

1 Simple computational methods can outperform deep learning 2 in designing diverse, binder-enriched antibody libraries

3 Lewis Chinery^{1,*}, Alissa M. Hummer^{1,*}, Brij Bhushan Mehta^{2,*}, Rahmad Akbar²,
4 Puneet Rawat², Andrei Slabodkin², Khang Le Quy², Fridtjof Lund-Johansen²,
5 Victor Greiff², Jeliasko R. Jeliaskov³, and Charlotte M. Deane^{1,**}

6 ¹*Department of Statistics, University of Oxford*, ²*Department of Immunology, University of Oslo*,

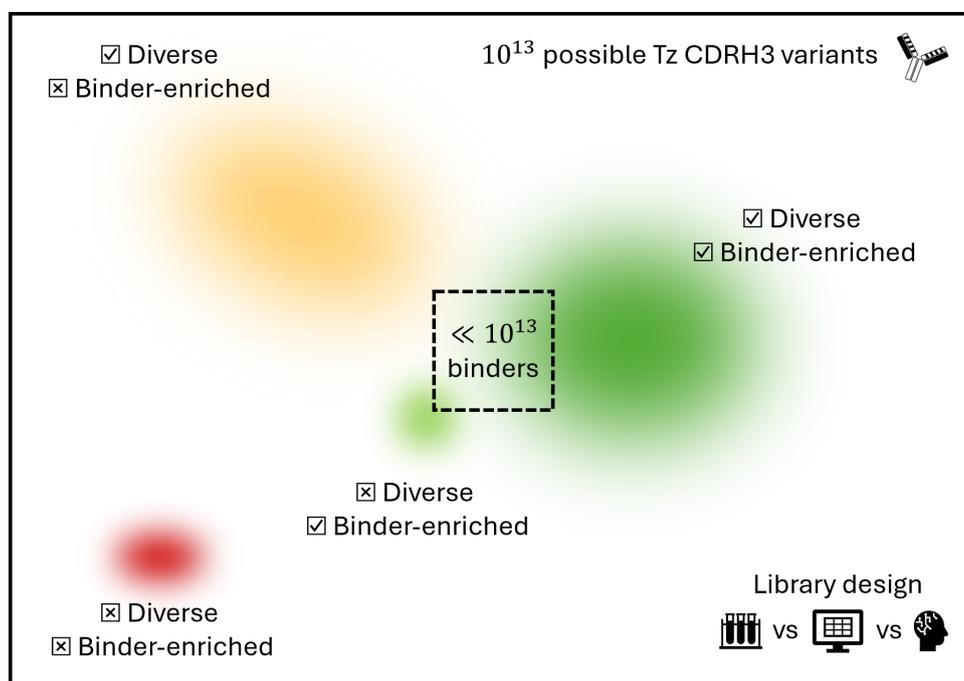
7 ³*Protein Design and Informatics, Computational Sciences, GSK R&D*

8 * *Equal contribution*

9 ** *To whom correspondence should be addressed*

11 Abstract

12 Strong antibody-antigen binding is the primary consideration when developing an effica-
13 cious therapeutic antibody. In recent years, much work has been devoted to applying complex
14 machine learning models to this cause, yet simple baselines are often lacking. Here, we show
15 that the widely used sequence alignment method, BLOSUM, can yield diverse, binder-enriched
16 libraries from a single starting antibody. Using Trastuzumab-HER2 as a model system, we ex-
17 perimentally validated 720 novel designs generated with five different computational methods
18 using surface plasmon resonance. The BLOSUM substitution matrix outperformed all four
19 deep learning design approaches tested, achieving an estimated minimum binder enrichment
20 of 12.5% and producing nine sub-nanomolar binders. These results underscore the importance
21 of comparing against simple baselines and set a benchmark to guide future computational
22 antibody library design.



23 Introduction

24 Therapeutic antibody development is a complex, multiparameter optimization challenge. There
25 are multiple, often competing, properties that influence the efficacy, safety, and developability of
26 an antibody. As a result, it is beneficial to be able to design diverse, binder-enriched antibody
27 libraries to allow for attrition during later stages of optimization.

28
29 Experimental methods for antibody library exploration (for example, Deep Mutational Scan-
30 ning, DMS) and screening (for example, surface display) are resource- and time-intensive. Compu-
31 tational strategies offer promise to more efficiently filter or traverse the search space. Recent work
32 has focused on two primary strategies: training antigen-specific models on experimental data and
33 one-shot library design.

34
35 Several previous studies have leveraged large experimental binding datasets for machine learn-
36 ing (ML) training. For example, Li *et al.* fine-tuned pre-trained large language models on antibody
37 variants for a coronavirus spike protein peptide [1]. The fine-tuning dataset consisted of 26.5k heavy
38 and 26.2k light chain sequences with up to three random Complementarity Determining Region
39 (CDR) mutations from an initial antibody candidate. The resulting models were used to design
40 libraries that consistently yielded >60% of sequences with improved empirical binding over the
41 starting candidate.

42
43 Many other studies have centered on Trastuzumab (Tz), an antibody that binds Human Epi-
44 dermal Growth Factor Receptor 2 (HER2), which is overexpressed in certain breast cancers. Mason
45 *et al.* built a library of 36.4k mutants, with up to 10 CDRH3 mutations, on the basis of the results
46 from a DMS screen [2]. This data was used to train a Convolutional Neural Network (CNN) clas-
47 sifier, which achieved areas under the receiver operating characteristic and precision-recall curves
48 (ROC AUC and PR AUC) of 0.91 and 0.83 for predicting binding versus non-binding variants.

49
50 Akbar *et al.* [3] and Frey *et al.* [4] implemented Recurrent Neural Networks with Long Short-
51 Term Memory and Discrete Walk-Jump Sampling, respectively, to design Trastuzumab variants
52 using this data. Each generative model was trained on ~9k experimentally confirmed binding
53 CDRH3 mutants from ref. [2]. Akbar *et al.* achieved a predicted enrichment of 68%, and wet lab
54 validation confirmed 70% of Frey *et al.*'s designs expressed and maintained affinity to HER2.

55
56 While the trained methods demonstrate strong performance, they are limited by experimental
57 data requirements. One-shot approaches aim to overcome this by leveraging models that have been
58 pre-trained on diverse antibody and/or general protein sequences.

59
60 Hie *et al.* applied the ESM-1b [5] and ESM-1v [6] language models to suggest 'plausible' single-
61 point mutations to antibody sequences, 14-71% of which were shown to improve affinity, depending
62 on the starting wild-type antibody and antigen [7]. To obtain multi-point mutations though, the
63 top predicted single-point mutations were first experimentally tested, similar to DMS. This hybrid
64 strategy reduced the number of single-point mutants that needed testing, but limited the sequence
65 space explored.

