaneous biotinylation during protein expression. Expi293F cells were

557    cultured in media containing 66% FreeStyle/33% Expi media (Thermo Fisher Scientific) and

558    grown in TriForest polycarbonate shaking flasks at 37 °C in 8% carbon dioxide. The day before

559    transfection, cells were pelleted by centrifugation and resuspended to a density of $3 \times 10^6$ cells

560    per milliliter in fresh media. The next day, cells were diluted and transfected at a density of

561    approximately $3–4 \times 10^6$ cells per milliliter. Transfection mixtures were made by adding the

562    following components: maxi-prepped DNA, culture media and FectoPRO (Polyplus) would be

563    added to cells to a ratio of 0.5 μg: 100 μl: 1.3 μl: 900 μl. For example, for a 100-ml transfection,

564    50 μg of DNA would be added to 10 ml of culture media, followed by the addition of 130 μl of

565    FectoPRO. For antibodies, we divided the transfection DNA equally among heavy and light

566    chains; in the previous example, 25 μg of heavy chain DNA and 25 μg of light chain DNA would

567    be added to 10 ml of culture media. After mixing and a 10-min incubation, the example

568    transfection cocktail would be added to 90 ml of cells. The cells were harvested 3–5 days after

569    transfection by spinning the cultures at 10,000$g$ for 10 min. Supernatants were filtered using a

570    0.45-μm filter.

571

572    *Antibody purification*

573       We purified antibodies using a 5-ml MabSelect Sure PRISM column on the ÄKTA pure

574    fast protein liquid chromatography (FPLC) instrument (Cytiva). The ÄKTA system was

575    equilibrated with line A1 in 20 mM 4-(2-hydroxyethyl)-1-piperazineethanesulfonic acid

576    (HEPES) pH 7.4, 150 mM sodium chloride (NaCl), line A2 in 100 mM glycine pH 2.8, line B1

30

577    in 0.5 M sodium hydroxide, Buffer line in 20 mM 4-(2-hydroxyethyl)-1-piperazineethanesulfonic

578    acid (HEPES) pH 7.4, 150 mM sodium chloride (NaCl) and Sample lines in water. The protocol

579    washes the column with A1, followed by loading of the sample in the Sample line until air is

580    detected in the air sensor of the sample pumps, followed by five column volume washes with A1,

581    elution of the sample by flowing of 20 ml of A2 directly into a 50-ml conical containing 2 ml of

582    1 M tris(hydroxymethyl)aminomethane (Tris) pH 8.0, followed by five column volumes of A1,

583    B1 and A1 and then a wash step of the fraction collector with A1. We concentrated the eluted

584    samples using 50-kDa cutoff centrifugal concentrators, followed by buffer exchange using a PD-

585    10 column (Sephadex) that had been pre-equilibrated into 20 mM 4-(2-hydroxyethyl)-1-

586    piperazineethanesulfonic acid (HEPES) pH 7.4, 150 mM sodium chloride (NaCl). Purified

587    antibodies were used directly in experiments or flash-frozen and stored at –20 °C.

588

589    *Antigen purification*

590         All RBD antigens were His-tagged and purified using HisPur Ni-NTA resin (Thermo

591    Fisher Scientific, 88222). Cell supernatants were diluted with 1/3 volume of wash buffer (20 mM

592    imidazole, 20 mM 4-(2-hydroxyethyl)-1-piperazineethanesulfonic acid (HEPES) pH 7.4,

593    150 mM sodium chloride (NaCl), and the Ni-NTA resin was added to diluted cell supernatants.

594    For all antigens, the samples were then incubated at 4 °C while stirring overnight.

595    Resin/supernatant mixtures were added to chromatography columns for gravity flow purification.

596    The resin in the column was washed with wash buffer (20 mM imidazole, 20 mM HEPES pH

597    7.4, 150 mM NaCl), and the proteins were eluted with 250 mM imidazole, 20 mM HEPES pH

598    7.4, 150 mM NaCl. Column elutions were concentrated using centrifugal concentrators at 10-

599    kDa cutoff, followed by size-exclusion chromatography on an ÄKTA pure system (Cytiva).

31

600  ÄKTA pure FPLC with a Superdex 200 Increase (S200) gel filtration column was used for

601  purification. Then, 1 ml of sample was injected using a 2-ml loop and run over the S200, which

602  had been pre-equilibrated in degassed 20 mM HEPES, 150 mM NaCl before use and flash-frozen

603  before storage at −20 °C.

604

605  *BLI binding experiments*

606  All reactions were run on an Octet RED96 at 30 °C, and samples were run in 1× PBS

607  with 0.1% BSA and 0.05% Tween 20 (Octet buffer). IgGs were assessed for binding to

608  biotinylated antigens using streptavidin biosensors (Sartorius/ForteBio). Antigen was loaded at a

609  concentration of 200nM. Tips were then washed and baselined in wells containing only Octet

610  buffer. Samples were then associated in wells containing IgG at 100 nM concentration. A control

611  well with loaded antigen but that was associated in a well containing only 200 μl of Octet buffer

612  was used as a baseline subtraction for data analysis. Association and dissociation binding curves

613  were fit in Octet System Data Analysis Software version 9.0.0.15 using a 1:2 bivalent model for

614  IgGs to determine apparent $K_d$. Fold-change in apparent $K_d$ were determined by computing the

615  ratio of wildtype $K_d$ to variant $K_d$. Averages of $K_d$ fold-change values from at least two

616  independent experiments are reported to two significant figures in **Supplementary Data 2**. To

617  estimate measurement error, we computed the standard deviation for each

618  antibody−antigen $K_d$ pair.

619

620  *Polyspecificity Particle assay*

621  Polyspecificity reagent (PSR) was obtained as described by Xu et al[63]. Soluble membrane

622  proteins were isolated from homogenized and sonicated Expi 293F cells followed by

32

623     biotinylation with Sulfo-NHC-SS-Biotin (Thermo Fisher Scientific, 21331) and stored in PBS at

624     −80 °C. The PolySpecificity Particle (PSP) assay was performed as described in Makowski et

625     al.[64]. Protein A magnetic beads (Invitrogen, 10001D) were washed three times in PBSB (PBS

626     with 1 mg ml$^{-1}$ BSA) and diluted to 54 µg ml$^{-1}$ in PBSB. Then, 30 µl of the solution containing

627     the beads was incubated with 85 µl of antibodies at 15 $\mu$g ml$^{-1}$ overnight at 4 °C with rocking.