66
67 Shanehsazzadeh *et al.* [8] developed structure-based design and inverse-folding models to identify
68 immediate multi-point mutations. This work, which focused on generating Trastuzumab CDRH3
69 sequences, achieved an estimated 10.6% binder-enriched library and resulted in some sequences
70 unlikely to be made following standard DMS. However, the majority of suggested mutations either
71 restored germline residues, or were Glycine and Tyrosine substitutions. Additionally, their models
72 are not publicly available.

73
74 Several open-source ML methods have been developed for protein design tasks, including lan-
75 guage models (for example, refs. [9, 10, 11]) and structure-based inverse folding models (for
76 example, refs. [12, 13, 14, 15]). These models could be applied for antibody library generation,
77 but have not yet been systematically evaluated.

78

79 In the age of deep learning, where models can be limited in generalizability and require exten-
80 sive computational resources for training and/or inference, it is particularly important to quantify
81 the value of newly developed methods. Benchmarks have advanced ML development for diverse
82 protein tasks, as they can expose limitations in existing methods and set a lower bound on per-
83 formance expectations (for example, refs. [12, 16, 17, 18, 19]). Due to rapid advances in antibody
84 library design and binder prediction, the field is lacking a clear baseline performance that new
85 models should aim to exceed.

86

87 Here, we applied computational methods to define baselines for exploring the antibody binder
88 search space. To aid this evaluation, we created the largest publicly available dataset of antibody-
89 antigen affinity measurements (>500,000), for variants of Trastuzumab, building on DMS results
90 [2]. We then evaluated over 700 designs from one-shot libraries generated with BLOSUM, AbLang,
91 ESM-2, and ProteinMPNN using Surface Plasmon Resonance (SPR).

92

93 Experimental validation of a subset of these libraries shows that the simplest method – the
94 BLOSUM substitution matrix – can outperform complex deep learning models for designing binder-
95 enriched antibody libraries. These results underscore the importance of comparing new methods
96 against simple approaches. Additionally, the resources to develop complex methods should be
97 focused to areas where they are likely to provide the greatest value.

98 Results

99 Hundreds of thousands of Deep Mutational Scanning-guided designs main- 100 tain binding to HER2

101 We used Trastuzumab-HER2 as a model system to explore computational antibody library design
102 and iterative binder enrichment. To first baseline our computational methods, we designed an
103 antibody library against HER2 based on the DMS experiment conducted by Mason *et al.* [2].
104 The DMS results – in which every single point mutation had been made at ten positions in the
105 Trastuzumab CDRH3 – were used to weight the design of more than half a million multi-point
106 Trastuzumab variants (further details in Methods).

107

108 To determine whether our DMS-designed variants maintained affinity to HER2, Fluorescence
109 Activated Cell Sorting (FACS) was used to gate the designs. This gating discarded cells with
110 very low expression or antigen binding. The remaining data was split into three bins, resulting in
111 178,160, 196,392, and 171,732 high-, medium-, and low-affinity designs, respectively (Figure S2).
112 We removed sequences that were assigned more than one label (4.0% of sequences), resulting in a
113 dataset of 524,346 Trastuzumab variants, 32.8% of which are high-affinity binders. This dataset
114 will from here on be referred to as ‘Trastuzumab_FACS_524346’.

115

116 To simplify subsequent analyses and align with the goal of selecting only the highest affinity
117 antibodies during lead optimisation, medium- and low-affinity variants were grouped into a single
118 negative class. This grouping aligns with prior work – simple classification methods trained on
119 data from Mason *et al.* offered low predicted binding probabilities for both ‘medium’ and ‘low’
120 classes, and high probabilities for the ‘high’ class (Figure S8).

121

122 The FACS results provide rich experimental data for training machine learning models – a
123 fact we used for prescreening novel computational designs before SPR. The DMS-guided FACS
124 success rate also provides a baseline binder enrichment, which can be used to compare computa-
125 tional library design methods against. At the upper end, 32.8% of the FACS-gated library were
126 high-affinity binders. However, when considering the full estimated library size of $1-2 \times 10^7$ se-
127 quences, including gated and non-gated sequences (based on theoretical diversity and expected
128 transformation efficiency), the binder enrichment rate is ca. 2.6-5.2%.

129 **Computational antibody library design explores diverse areas of sequence**
130 **space**

131 While DMS provides valuable information from which an antibody library can be built, it adds
132 time and costs to the start of an antibody affinity optimisation campaign. Computational meth-
133 ods, conversely, offer immediate, resource-efficient alternatives.

134
135 We explored four existing open-source tools for library design: BLOSUM[20], AbLang[10],
136 ESM-2[9], and ProteinMPNN[12]. Briefly, BLOSUM matrices measure whether or not an amino
137 acid substitution is conservative, AbLang and ESM-2 are protein language models that output
138 masked-residue likelihoods, and ProteinMPNN designs sequences for a given structure (see Meth-
139 ods and original papers for more details). For AbLang, we generated designs masking one CDRH3
140 residue at a time, as well as masking all ten CDRH3 residues at once. For ESM-2, we masked
141 one residue at a time. For each method, we generated 1,000,000 sequences and aimed to retain
142 1,000 sequences that matched the edit distance distribution of Trastuzumab_FACS_524346 (for
143 more details and the final library sizes, see Methods and Figures S11 & S13). This sub-sampling
144 allowed a fairer comparison between methods, as smaller edit distances from Trastuzumab contain
145 proportionally more high-affinity variants than large edit distances.