628     The coated beads were then washed twice with PBSB using a magnetic plate stand (Invitrogen,

629     12027) and resuspended in PBSB. We then incubated 50 µl of 0.1 mg ml$^{-1}$ PSR with the washed

630     beads at 4 °C with rocking for 20 min. Beads were then washed with PBSB and incubated with

631     0.001× streptavidin-APC (BioLegend, 405207) and 0.001× goat anti-human Fab fragment FITC

632     (Jackson ImmunoResearch, 109-097-003) at 4 °C with rocking for 15 min. Beads were then

633     washed and resuspended with PBSB. Beads were profiled via flow cytometry using a Sony

634     SH800 cell sorter. Data analysis was performed with FlowJo software version 10.9.0 to obtain

635     median fluorescence intensity (MFI) values, which are reported for each antibody across three or

636     more replicate wells. Elotuzumab (Fisher Scientific) and ixekizumab (Fisher Scientific) are also

637     included in each assay as controls.

638

639     *Lentivirus production*

640         We produced SARS-CoV-2 Spike (Wuhan-Hu-1, BA.1, and BQ.1.1 variants)

641     pseudotyped lentiviral particles. Viral transfections were done in HEK293T cells (American

642     Type Culture Collection, CRL-3216) using BioT (BioLand) transfection reagent. Six million

643     cells were seeded in D10 media (DMEM + additives: 10% FBS, L-glutamate, penicillin,

644     streptomycin and 10 mM HEPES) in 10-cm plates one day before transfection. A five-plasmid

645     system was used for viral production, as described in Crawford et al[65]. The Spike vector

33

646    contained the 21-amino-acid truncated form of the SARS-CoV-2 Spike sequence from the

647    Wuhan-Hu-1 strain of SARS-CoV-2 (GenBank: BCN86353.1), BA.1 variant of concern

648    (GenBank: OL672836.1), or BQ.1.1 variant of concern (GenBank: OP412163.1. The other viral

649    plasmids, used as previously described[65], are pHAGE-Luc2-IRS-ZsGreen (NR-52516), HDM-

650    Hgpm2 (NR-52517), pRC-CMV-Rev1b (NR-52519) and HDM-tat1b (NR-52518). These

651    plasmids were added to D10 medium in the following ratios: 10 μg pHAGE-Luc2-IRS-ZsGreen,

652    3.4 μg FL Spike, 2.2 μg HDM-Hgpm2, 2.2 μg HDM-Tat1b and 2.2 μg pRC-CMV-Rev1b in a

653    final volume of 1,000 μl.

654        After adding plasmids to medium, we added 30 μl of BioT to form transfection

655    complexes. Transfection reactions were incubated for 10 min at room temperature, and then 9 ml

656    of medium was added slowly. The resultant 10 ml was added to plated HEK cells from which the

657    medium had been removed. Culture medium was removed 24 h after transfection and replaced

658    with fresh D10 medium. Viral supernatants were harvested 72 h after transfection by spinning at

659    300$g$ for 5 min, followed by filtering through a 0.45-μm filter. Viral stocks were aliquoted and

660    stored at −80 °C.

661

662    *Pseudovirus neutralization*

663        The target cells used for infection in SARS-CoV-2 pseudovirus neutralization assays are

664    from a HeLa cell line stably overexpressing human angiotensin-converting enzyme 2 (ACE2) as

665    well as the protease known to process SARS-CoV-2: transmembrane serine protease 2

666    (TMPRSS2). Production of this cell line is described in detail by Rogers et al[66]. with the addition

667    of stable TMPRSS2 incorporation. ACE2/TMPRSS2/HeLa cells were plated 1 day before

34

668  infection at 8,000 cells per well. Ninety-six-well, white-walled, white-bottom plates were used

669  for neutralization assays (Thermo Fisher Scientific).

670  On the day of the assay, purified IgGs in 1× PBS were made into D10 medium (DMEM +

671  additives: 10% FBS, L-glutamate, penicillin, streptomycin and 10 mM HEPES). A virus mixture

672  was made containing the virus of interest (for example, SARS-CoV-2) and D10 media. Virus

673  dilutions into media were selected such that a suitable signal would be obtained in the virus-only

674  wells. A suitable signal was selected such that the virus-only wells would achieve a

675  luminescence of at least >1,000,000 relative light units (RLU). Then, 60 μl of this virus mixture

676  was added to each of the antibody dilutions to make a final volume of 120 μl in each well. Virus-

677  only wells were made, which contained 60 μl of D10 and 60 μl of virus mixture. Cells-only wells

678  were made, which contained 120 μl of D10 media.

679  The antibody/virus mixture was left to incubate for 1 h at 37 °C. After incubation, the

680  medium was removed from the cells on the plates made one day prior. This was replaced with

681  100 μl of antibody/virus dilutions and incubated at 37 °C for approximately 48 h. Infectivity

682  readout was performed by measuring luciferase levels. Medium was removed from all wells, and

683  cells were lysed by the addition of 100 μl of BriteLite assay readout solution (PerkinElmer) into

684  each well. Luminescence values were measured using an Infinite 200 PRO Microplate Reader

685  (Tecan) using i-control version 2.0 software (Tecan) after shaking for 30 sec. Each plate was

686  normalized by averaging the cells-only (0% infection) and virus-only (100% infection)

687  wells. Neutralization titer was defined as the sample dilution at which the RLU was decreased by

688  50% as compared with the RLU of virus-only control wells after subtraction of background

689  RLUs in wells containing cells only. Normalized values were fitted with a three-parameter

690  nonlinear regression inhibitor curve in GraphPad Prism 9.1.0 to determine the half-maximal

35

691    inhibitory concentration ($IC_{50}$) and are reported in **Supplementary Data 1**. Neutralization assays

692    were performed in biological duplicates with technical duplicates.


693    *Computing frequency of changes to antibody protein sequences*

694    We computed the frequency of residues involved in affinity-enhancing substitutions using

695    the abYsis webtool, which also computes the frequency of amino acids at each position based on

696    a multiple sequence alignment. We aligned VH and VL protein sequences using the default

697    settings provided in the 'Annotate' tool, using the database of 'All' sequences as of April 1,

698    2023. We also used the Kabat region definition provided by abYsis webtool version 3.4.1 to

699    annotate the framework regions and CDRs within the VH and VL sequences which are reported

700    in **Supplementary Table 2**.

701

702    *Comparing efficiency of machine learning-guided directed evolution methods*

703    To compare inverse folding against other machine learning methods for protein

704    evolution, we compared the fraction of variants tested in the protein engineering campaign to the

705    number of assay-labeled training data points used to inform the predictions. Data was sourced

706    from Biswas et al.[17] and made contemporaneous by the addition of recently published studies as

707    indicated in **Supplementary Data 5**. The fraction improved, or hit rate, refers to experimentally

708    tested predictions which have improved functional activity relative to either a wildtype protein

709    that is used as a starting point for directed evolution or the protein used as a reference template

710    for design.

711

**Acknowledgments**

**Author contributions**

Conceptualization, methodology, interpretation: V.R.S., B.L.H., P.S.K.; Computational experiments and software development: V.R.S.; Antibody and antigen cloning, expression, and purification: V.R.S., T.U.J.B.; Lentivirus production and pseudovirus neutralization: T.U.J.B; Binding assays: V.R.S.; Writing (original draft): V.R.S with assistance from B.L.H and P.S.K.; Writing (final draft): all authors

**Competing interests**

V.R.S., B.L.H., and P.S.K. are named as inventors on a patent application applied for by Stanford University and the Chan Zuckerberg Biohub entitled "Antibody Compositions and Optimization Methods".

## References

1. Chothia, C. & Lesk, A. M. The relation between the divergence of sequence and structure in proteins. *EMBO J.* **5**, 823–826 (1986).

2. Bloom, J. D., Labthavikul, S. T., Otey, C. R. & Arnold, F. H. Protein stability promotes evolvability. *Proc. Natl. Acad. Sci.* **103**, 5869–5874 (2006).

3. Axe, D. D., Foster, N. W. & Fersht, A. R. A Search for Single Substitutions That Eliminate Enzymatic Function in a Bacterial Ribonuclease. *Biochemistry* **37**, 7157–7166 (1998).

4. Romero, P. A. & Arnold, F. H. Exploring protein fitness landscapes by directed evolution. *Nat. Rev. Mol. Cell Biol.* **10**, 866–876 (2009).

5. Shafikhani, S., Siegel, R. A., Ferrari, E. & Schellenberger, V. Generation of large libraries of random mutants in Bacillus subtilis by PCR-based plasmid multimerization. *BioTechniques* **23**, 304–310 (1997).

6. Guo, H. H., Choe, J. & Loeb, L. A. Protein tolerance to random amino acid change. *Proc. Natl. Acad. Sci.* **101**, 9205–9210 (2004).

7. Rennell, D., Bouvier, S. E., Hardy, L. W. & Poteete, A. R. Systematic mutation of bacteriophage T4 lysozyme. *J. Mol. Biol.* **222**, 67–88 (1991).

8. Ogden, P. J., Kelsic, E. D., Sinai, S. & Church, G. M. Comprehensive AAV capsid fitness landscape reveals a viral gene and enables machine-guided design. *Science* **366**, 1139–1143 (2019).

9. Poelwijk, F. J., Kiviet, D. J., Weinreich, D. M. & Tans, S. J. Empirical fitness landscapes reveal accessible evolutionary paths. *Nature* **445**, 383–386 (2007).

10. Wittmann, B. J., Yue, Y. & Arnold, F. H. Informed training set design enables efficient machine learning-assisted directed protein evolution. *Cell Syst.* **12**, 1026-1045.e7 (2021).

755   11. Hsu, C. *et al*. Learning inverse folding from millions of predicted structures. in *Proceedings*

756       *of the 39th International Conference on Machine Learning* 8946–8970 (PMLR, 2022).

757   12. Jumper, J. *et al*. Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,

758       583–589 (2021).

759   13. Lin, Z. *et al*. Evolutionary-scale prediction of atomic-level protein structure with a language

760       model. *Science* **379**, 1123–1130 (2023).

761   14. Carter, P. J. & Lazar, G. A. Next generation antibody drugs: pursuit of the 'high-hanging

762       fruit'. *Nat. Rev. Drug Discov.* **17**, 197–223 (2018).

763   15. Schroeder, H. W. & Cavacini, L. Structure and function of immunoglobulins. *J. Allergy Clin.*

764       *Immunol.* **125**, S41–S52 (2010).

765   16. Makowski, E. K. *et al*. Co-optimization of therapeutic antibody affinity and specificity using

766       machine learning models that generalize to novel mutational space. *Nat. Commun.* **13**, 3788

767       (2022).

768   17. Biswas, S., Khimulya, G., Alley, E. C., Esvelt, K. M. & Church, G. M. Low-N protein

769       engineering with data-efficient deep learning. *Nat. Methods* **18**, 389–396 (2021).

770   18. Bedbrook, C. N. *et al*. Machine learning-guided channelrhodopsin engineering enables

771       minimally invasive optogenetics. *Nat. Methods* **16**, 1176–1184 (2019).

772   19. Wu, Z., Kan, S. B. J., Lewis, R. D., Wittmann, B. J. & Arnold, F. H. Machine learning-

773       assisted directed protein evolution with combinatorial libraries. *Proc. Natl. Acad. Sci.* **116**,

774       8852–8858 (2019).

775   20. Cadet, F. *et al*. A machine learning approach for reliable prediction of amino acid

776       interactions and its application in the directed evolution of enantioselective enzymes. *Sci.*

777       *Rep.* **8**, 16757 (2018).

778    21. Saito, Y. *et al*. Machine-Learning-Guided Mutagenesis for Directed Evolution of Fluorescent

779        Proteins. *ACS Synth. Biol.* **7**, 2014–2022 (2018).

780    22. Biswas, S. *et al*. Toward machine-guided design of proteins. 337154 Preprint at

781        https://doi.org/10.1101/337154 (2018).

782    23. Romero, P. A., Krause, A. & Arnold, F. H. Navigating the protein fitness landscape with

783        Gaussian processes. *Proc. Natl. Acad. Sci.* **110**, E193–E201 (2013).

784    24. Bedbrook, C. N., Yang, K. K., Rice, A. J., Gradinaru, V. & Arnold, F. H. Machine learning

785        to design integral membrane channelrhodopsins for efficient eukaryotic expression and

786        plasma membrane localization. *PLOS Comput. Biol.* **13**, e1005786 (2017).

787    25. Liao, J. *et al*. Engineering proteinase K using machine learning and synthetic genes. *BMC

788        Biotechnol.* **7**, 16 (2007).

789    26. Repecka, D. *et al*. Expanding functional protein sequence spaces using generative adversarial

790        networks. *Nat. Mach. Intell.* **3**, 324–333 (2021).

791    27. Shin, J.-E. *et al*. Protein design and variant prediction using autoregressive generative

792        models. *Nat. Commun.* **12**, 2403 (2021).

793    28. Hie, B. L. *et al*. Efficient evolution of human antibodies from general protein language

794        models. *Nat. Biotechnol.* 1–9 (2023) doi:10.1038/s41587-023-01763-2.