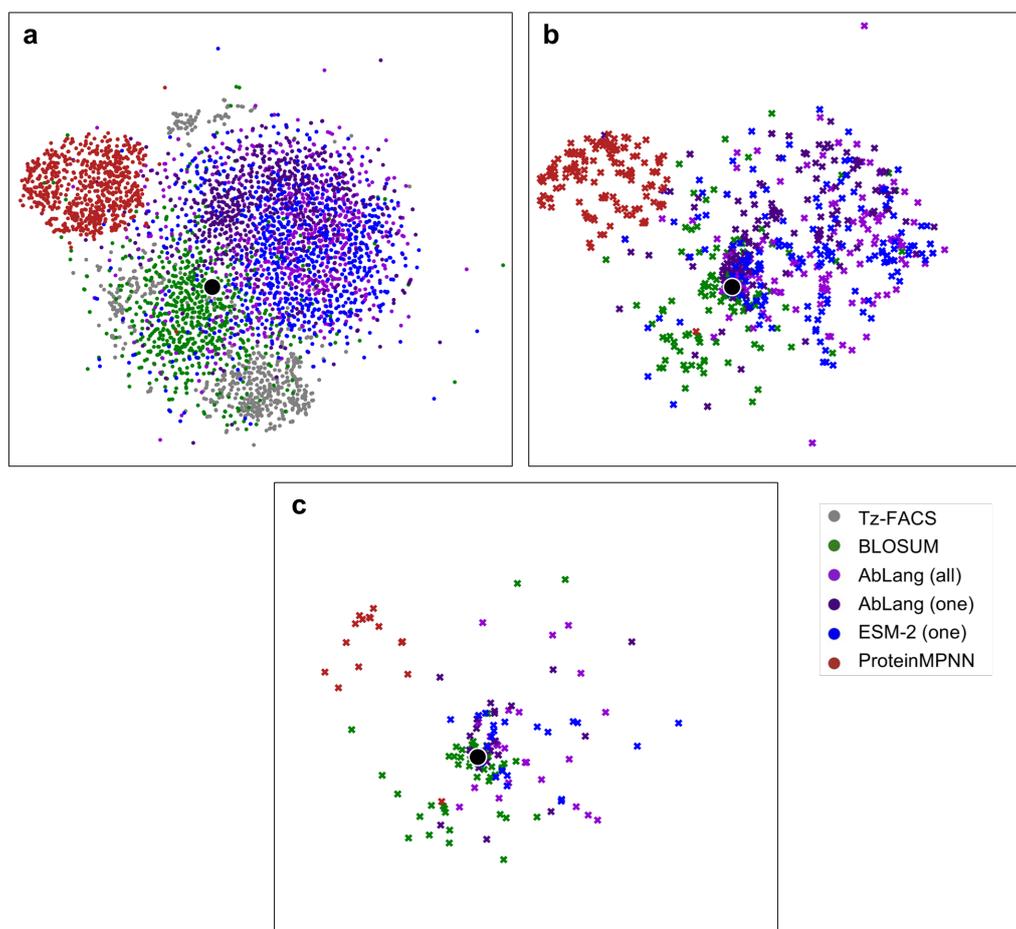


Figure 1: The sequence space explored by different library design methods. (a-c) tSNE plots visualizing (a) approximately 800 randomly selected sequences (to avoid overcrowding the plot) from the DMS-informed library Trastuzumab_FACS_524346 (Tz-FACS) and each of the computational library design methods, (b) sequences selected for experimental validation, (c) experimentally confirmed binders. All t-SNE plots use the same scale. Outlying data points are cropped from (a) for plotting clarity. Wild-type Trastuzumab is shown in black.

146

147 Each computational method produced a diverse library, differing from the sequence spaces
148 explored by DMS and the other design methods (Figure 1a; see Figure S10 for other dimensionality
149 reduction methods).

150 Experimental validation of one-shot computationally designed antibody 151 libraries

152 We experimentally validated designs from the computationally-designed Trastuzumab CDRH3 li-
153 braries using SPR. For each of the five design methods – BLOSUM, AbLang (one residue masked at
154 a time), AbLang (all residues masked), ESM-2 (one residue masked at a time), and ProteinMPNN
155 – we selected 140 sequences for testing.

156

157 The sequences were chosen on the basis of predicted binding, sequence liabilities, and predicted
158 structure retention. As the designed libraries were far larger than could feasibly be tested, we
159 aimed to further enrich the selected subset for binders. We trained a CNN classifier, adapted from
160 ref. [2], on one-hot CDRH3 encodings of the Trastuzumab_FACS_524346 dataset to predict the
161 probability a sequence will retain binding to HER2 (for more details, see Methods and SI Sections
162 5 & 6). The CNN achieved near-perfect test-set accuracy across edit distances one to nine (see
163 Table S2 and Figure S5) and sequences with a binding probability exceeding 90% were shortlisted.
164 Subsequently, constructs containing sequence liabilities (for example, Cysteine, GPR motifs, and
165 long amino acid repeats) and whose predicted CDRH3 structures deviated from an ABodyBuilder2
166 [21] model of wild-type Trastuzumab by $>3.5\text{\AA}$ were filtered out. From the designs that passed all
167 pre-screening steps, sequences for testing were sampled equitably from edit distances one to ten
168 from Trastuzumab, where possible. Further details on the selection of our experimentally tested
169 designs can be found in the Methods. The selected designs covered a diverse area of sequence space
170 (Figure 1b).

171

172 Additionally, AbLang (all residues masked) was used to design 20 sequences with one-
173 residue-shortened CDRH3s to test the loop’s criticality in HER2 binding. Finally, we in-
174 cluded positive controls (wild-type Trastuzumab and true positives from the CNN trained on
175 Trastuzumab_FACS_524346) and negative controls (anti-Respiratory Syncytial Virus antibody
176 and true negatives) (see Methods). Binding affinities measured for wild-type Trastuzumab ranged
177 from 0.029-0.18 nM, overlapping with previous measurements [22, 23]. When exact affinities
178 could not be measured for antibody designs, the sequences were assigned ‘inconclusive binding’ or
179 ‘non-binding’ labels.

180

181 Every computational library design method achieved double-digit success rates, defined as
182 the percent of computationally pre-filtered designs which retained measurable binding to HER2.
183 BLOSUM (48% success rate) significantly outperformed the ML methods (ProteinMPNN 11%,
184 ESM-2 16%, AbLang (one) 16%, AbLang (all) 21%). AbLang achieved a higher success rate when
185 all ten CDRH3 positions were masked than when sequences were designed with more sequence
186 context (one residue masked at a time), although this difference was not statistically significant
187 (see Table S4).

188

189 Binders from the five design methods were observed across a range of affinities, edit distances,
190 and sequence search space (Figure 1c, Figure 2). While the strongest binders typically had lower
191 edit distances, twelve sub-nanomolar affinity binders (nine from BLOSUM) were identified with
192 up to six mutations from Trastuzumab. The highest-affinity antibody tested in our assay, which
193 achieved tighter binding than Trastuzumab, was a design generated by BLOSUM (0.023 nM, edit
194 distance = 3).