795    29. Livesey, B. J. & Marsh, J. A. Using deep mutational scanning to benchmark variant effect

796        predictors and identify disease mutations. *Mol. Syst. Biol.* **16**, e9380 (2020).

797    30. Weile, J. *et al*. A framework for exhaustively mapping functional missense variants. *Mol.

798        Syst. Biol.* **13**, 957 (2017).

799    31. Bandaru, P. *et al*. Deconstruction of the Ras switching cycle through saturation mutagenesis.

800        *eLife* **6**, e27810 (2017).

801    32. Brenan, L. *et al*. Phenotypic Characterization of a Comprehensive Set of MAPK1/ERK2

802        Missense Mutants. *Cell Rep*. **17**, 1171–1183 (2016).

803    33. Matreyek, K. A. *et al*. Multiplex assessment of protein variant abundance by massively

804        parallel sequencing. *Nat. Genet*. **50**, 874–882 (2018).

805    34. Mishra, P., Flynn, J. M., Starr, T. N. & Bolon, D. N. A. Systematic Mutant Analyses

806        Elucidate General and Client-Specific Aspects of Hsp90 Function. *Cell Rep*. **15**, 588–598

807        (2016).

808    35. Roscoe, B. P. & Bolon, D. N. A. Systematic exploration of ubiquitin sequence, E1 activation

809        efficiency, and experimental fitness in yeast. *J. Mol. Biol*. **426**, 2854–2870 (2014).

810    36. Kitzman, J. O., Starita, L. M., Lo, R. S., Fields, S. & Shendure, J. Massively parallel single-

811        amino-acid mutagenesis. *Nat. Methods* **12**, 203–206, 4 p following 206 (2015).

812    37. Stiffler, M. A., Hekstra, D. R. & Ranganathan, R. Evolvability as a Function of Purifying

813        Selection in TEM-1 β-Lactamase. *Cell* **160**, 882–892 (2015).

814    38. Rockah-Shmuel, L., Tóth-Petróczy, Á. & Tawfik, D. S. Systematic Mapping of Protein

815        Mutational Space by Prolonged Drift Reveals the Deleterious Effects of Seemingly Neutral

816        Mutations. *PLOS Comput. Biol*. **11**, e1004421 (2015).

817    39. Hewitt, W. M. *et al*. Insights Into the Allosteric Inhibition of the SUMO E2 Enzyme Ubc9.

818        *Angew. Chem. Int. Ed*. **55**, 5703–5707 (2016).

819    40. Cho, L. T.-Y. *et al*. An Intracellular Allosteric Modulator Binding Pocket in SK2 Ion

820        Channels Is Shared by Multiple Chemotypes. *Structure* **26**, 533-544.e3 (2018).

821    41. Klink, B. U., Goody, R. S. & Scheidig, A. J. A newly designed microspectrofluorometer for

822        kinetic studies on protein crystals in combination with x-ray diffraction. *Biophys. J*. **91**, 981–

823        992 (2006).

42. Ward, R. A. *et al*. Structure-Guided Design of Highly Selective and Potent Covalent Inhibitors of ERK1/2. *J. Med. Chem.* **58**, 4790–4801 (2015).

43. Wu, H. *et al*. Structural basis of allele variation of human thiopurine-S-methyltransferase. *Proteins Struct. Funct. Bioinforma.* **67**, 198–208 (2007).

44. Meyer, P. *et al*. Structural basis for recruitment of the ATPase activator Aha1 to the Hsp90 chaperone machinery. *EMBO J.* **23**, 511–519 (2004).

45. Grishin, A. M. *et al*. Structural Basis for the Inhibition of Host Protein Ubiquitination by Shigella Effector Kinase OspG. *Structure* **22**, 878–888 (2014).

46. Hong, M. *et al*. Structural Basis for Dimerization in DNA Recognition by Gal4. *Structure* **16**, 1019–1026 (2008).

47. Minasov, G., Wang, X. & Shoichet, B. K. An Ultrahigh Resolution Structure of TEM-1 β-Lactamase Suggests a Role for Glu166 as the General Base in Acylation. *J. Am. Chem. Soc.* **124**, 5333–5340 (2002).

48. Didovyk, A. & Verdine, G. L. Structural Origins of DNA Target Selection and Nucleobase Extrusion by a DNA Cytosine Methyltransferase *. *J. Biol. Chem.* **287**, 40099–40105 (2012).

49. Meier, J. *et al*. Language models enable zero-shot prediction of the effects of mutations on protein function. 2021.07.09.450648 Preprint at https://doi.org/10.1101/2021.07.09.450648 (2021).

50. Phillips, A. M. *et al*. Binding affinity landscapes constrain the evolution of broadly neutralizing anti-influenza antibodies. *eLife* **10**, e71393 (2021).

51. Koenig, P. *et al*. Mutational landscape of antibody variable domains reveals a switch modulating the interdomain conformational dynamics and antigen binding. *Proc. Natl. Acad. Sci.* **114**, E486–E495 (2017).

847    52. Dreyfus, C. *et al*. Highly Conserved Protective Epitopes on Influenza B Viruses. *Science*

848        **337**, 1343–1348 (2012).

849    53. Ekiert, D. C. *et al*. Antibody Recognition of a Highly Conserved Influenza Virus Epitope.

850        *Science* **324**, 246–251 (2009).

851    54. Fuh, G. *et al*. Structure-Function Studies of Two Synthetic Anti-vascular Endothelial Growth

852        Factor Fabs and Comparison with the Avastin™ Fab *. *J. Biol. Chem*. **281**, 6625–6631

853        (2006).

854    55. Swindells, M. B. *et al*. abYsis: Integrated Antibody Sequence and Structure—Management,

855        Analysis, and Prediction. *J. Mol. Biol*. **429**, 356–364 (2017).

856    56. Westendorf, K. *et al*. LY-CoV1404 (bebtelovimab) potently neutralizes SARS-CoV-2

857        variants. *Cell Rep*. **39**, 110812 (2022).

858    57. Takashita, E. *et al*. Efficacy of Antibodies and Antiviral Drugs against Omicron BA.2.12.1,

859        BA.4, and BA.5 Subvariants. *N. Engl. J. Med*. **387**, 468–470 (2022).

860    58. Research, C. for D. E. and. FDA Announces Bebtelovimab is Not Currently Authorized in

861        Any US Region. *FDA* (2022).

862    59. Cao, Y. *et al*. BA.2.12.1, BA.4 and BA.5 escape antibodies elicited by Omicron infection.

863        *Nature* **608**, 593–602 (2022).