195

196 From the experimentally measured binding rates, we calculated ‘floor’ one-shot binder
197 enrichment rates: BLOSUM: $12.5^{+2.2}_{-2.2}\%$, AbLang (all): $5.7^{+1.9}_{-1.7}\%$, AbLang (one): $4.9^{+1.9}_{-1.6}\%$, ESM-2
198 (one): $4.9^{+1.9}_{-1.7}\%$, ProteinMPNN: $2.1^{+1.1}_{-0.8}\%$. These values are equal to the proportion of designs
199 with CNN-predicted binding probabilities above 90% (Figure S13) multiplied by the proportion

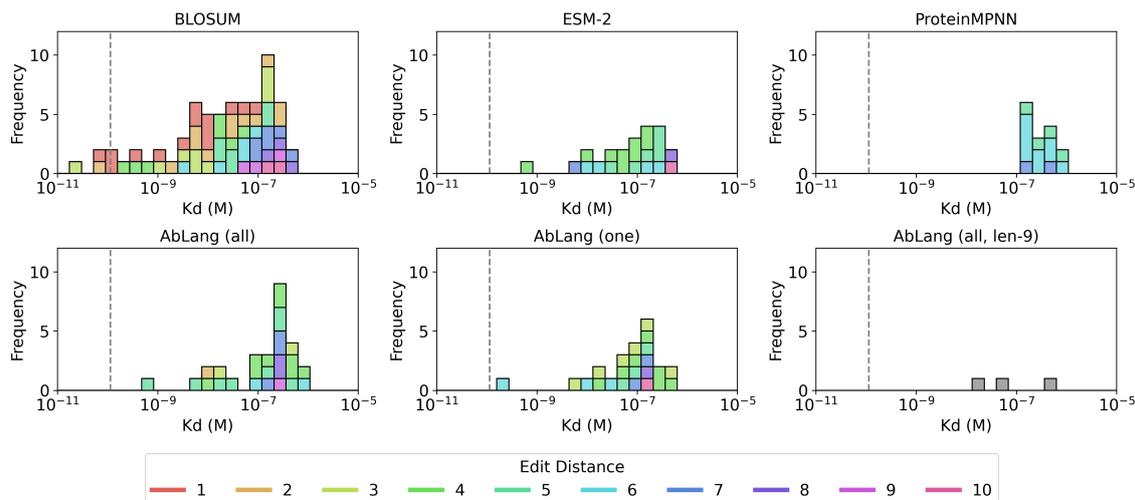


Figure 2: The binding affinities of antibody constructs with measurable HER2 binding, as assessed by SPR. Stacked bars are coloured by the edit distances to Trastuzumab. The AbLang length-9 designs are shown in grey, as edit distances cannot be calculated for length-mismatched sequences. The dashed vertical line shows the mean Trastuzumab affinity (0.1 nM) measured in our SPR experiments.

200 of these that were confirmed as binding by SPR. These rates are described as ‘floor’ enrichment
201 rates as sequences with predicted binding probabilities below 90% may still bind HER2. 95%
202 confidence intervals were calculated using beta functions assuming no prior knowledge (Figure
203 S19). AbLang’s 20 length-nine CDRH3 designs were subject to no *in silico* filtering given the
204 required fixed input size of the CNN, meaning it achieved the highest mean success rate but with
205 the largest uncertainty due to the smaller sample size: $15.0^{+21.3\%}_{-9.6\%}$.

206
207 It is possible to increase these success rates further. Assuming a minimum underlying
208 enrichment rate of 10%, we used our Trastuzumab_FACS_524346 dataset and CNN to simulate a
209 design-experiment-refine continuous learning cycle (Figure S20). After testing just 540 sequences,
210 subsequent rounds of experimentation contained libraries with simulated enrichments above 30%.

211
212 Combined, the computational and experimental results demonstrate successful one-shot anti-
213 body library design is possible with open-source, general-purpose tools.

214 Discussion

215 Antibody optimization is a challenging problem. ML shows great potential for efficiently exploring
216 the huge potential search space and identifying regions that maintain antigen binding, the primary
217 objective in therapeutic development. This has been an area with active development in recent
218 years, but it remains challenging to compare and assess the value of methods. The lack of
219 consistent benchmarks and baselines can result in a poor understanding of the state of the field,
220 as well as an inefficient use of computational resources.

221
222 We applied commonly-used computational tools for sequence generation, in an aim to design
223 antibody libraries that are enriched for binders from a single starting sequence. Conservative
224 estimates based on experimental validation of resulting designs showed that these methods
225 could yield floor binder enrichment rates of 2.1-12.5%. While lower than the rates observed
226 for some other studies which trained on antigen-specific data [3, 4], these enrichment rates
227 represent one-shot expected binding rates with no experimental data or further training required.
228 Additionally, these rates approach and, in the case of BLOSUM, even exceed the binder rates
229 estimated for unfiltered DMS libraries (ca. 2.6-5.2%).

230

231 Of the tested library design methods, BLOSUM proved the most effective by a significant
232 margin, demonstrating that simple methods can be used to design binder-enriched libraries.
233 This supports recent findings that BLOSUM scores correlate with antibody binding affinity [24].
234 Our results may be impacted by the starting point (Table S3) – as BLOSUM’s quantification of
235 mutational tolerance may be most suitable for designing binders from a strong initial construct
236 such as Trastuzumab – and the pre-filtering step – as the BLOSUM designs exhibited greater
237 overlap with the sequence space of the Trastuzumab_FACS_524346 dataset. More complex
238 design methods may offer other advantages and access more distal sequences for weaker starting
239 antibodies.

240

241 The different methods’ designs occupied different areas of the sequence space and, together,
242 the designs resulted in broad sequence diversity. There could therefore be advantages in employing
243 multiple methods in the initial stages of library design, particularly if maximising diversity is a
244 priority. An additional design consideration for future method development, balancing exploration
245 and exploitation, includes how much context to provide the model. While information about
246 the starting sequence appears to be useful for affinity prediction, the higher success rates from
247 AbLang designs generated with no CDRH3 knowledge, rather than single-residue masking,
248 indicate that there is utility in some cases in providing greater design freedom. Notably, the
249 success of AbLang’s shortened designs also suggests that optimising Trastuzumab’s CDRH3 loop
250 may be an oversimplified problem, where avoiding disrupting the position of the other CDR loops
251 is sufficient to retain target affinity.

252

253 Integrating experiments with ML offers a route that leverages simpler, but trainable computa-
254 tional methods. Our simulations of a continuous learning strategy explored the full high-affinity
255 sequence space, but did so more efficiently by avoiding testing of low-affinity designs. This
256 strategy is fast and independent of the target or library design method, and can therefore be easily
257 adapted for any desired property. Effectively integrating lightweight *in silico* and experimental
258 methods into a continuous learning cycle will allow high-affinity antibody variants to be explored
259 at low costs and timescales in most research settings.