864    60. Cao, Y. *et al*. Rational identification of potent and broad sarbecovirus-neutralizing antibody

865        cocktails from SARS convalescents. *Cell Rep*. **41**, (2022).

866    61. Song, R. *et al*. Post-exposure prophylaxis with SA58 (anti-SARS-COV-2 monoclonal

867        antibody) nasal spray for the prevention of symptomatic COVID-19 in healthy adult

868        workers: a randomized, single-blind, placebo-controlled clinical study*. *Emerg. Microbes*

869        *Infect*. **12**, 2212806 (2023).

870    62. Starr, T. N. *et al*. Shifting mutational constraints in the SARS-CoV-2 receptor-binding

871        domain during viral evolution. *Science* **377**, 420–424 (2022).

872    63. Xu, Y. *et al*. Addressing polyspecificity of antibodies selected from an in vitro yeast

873        presentation system: a FACS-based, high-throughput selection and analytical tool. *Protein*

874        *Eng. Des. Sel.* **26**, 663–670 (2013).

875    64. Makowski, E. K., Wu, L., Desai, A. A. & Tessier, P. M. Highly sensitive detection of

876        antibody nonspecific interactions using flow cytometry. *mAbs* **13**, 1951426 (2021).

877    65. Crawford, K. H. D. *et al*. Protocol and Reagents for Pseudotyping Lentiviral Particles with

878        SARS-CoV-2 Spike Protein for Neutralization Assays. *Viruses* **12**, 513 (2020).

879    66. Rogers, T. F. *et al*. Isolation of potent SARS-CoV-2 neutralizing antibodies and protection

880        from disease in a small animal model. *Science* **369**, 956–963 (2020).

881


882


883

884     **Supplementary Figures, Tables, Information, & Data**

885

886     **Supplementary Table 1:** List of proteins, protein structures, and assay information for deep

887     mutational scanning experiments

888     **Supplementary Table 2**: Analysis of neutralization-enhancing mutations

889

890     **Supplementary Information**: Antibody sequences

891

892     **Supplementary Data 1**: Neutralization data with $IC_{50}$ values of evolved antibodies across both

893     evolutionary campaigns

894     **Supplementary Data 2:** Binding data with IgG $K_D$ values of evolved antibodies

895     **Supplementary Data 3**: Antibody variant prediction benchmarking results

896     **Supplementary Data 4**: MFI values for polyspecificity experiments

897     **Supplementary Data 5**: Efficiency comparison of machine learning-guided directed evolution

898     methods

899
900

**Supplementary Figure 1: Inverse folding identifies high fitness variants across proteins with diverse functions**

In addition to higher hit rates of high fitness variants, inverse folding generally identifies variants with greater magnitude of improvements in fitness. The top ten predicted variants with experimental fitness values ranking in the 20th percentile of all variants profiled in the deep mutational s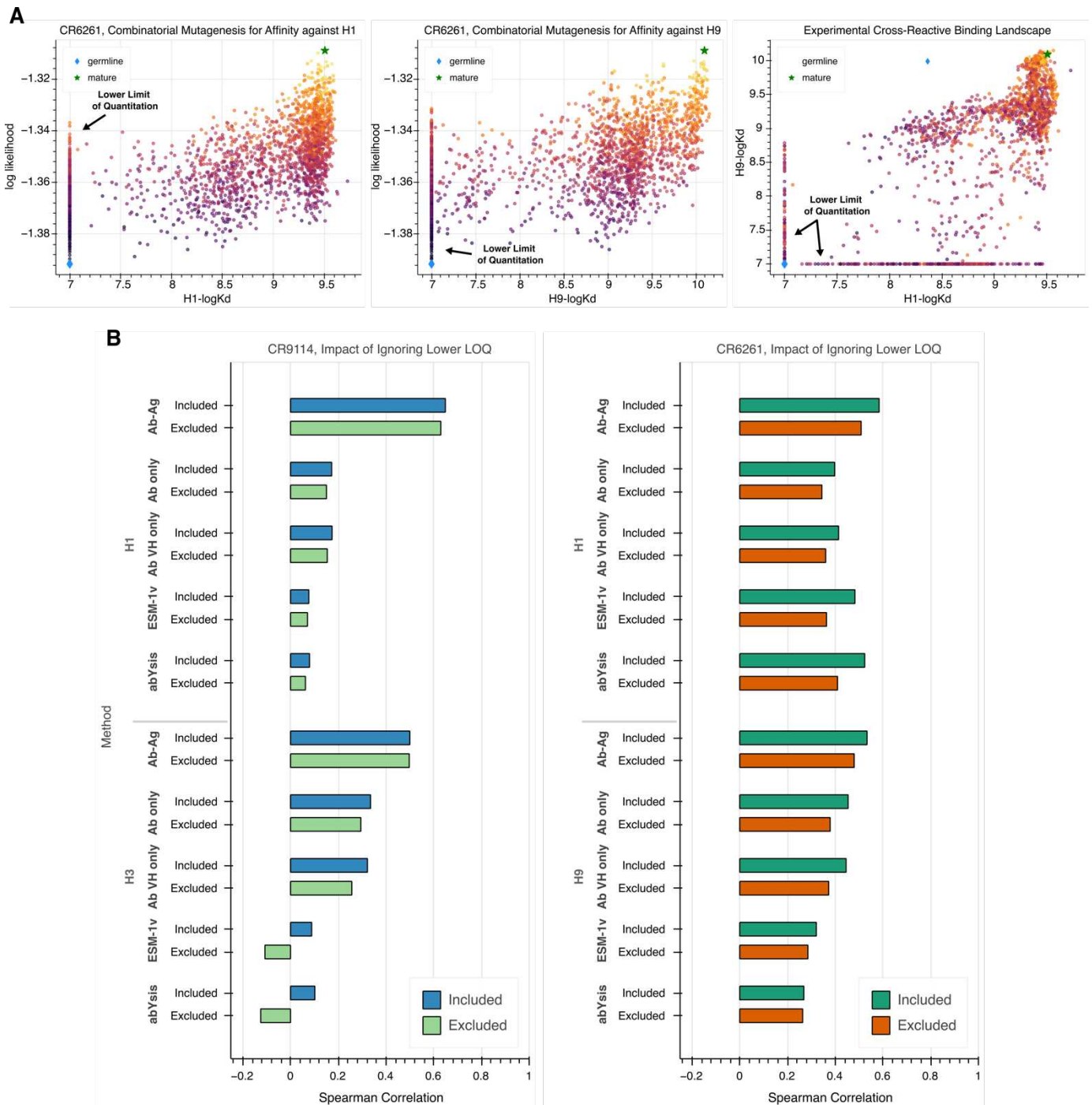creen are shown. The grey curve shows the empirical cumulative distribution function (ECDF) of all experimental fitness values determined in the screen. The dotted lines correspond to the three percentile-based thresholds used in the sensitivity analysis (**Figure 1d**) to classify high fitness variants. bla, Beta-lactamase TEM; CALM1, Calmodulin-1; haeIIIM, Type II methyltransferase M.HaeIII; HRAS, GTPase HRas; MAPK1, Mitogen-activated protein