260

261 Recent breakthroughs in deep learning, such as generative AI, promise revolutionary changes
262 for drug discovery. We demonstrate, however, that comparatively simple methods may achieve
263 similar performances. It is critical for new methods to be compared against simple baselines to
264 accurately quantify the added value. The utility of simple methods offers a great advantage to the
265 field though as they require drastically fewer computational resources to train and implement.

266 Methods

267 Trastuzumab_FACS_524346 collection methods

268 The Trastuzumab single chain variable fragment (scFv) CDRH3 dataset used in this study,
269 Trastuzumab_FACS_524346, was guided by the site-specific Deep Mutational Scanning results
270 generated previously by Mason *et al.*[2].

271
272 Briefly, a Trastuzumab scFv antibody library was cloned into a pSYD yeast display vector, a
273 variant of the pDNL6 yeast display vector (pSYD uses N-terminal fusion for scFv-aga2 display,
274 while pDNL6 uses a C-terminal fusion of aga2-scFv). The Trastuzumab scFv antibody library
275 cloned in pSYD vector was transformed in EBY100 yeast cells (ATCC #MYA-4941DQ) selected
276 on SD + CAA plates (2% dextrose, 0.67% yeast nitrogen base, and 0.5% casamino acids yeast
277 selection media) at 30°C for 48-72 hours. Yeast display analysis of the Trastuzumab scFv library
278 was performed as described previously by Ferrara *et al.*[25] and Chao *et al.*[26].

279
280 The next day, the cell pellet was resuspended in SG + CAA (containing 2% galactose and
281 0.1% dextrose) at 0.5 OD/ml and incubated at 20°C with shaking for one to two doublings, as
282 determined by OD. The cells were washed with the wash buffer and processed for staining to
283 check HER2 binding. Around 1-10×10⁷ cells were labelled with 100µg/ml anti-V5 tag antibody,
284 followed by the addition of 100nM HER2 and incubated for 30 minutes on ice. The cells were
285 then washed twice more with wash buffer and labelled with a 1:200 dilution of secondary reagents
286 (goat anti-mouse - Alexa 488 and streptavidin-PE).

287
288 Finally, the cells were incubated for 30 minutes on ice, washed twice with a wash buffer,
289 and resuspended in 1ml of sorting buffer. To determine their affinity, the cells were sorted for
290 the brightest V5 FITC0-positive (scFv expression) antigen binding population (PE positive)
291 and labelled as high-affinity binders, as shown in Figure S2. The cells were further sorted
292 for the brightest V5 FITC-positive medium and low-affinity antigen-binding populations. The
293 populations were sorted into tubes containing YPD media and grown in SD + CAA liquid media
294 at 30°C with shaking overnight as described previously[25].

295
296 Plasmid DNA was isolated using a yeast plasmid isolation kit (Zymoprep Yeast Plasmid
297 Miniprep I #D2100) following the user protocol. The variable heavy (VH) gene containing the
298 CDRH3 sequence for each population was PCR amplified using in-house NGS-specific primers.
299 The amplicons were PCR-cleaned and prepared for NGS. The DNA libraries were sequenced on
300 Illumina using NovaSeq 6000 S2 Reagent Kit v1.5 (300 cycles), and the raw data has been de-
301 posited on Zenodo - doi.org/10.5281/zenodo.10549114. The primers used for generating variable
302 heavy amplicons were:

303 NGSVH Fwd: 5´CACCCGTTATGCCGACAG3´
304 NGSVH Rev: 5´GGGATTGGTTGCCGCTAG3´

305
306 The raw paired NGS reads were merged using PEAR (v0.9.6). The subsequent dataset con-
307 sisted of 618,585, 799,368, and 663,397 high, medium and low-affinity unique CDRH3 sequences,
308 respectively. Singleton (count=1) sequences were removed from the dataset to improve the qual-
309 ity of the data. The final Trastuzumab variant dataset comprised 178,160, 196,392, and 171,732
310 sequences in ‘high’, ‘medium’, and ‘low’ affinity binder classes, respectively. The heavy and light
311 chain sequences (from *In8z*) were numbered according to the IMGT scheme. Heavy chain insertion
312 start and stop positions are 107 and 116, respectively (spliced positions are shown in bold).

313 Heavy chain sequence:

314 EVQLVESGGGLVQPGGSLRLSCAASGFNIKDTYIHWVRQAPGKLEWVARIYPTNGYTRYADSVKG
315 RFTISADTSKNTAYLQMNSLRAEDTAVYYCSRWGGDGFYAMDYWGQGTLLVTVSSA

316 Light chain sequence:

317 DIQMTQSPSSLSASVGDRTITCRASQDVNTAVAWYQQKPGKAPKLLIYSASFLYSGVPSRFSGSR
318 SGTDFTLTISLQPEDFATYYCQHYTTPPTFGGKTKVEIKR

319 Computational library design methods

320 We used BLOSUM[20], AbLang[10], ESM-2[9], and ProteinMPNN[12] to generate antibody
321 libraries against HER2, using Trastuzumab as an initial lead. For each method, we generated
322 1,000,000 sequences and aimed to retain 1,000 sequences that matched the edit distance distri-
323 bution of Trastuzumab_FACS_524346. This sub-sampling allowed a fairer comparison between
324 methods, as smaller edit distances from Trastuzumab contain proportionally more high-affinity
325 variants than large edit distances.

326
327 Some methods tended to generate sequences with larger edit distances from Trastuzumab due
328 to the nature of their sampling distributions (Figure S11). In these instances, the true number of
329 sequences sampled at shorter edit distances was sometimes below the target number. The true
330 number of sequences sampled for each method is shown in Figure S13.

331
332 All design methods were run with default parameters unless stated otherwise. Code to recreate
333 these libraries can be found at github.com/oxpig/Tz_her2_affinity_and_beyond.

334 Random library design

335 As a baseline, we randomly mutated each of the ten Trastuzumab residues between positions 107
336 and 116 to each of the 20 standard amino acids with equal probability. The distribution of amino
337 acids to sample from for each of the ten sequence positions was the same (Figure S9).

338 BLOSUM library design

339 BLOSUM matrices (BLOcks SUBstitution Matrices) provide information on which amino acid
340 substitutions are most likely to be observed[20]. These matrices can be used to obtain the
341 frequencies with which we expect to observe each amino acid type replaced with any of the
342 standard 20 amino acids[27].

343
344 For our BLOSUM library design, we used the BLOSUM-45 matrix and, as background, the
345 amino acid frequencies observed in CRDH3s in SAbDab (the Structural Antibody Database)[28, 29]
346 in our reverse calculation (see SI Section 7). These background frequencies describe how often
347 each amino acid is found in a CDRH3, not how often each is likely to be replaced by another.

348
349 Once our final BLOSUM frequencies had been obtained, we used these to weight the sampling
350 of each amino acid type based on the original Trastuzumab residue between positions IMGT 107
351 and 116. The distribution of amino acids sampled from was identical for matching amino acid
352 types in Trastuzumab, e.g. the Glycine residues at positions 108, 109, and 111 (Figure S9).

354 AbLang library design

355 AbLang is an antibody language model designed to restore missing residues[10]. AbLang was
356 trained on over 14m sequences from the Observed Antibody Space database, OAS[30, 31]. This set
357 was dominated by germline sequences, so the restored residues often reflect germline observations.

358
359 We used AbLang to obtain amino acid likelihoods at each sequence position, independent of
360 the original Trastuzumab residues. We tested two methods – masking the entire ten residues
361 between IMGT positions 107 and 116 at once and masking just one residue at a time. In both
362 instances, we limited AbLang to predicting CDRH3s of the same length as Trastuzumab.

363
364 The likelihoods returned by AbLang were used to weight the sampling of each amino acid type
365 between positions IMGT 107 and 116. AbLang’s likelihoods are position-specific, meaning the
366 sampling weights were unique for each sequence position, unlike BLOSUM (Figure S9).

367 **ESM-2 library design**

368 Evolutionary Scale Modeling (ESM) is a general-purpose protein language model (not fine-tuned
369 for studying antibodies) from Meta’s Fundamental AI Research (FAIR) Protein Team[9]. Recent
370 methods have found success in using ESM to suggest affinity-improving single-point mutations[7];
371 here we test its efficacy for multi-site mutations.

372
373 Like AbLang, ESM-2 can be used to obtain amino acid likelihoods at masked sequence
374 positions. We masked each residue between IMGT positions 107 and 116 one at a time and used
375 the 33-layer, 650m parameter implementation of ESM-2 (esm2_t33.650M_UR50D) to suggest
376 residue likelihoods. We also tested masking the entire ten residues (IMGT positions 107 to 116)
377 at once, limiting ESM-2 to predicting CDRH3s of the same length as Trastuzumab. However,
378 predicted enrichments when masking all residues at once were lower than when masking residues
379 one at a time (Figure S14), so we only experimentally tested the latter.

380
381 The likelihoods returned by ESM-2 were used to weight the sampling of each amino acid type
382 between positions IMGT 107 and 116. ESM-2’s likelihoods are position-specific, like AbLang
383 (Figure S9).

384 **ProteinMPNN library design**

385 ProteinMPNN[12] is a deep-learning method for predicting protein sequences for a given structure.
386 We used the ABodyBuilder2[21] predicted structure of Trastuzumab as the base structure to
387 better recapitulate a general design process, as crystal structures are not readily available for
388 most antibodies. CDRH3 residues between IMGT positions 107 and 116 were then masked,
389 and ProteinMPNN was tasked with generating sequences predicted to fit the modelled CDRH3
390 conformation.

391
392 We used a high sampling temperature of 0.3 for ProteinMPNN to produce diverse sequences.
393 Unlike the previous methods described which were used to generate likelihoods to sample from,
394 we used the exact sequences predicted by ProteinMPNN for our generated library. Due to
395 speed limitations and an observed lack of diversity, we generated only 3,000 sequences with
396 ProteinMPNN, which resulted in 2,331 non-redundant outputs.

397
398 Considering the constrained distribution of sequences created by ProteinMPNN, it
399 was not possible to sub-sample from these to match the edit distance distribution of
400 Trastuzumab.FACS.524346. Instead, 1k sequences were randomly sampled from the 2,331 de-
401 signed by ProteinMPNN for comparison against other methods. Further details can be found in
402 SI Section 15.

403 **Prescreening of sequences for liabilities before experimental testing**

404 We experimentally validated 140 designs from each of the five library design methods – BLOSUM,
405 AbLang (masking all positions at once and one at a time), ESM-2, and ProteinMPNN, using
406 Surface Plasmon Resonance (SPR, Twist Bioscience). These 140 shortlisted sequences were
407 subject to a number of filters, described below.

408
409 First, 1,000,000 sequences were generated using each method, excluding ProteinMPNN due
410 to the speed and diversity reasons explained above. Next, we removed redundant sequences and
411 filtered the remainder by their CNN predicted binding probabilities, keeping only sequences with
412 $P_{bind} > 90\%$. For details on the CNN training, see SI Sections 5 & 6.

413
414 Designs were then checked for sequence liabilities: designs containing Cystines or ‘GPR’ motifs
415 (indicative of CD11c/CD18 cross-reactive binding) were removed. Sequences containing four or
416 more of the liabilities listed below were also removed.

417
418 Liabilities (regex) present anywhere in the Fv:

419 N-linked glycosylation: ["N[~P][ST]"]
420 Integrin binding: ["RGD", "RYD", "LDV"]
421
422 Liabilities (regex) present anywhere in the CDRs or Vernier zones:
423 Methionine oxidation: ["M"]
424 Tryptophan oxidation: ["W"]
425 Aspartic acid isomerisation: ["DG", "DS", "DD", "DT", "DH", "DN"]
426 Asparagine deamidation: ["NG", "NS", "NT", "NN"]
427 Lysine glycation: ["KE", "KD", "EK", "ED"]
428 Fragmentation: ["DP", "DQ"]

429
430 To avoid testing Glycine and/or Tyrosine-dominated designs, which appeared often in
431 Shanehsazzadeh *et al.*[8], sequences containing five or more repeats of the same amino acid were
432 dropped as well.

433
434 Finally, designs were modelled using ABodyBuilder2[21], and their RMSDs from a model of
435 Trastuzumab were measured. Models with CDRH3 RMSDs above 3.5Å or total CDR RMSDs
436 above 1.5Å were excluded. From the designs that passed all prescreening steps, sequences for
437 testing were sampled equitably from edit distances one to ten from Trastuzumab, where possible.

438
439 For filtering length-shortened CDRH3 designs, all the above steps, with the exception of the
440 CNN prediction cutoff, were used. To calculate RMSDs, IMGT position 111 in the model of
441 Trastuzumab was ignored.

442
443 All of the sequence liability and structural similarity filtering steps described above are purely
444 computational and quick to run.

445 Selection of positive and negative experimental controls

446 Alongside the 700 new designs (140 for the five library design methods), we tested 68 further
447 sequences, bringing the total to 768 sequences, or eight 96-well plates.