912      kinase; TMPT, Thiopurine S-methyltransferase; TPK1, Thiamin pyrophosphokinase 1; UBI4,

913      Polyubiquitin; UBE2I, SUMO-conjugating enzyme UBC9

**Supplementary Figure 2: Impact of lower limit of quantitation of binding assay on**

**predictive performance**

(**A**) Scatter plots showing CR6261 variant sequences scored with inverse folding compared to

experimental binding data and inclusive of the assay's lower limit of quantitation, which is

918    omitted for visualization in **Figure 3b**. **(B)** Comparative bar plots showing the impact of

919    removing sequences with experimental measurements bounded artificially by the assay to

920    dataset-wide correlation. While Spearman correlations shown in Figure 3a are computed without

921    any modification to the data, trends in prediction and comparison among modeling methods are

922    robust to filtering sequences affected by this assay artifact.
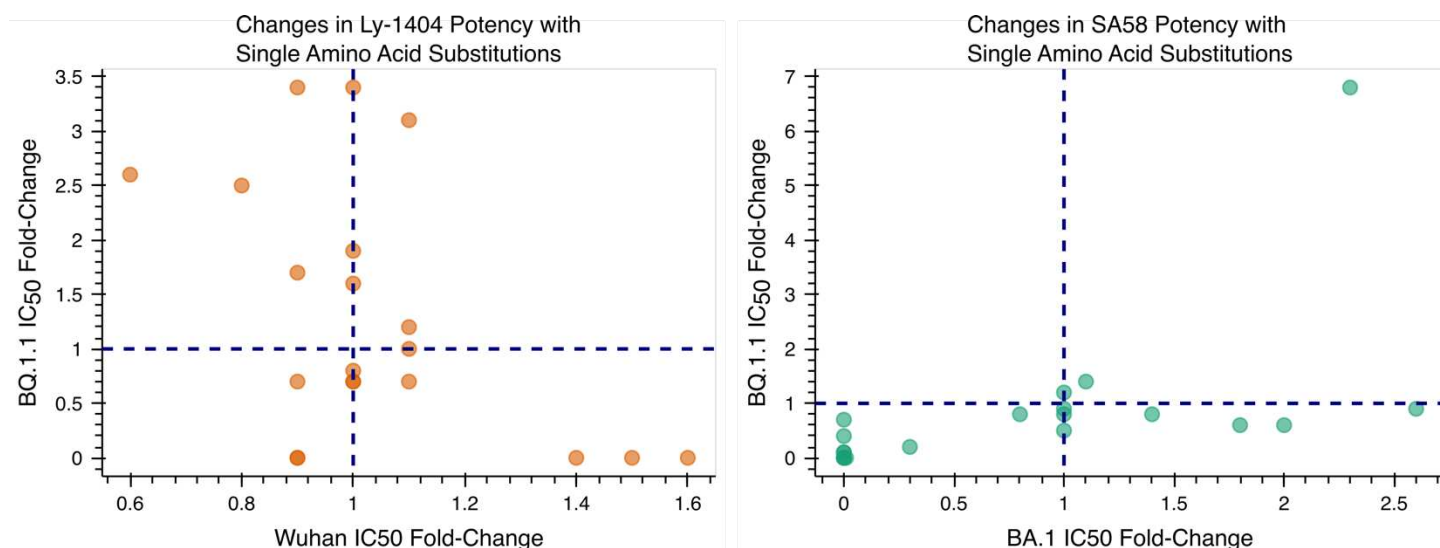
923

**HA of H5N1 influenza,** PDB 4FQI
**HA of H1N1 infleunza,** PDB 7SCO

Backbone RMSD =2.1Å (1416 atoms)
Hamming Distance = 183/499 Amino Acids

924     **Supplementary Figure 3: Structural and sequence similarity of H5 and H1**

925     For cross-reactive antibodies, inclusion of the antigen structure is informative even for predicting

926     binding to a different antigen. In Figure 3a, we report a correlation of 0.65 between inverse

927     folding log likelihoods of CR9114 variants and experimental affinity measurements to H1

928     despite using a structure solved with CR9114 in complex with H5. Inverse folding uses both the

929     protein sequence and backbone structure coordinates as input. Across both HA subunits, H5 and