448
449 We included three positive and three negative controls per plate, totaling 48 sequences. Our
450 experimental controls included eight repeats (one per plate) of Trastuzumab as a strict positive
451 control. Eight repeats of an anti-Respiratory Syncytial Virus antibody were included as a strict
452 negative control.

453
454 Alongside these strict controls, 32 sequences from Trastuzumab_FACS_524346 were included
455 as additional controls. Eight sequences labelled as having high affinity and with CNN binding
456 probabilities above 90% (true positives) were included as positive controls. Eight sequences
457 labelled as having low or medium affinity and with CNN binding probabilities below 10% (true
458 negatives) were included as negative controls.

459
460 Eight sequences labelled as having high affinity but with CNN binding probabilities below 10%
461 (false negatives) were included both as softer positive controls and to test the CNN's capability
462 at identifying incorrectly labelled data. Similarly, eight sequences labelled as having low or
463 medium affinity but with CNN binding probabilities above 90% (false positives) were included
464 as our final softer negative controls. All edit distances from one to ten were sampled for our
465 Trastuzumab_FACS_524346 control sequences.

466
467 The remaining 20 sequences were length-shortened CDRH3 loops (shortened by one amino acid,
468 total length 9), designed using AbLang.

469 **Antibody Purification and Surface Plasmon Resonance**

470 The antibody constructs were expressed in human embryonic kidney 293 (HEK293) or Chinese
471 hamster ovary (CHO) cells, using transfection media consisting of Opti-MEM I Reduced Serum
472 Medium, ExpiFectamine Enhancer 1, and ExpiFectamine Enhancer 2. Four days (for HEK293
473 cells) or six days (for CHO cells) post-transcription, cell viability was assessed using the ViaCell
474 instrument to ensure viability remained above 70%. The deep-well plate-containing cells were
475 removed from the incubator and spun down at 500xg for 5 minutes to pellet the cells. The
476 supernatant was transferred to the Qpix deep-well plate. Batch purification was performed using
477 the MabSelect SuRe antibody purification resin. The proteins were eluted in either 0.1 M Glycine,
478 pH 2.7, neutralized with 1 M Tris pH 7.5 or 0.1 M Citrate, pH 2.5, neutralized with 1 M HEPES,
479 pH 7.6. Immediately after purification, the elution plate was sealed with PCR film sealed and
480 stored at 4°C. The affinity of the purified antibody constructs for HER2 was measured using SPR
481 with Carterra LSA screening as in ref. [32].

482 **Data Availability**

483 The raw FACS data can be found on Zenodo - doi.org/10.5281/zenodo.10549114.
484 The SPR results can be found on GitHub - github.com/oxpig/Tz_her2_affinity_and_beyond

485 **Code Availability**

486 Code for designing antibody libraries and classifying binders using all methods presented can be
487 found on GitHub - github.com/oxpig/Tz_her2_affinity_and_beyond

488 **Author contributions**

489 L.C. prepared the training data, investigated the library design methods, and evaluated the FLAML
490 and CNN affinity classification approaches. A.M.H. performed the EGNN analysis. B.B.M. led
491 the Trastuzumab_FACS_524346 data collection. R.A., P.R., A.S., K.L.Q., and F.L.J. supported
492 the data collection and pre-processing. L.C. and A.M.H. analysed the Surface Plasmon Resonance
493 results. L.C., A.M.H., and B.B.M. wrote the manuscript. C.M.D., J.R.J., and V.G. supervised the
494 work. All authors read and approved the final version of the manuscript.

495 **Funding**

496 This work was supported by the Biotechnology and Biological Sciences Research Council (BBSRC)
497 [#BB/V509681/1 awarded to L.C.], the Medical Research Council [#MR/N013468/1 awarded to
498 A.M.H.], the Leona M. and Harry B. Helmsley Charitable Trust [#2019PG-T1D011 awarded to
499 V.G.], UiO World-Leading Research Community [awarded to V.G.], UiO: LifeScience Convergence
500 Environment Immunolingo [awarded to V.G. and F.L.J.], the European Union's Horizon 2020 re-
501 search and innovation programme under the Marie Skłodowska-Curie grant agreement [#801133
502 awarded to P.R.], EU Horizon 2020 iReceptorplus [#825821 awarded to V.G.], a Norwegian Can-
503 cer Society Grant [#215817 awarded to V.G.], Research Council of Norway projects [#300740
504 awarded to V.G., and #331890 awarded to V.G. and F.L.J.], a Research Council of Norway IK-
505 TPLUSS project [#311341 awarded to V.G.], GlaxoSmithKline (GSK), the Innovative Medicines
506 Initiative 2 Joint Undertaking (Inno4Vac, supported by the European Union's Horizon 2020 re-
507 search and innovation programme and EFPIA) [#101007799], and the European Union (ERC,
508 AB-AG-INTERACT) [#101125630].

509 **Ethics declarations**

510 *Competing interests:* C.M.D. discloses membership of the Scientific Advisory Board of Fusion An-
511 tibodies and AI Proteins, as well as a founder of Dalton. V.G. declares advisory board positions in

512 aiNET GmbH, Enpicom B.V, Absci, Omniscope, and Diagonal Therapeutics. V.G. is a consultant
513 for Adaptive Biosystems, Specifica Inc, Roche/Genentech, immunai, LabGenius, and FairJourney
514 Biologics. J.R.J. is employed by Profluent. The remaining authors declare no competing interests.