930     H1 have considerable sequence differences and a 2.1 Å root mean square deviation (RMSD)
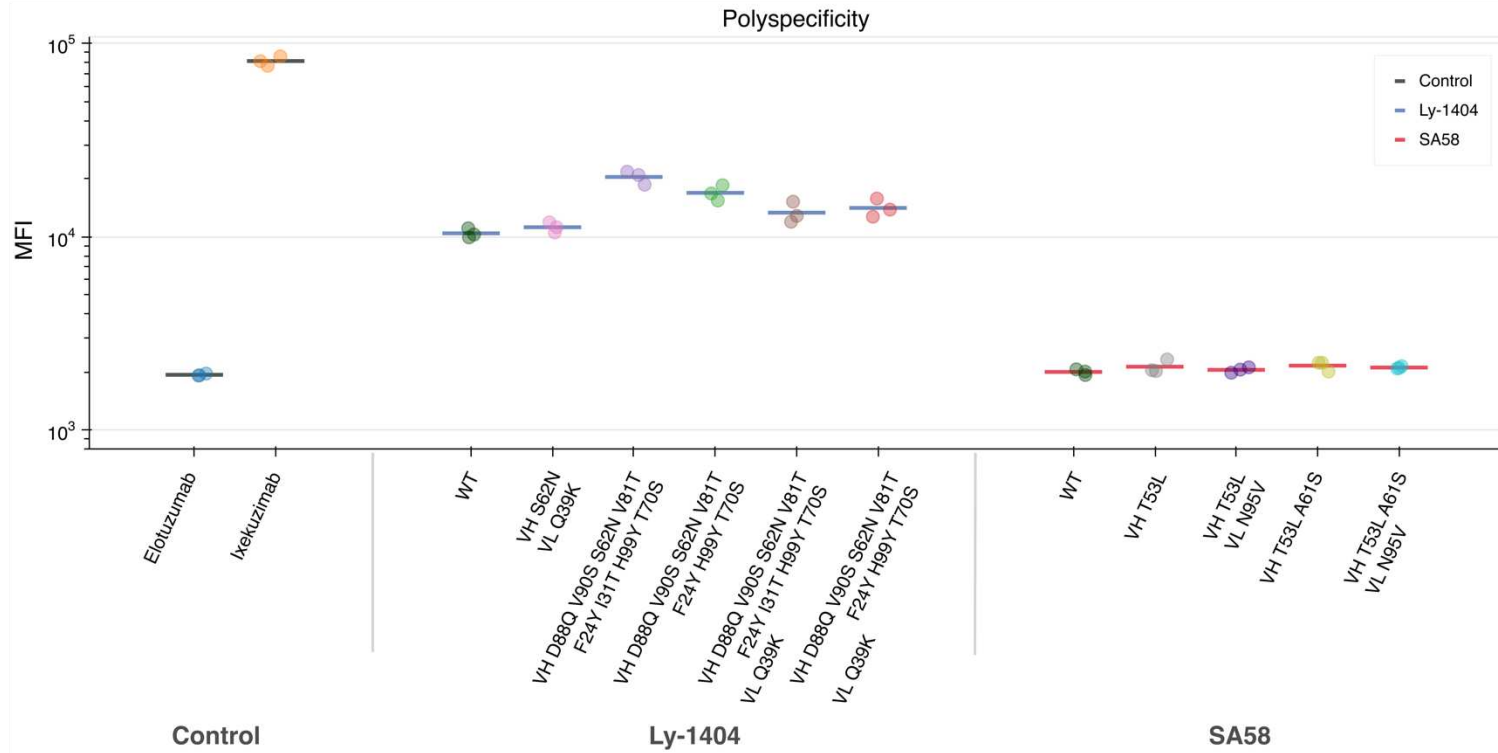
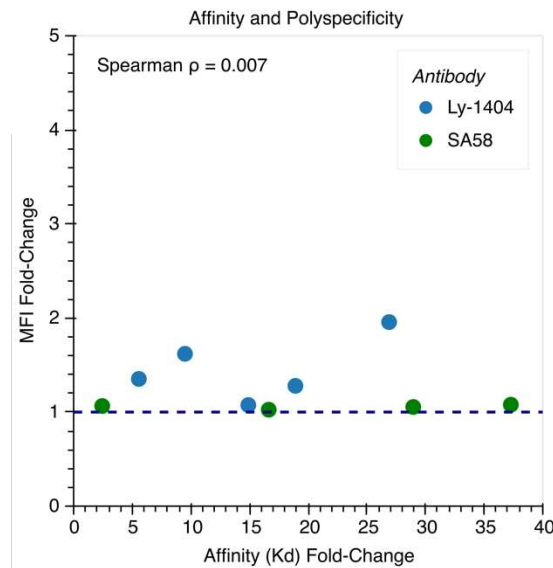931     across the entire protein backbone.

932

**Supplementary Figure 4: Functional diversity of inverse folding-recommended mutations**

Among the 20 single amino acid substitutions tested for Ly-1404, 14 of 20 = 70% improve

neutralization against at least one of the two strains tested. Similarly, 7 of 20 = 35% of the single

amino acid substitutions tested for SA58 improve neutralization. While some variants improve

function against both pseudovirus strains, others overwhelmingly only improve against one. This

suggests that focusing sequence exploration to structurally compatible mutations does not
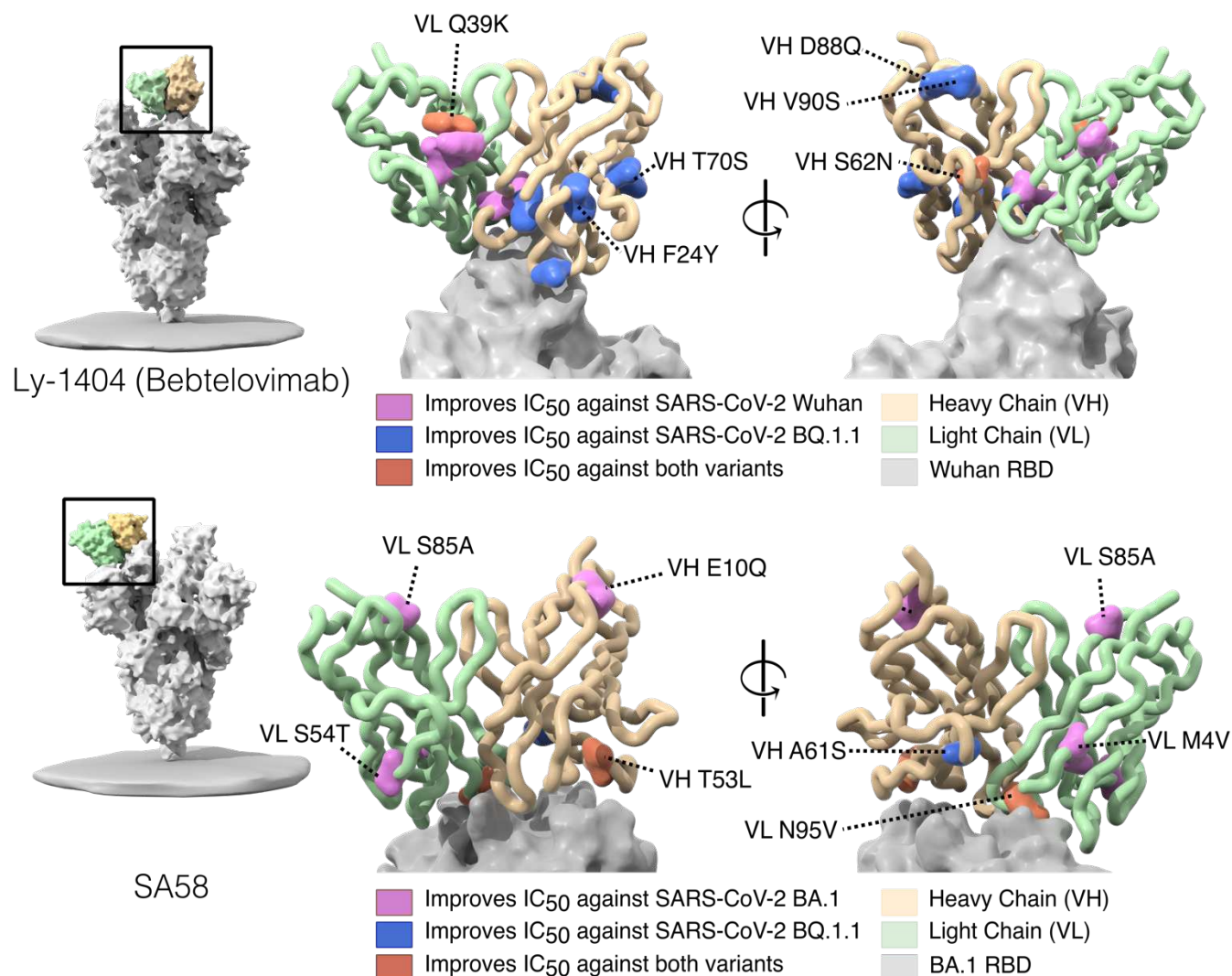
compromise functional diversity.

**Supplementary Figure 5: Polyspecificity of evolved antibodies**

**(A)** The median fluorescence intensity (MFI) signal obtained from flow cytometry is shown for

several evolved antibodies with improved affinity and compared to two clinical monoclonal

945    antibodies with high and low polyspecificity used to define a clinically viable range. **(B)** Fold-

946    change in polyspecificity signal is plotted against fold-change in affinity as IgG against BQ.1.1

947    for Ly-1404 and XBB.1.5 for SA58. There is no correlation between the improvements in on-

948    target improvements in affinity and off-target nonspecific changes in polyspecificity (Spearman ρ

949    = 0.007).

950

951

**Supplementary Figure 6: Mapping neutralization-enhancing substitutions**

Neutralization-enhancing mutations are labeled on the structure of the wild-type antibody in

complex with the RBD of SARS-CoV-2 spike protein (Ly-1404: PDB 7mmo; SA58: PDB

7y0w). Notably, several mutations are identified to have significant beneficial impacts on binding

neutralization and affinity (**Supplementary Data 1 & 2**) despite located away from the binding

interface.

958

**Supplementary Table 1.** Summary of the DMS datasets used in this analysis, including functional assay, method of mutagenesis, and structure used for inverse folding scoring. We also note the specific DMS assay from each study we use for calculating correlation with inverse folding log likelihoods.

| Protein(s) (Uniprot ID) | Organism | Functional Assay | Mutagenesis Method | Utilized assay | PDB Structure | Total coverage of DMS (%) | Access date* | Reference |
|---|---|---|---|---|---|---|---|---|
| UBE2I (P63279) | Human | POPCode, a variant of multiple-site directed mutagenesis. | Competitive growth assay in yeast. | score | 5F6E chain A | 100 | 12/10/2018 | (Weile *et al*, 2017) |
| TPK1 (Q9H3S3) | | | | score | 3S4Y chain A | 92.46 | | |
| CALM1 (P0DP23) | | | | score | 5V03 chain R | 100 | | |
| HRas (P01112) | Human | Systematic site-directed mutagenesis. | Two-hybrid assay. | unregulated | 2CE2 chain X | 100 | 12/10/2018 | (Bandaru *et al*, 2017) |
| MAPK1 (P28482) | Human | Systematic site-directed mutagenesis. | Competitive growth assay. | VRT | 4ZZN chain A | 99.44 | 12/10/2018 | (Brenan *et al*, 2016) |
| TPMT (P51580) | Human | Systematic site-directed mutagenesis. | Fluorescence of a GFP fusion protein. | score | 2BZG chain A | 92.9 | 12/10/2018 | (Matreyek *et al*, 2018) |
| UBI4(b) (P0CG63) | Yeast | Site directed mutagenesis by cassette ligation. | Fluorescence activated cell sorting (FACS). | Relative_E1-activity_limiting | 4Q5E chain B | 100 | 12/10/2018 | (Roscoe & Bolon, 2014) |
| GAL4 (P04386) | Yeast | Systematic site-directed mutagenesis. | Two-hybrid assay. | Nonselection_24 | 3COQ chain B | 90.64 | 12/10/2018 | (Kitzman *et al*, 2015) |
| bla(b) (P62593) | E. coli | Systematic site-directed mutagenesis. | Antibiotic resistance. | Ampicillin_2500 | 1M40 chain A | 100 | 12/10/2018 | (Stiffler *et al*, 2015) |
| haeIIIM (P20589) | H. aegyptius | Random mutagenesis. | Competitive growth assay. | DMS_G3 | 3UBT chain B | 99.37 | 12/10/2018 | (Rockah-Shmuel *et al*, 2015) |

*Access date is as reported in *Livesey & Marsh*, 2020 study from which these data were sourced and this table was adapted

**Supplementary Table 2.** Single amino acid substitutions with beneficial effects on neutralization are reported alongside the region of the variable domain they are located within, as well as the wild-type and mutant amino acid frequencies in observed human antibody sequences.

### Ly-1404

| Chain Mutated | Design | Region | WT Amino Acid Frequency | Mutant Amino Acid Frequency |
|---|---|---|---|---|
| HC | D88Q | HFR3 | 0.03333 | 0.00382 |
| HC | V90S | HFR3 | 0.03316 | 0.05155 |
| HC | S62N | CDR-H2 | 0.13159 | 0.16299 |
| HC | V81T | HFR3 | 0.03432 | 0.00205 |
| HC | F24Y | HFR1 | 0.01738 | 0.00002 |
| HC | I31T | CDR-H1 | 0.00933 | 0.09048 |
| HC | H99Y | HFR3 | 0.01593 | 0.00138 |
| HC | T70S | HFR3 | 0.88405 | 0.06153 |
| HC | I105L | CDR-H3 | 0.02764 | 0.05760 |
| LC | A98I | CDR-L3 | 0.02297 | 0.03198 |
| LC | Q39K | LFR2 | 0.92316 | 0.00238 |
| LC | T5Q | LFR1 | 0.89340 | 0.00933 |
| LC | K47E | LFR2 | 0.52285 | 0.01490 |
| LC | M49L | LFR2 | 0.05585 | 0.77076 |

### SA58

| Chain Mutated | Design | Region | WT Amino Acid Frequency | Mutant Amino Acid Frequency |
|---|---|---|---|---|
| HC | T53L | CDR-H2 | 0.03814 | 0.00963 |
| HC | A61S | CDR-H2 | 0.59797 | 0.13159 |
| HC | E10Q | HFR1 | 0.24182 | 0.01366 |
| LC | N95V | CDR-L3 | 0.13399 | 0.00685 |
| LC | S85A | LFR3 | 0.01109 | 0.00698 |
| LC | S54T | CDR-L2 | 0.65138 | 0.05372 |
| LC | M4V | LFR1 | 0.29424 | 0.03348 |