515 References

- 516 [1] Li, L. *et al.* Machine learning optimization of candidate antibody yields highly diverse sub-
517 nanomolar affinity antibody libraries. *Nature Communications* 2023 14:1 **14**, 1–12 (2023).
518 URL <https://www.nature.com/articles/s41467-023-39022-2>.
- 519 [2] Mason, D. M. *et al.* Optimization of therapeutic antibodies by predicting antigen specificity
520 from antibody sequence via deep learning. *Nature Biomedical Engineering* 2021 5:6 **5**, 600–612
521 (2021). URL <https://www.nature.com/articles/s41551-021-00699-9>.
- 522 [3] Akbar, R. *et al.* In silico proof of principle of machine learning-based antibody design at
523 unconstrained scale. *mAbs* **14** (2022). URL [https://www.tandfonline.com/doi/abs/10.](https://www.tandfonline.com/doi/abs/10.1080/19420862.2022.2031482)
524 [1080/19420862.2022.2031482](https://www.tandfonline.com/doi/abs/10.1080/19420862.2022.2031482).
- 525 [4] Frey, N. C. *et al.* Protein discovery with discrete walk-jump sampling (2023). URL <https://arxiv.org/abs/2306.12360v2>.
- 526
- 527 [5] Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning
528 to 250 million protein sequences. *Proceedings of the National Academy of Sciences* **118**,
529 e2016239118 (2021). URL <https://www.pnas.org/doi/abs/10.1073/pnas.2016239118>.
- 530 [6] Meier, J. *et al.* Language models enable zero-shot prediction of the effects of mutations on
531 protein function. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P. & Vaughan, J. W.
532 (eds.) *Advances in Neural Information Processing Systems*, vol. 34, 29287–29303 (Adv. Neural.
533 Inf. Process. Syst. 34, 2021). URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2021/file/f51338d736f95dd42427296047067694-Paper.pdf)
534 [2021/file/f51338d736f95dd42427296047067694-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/f51338d736f95dd42427296047067694-Paper.pdf).
- 535 [7] Hie, B. L. *et al.* Efficient evolution of human antibodies from general protein language
536 models. *Nature Biotechnology* 2023 1–9 (2023). URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41587-023-01763-2)
537 [s41587-023-01763-2](https://www.nature.com/articles/s41587-023-01763-2).
- 538 [8] Shanehsazzadeh, A. *et al.* Unlocking de novo antibody design with generative artificial intel-
539 ligence. *bioRxiv* 2023.01.08.523187 (2023). URL [https://www.biorxiv.org/content/10.](https://www.biorxiv.org/content/10.1101/2023.01.08.523187v3)
540 [1101/2023.01.08.523187v3](https://www.biorxiv.org/content/10.1101/2023.01.08.523187v3).
- 541 [9] Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language
542 model. *Science* **379**, 1123–1130 (2023). URL [https://www.science.org/doi/abs/10.1126/](https://www.science.org/doi/abs/10.1126/science.ade2574)
543 [science.ade2574](https://www.science.org/doi/abs/10.1126/science.ade2574). Earlier versions as preprint: bioRxiv 2022.07.20.500902.
- 544 [10] Olsen, T. H., Moal, I. H. & Deane, C. M. Ablang: an antibody language model for completing
545 antibody sequences. *Bioinformatics Advances* **2** (2022). URL [https://academic.oup.com/](https://academic.oup.com/bioinformaticsadvances/article/2/1/vbac046/6609807)
546 [bioinformaticsadvances/article/2/1/vbac046/6609807](https://academic.oup.com/bioinformaticsadvances/article/2/1/vbac046/6609807).
- 547 [11] Olsen, T. H., Moal, I. H. & Deane, C. M. Addressing the antibody germline bias and
548 its effect on language models for improved antibody design. *bioRxiv* 2024.02.02.578678
549 (2024). URL [https://www.biorxiv.org/content/10.1101/2024.02.02.578678](https://www.biorxiv.org/content/10.1101/2024.02.02.578678v1https://www.biorxiv.org/content/10.1101/2024.02.02.578678v1.abstract)
550 [v1https://www.biorxiv.org/content/10.1101/2024.02.02.578678v1.abstract](https://www.biorxiv.org/content/10.1101/2024.02.02.578678v1https://www.biorxiv.org/content/10.1101/2024.02.02.578678v1.abstract).
- 551 [12] Dauparas, J. *et al.* Robust deep learning-based protein sequence design using protein-
552 mpnn. *Science* **378**, 49–56 (2022). URL [https://www.science.org/doi/10.1126/science.](https://www.science.org/doi/10.1126/science.add2187)
553 [add2187](https://www.science.org/doi/10.1126/science.add2187).
- 554 [13] Høie, M. H. *et al.* Antifold: improved structure-based antibody design using inverse fold-
555 ing. *Bioinformatics Advances* **5**, vbae202 (2025). URL [https://doi.org/10.1093/bioadv/](https://doi.org/10.1093/bioadv/vbae202)
556 [vbae202](https://doi.org/10.1093/bioadv/vbae202).

- 557 [14] Dreyer, F. A., Cutting, D., Schneider, C., Kenlay, H. & Deane, C. M. Inverse folding for
558 antibody sequence design using deep learning. *ICML Workshop on Computational Biology*
559 (2023).
- 560 [15] Sun, Z.-Y. *et al.* Antbmpnn: Structure-guided graph neural networks for precision anti-
561 body engineering. *Advanced Science* **n/a**, e04278. URL <https://advanced.onlinelibrary.wiley.com/doi/abs/10.1002/advs.202504278>.
562
- 563 [16] Pereira, J. *et al.* High-accuracy protein structure prediction in casp14. *Proteins: Structure,*
564 *Function, and Bioinformatics* **89**, 1687–1699 (2021). URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/prot.26171>.
565
- 566 [17] Sirin, S., Apgar, J. R., Bennett, E. M. & Keating, A. E. AB-Bind: Antibody binding muta-
567 tional database for computational affinity predictions. *Protein Science* **25**, 393–409 (2016).
568 URL <https://onlinelibrary.wiley.com/doi/10.1002/pro.2829>.
- 569 [18] Jankauskaite, J., Jiménez-García, B., Dapkunas, J., Fernández-Recio, J. & Moal, I. H.
570 SKEMPI 2.0: An updated benchmark of changes in protein-protein binding energy, ki-
571 netics and thermodynamics upon mutation. *Bioinformatics* **35**, 462–469 (2019). URL
572 <https://academic.oup.com/bioinformatics/article/35/3/462/5055583>.
- 573 [19] Chungyoun, M., Ruffolo, J. & Gray, J. J. Flab: Benchmarking tasks in fitness landscape in-
574 ference for antibodies. *bioRxiv* (2023). URL [https://www.biorxiv.org/content/10.1101/](https://www.biorxiv.org/content/10.1101/2024.01.13.575504v1)
575 [2024.01.13.575504v1](https://www.biorxiv.org/content/10.1101/2024.01.13.575504v1).
- 576 [20] Henikoff, S. & Henikoff, J. G. Amino acid substitution matrices from protein blocks. *Proceed-*
577 *ings of the National Academy of Sciences of the United States of America* **89**, 10915 (1992).
578 URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC50453/>.
- 579 [21] Abanades, B. *et al.* Immunebuilder: Deep-learning models for predicting the
580 structures of immune proteins. *bioRxiv* 2022.11.04.514231 (2022). URL <https://www.biorxiv.org/content/10.1101/2022.11.04.514231v1https://www.biorxiv.org/content/10.1101/2022.11.04.514231v1.abstract>.
- 583 [22] Baselga, J. Clinical trials of herceptin® (trastuzumab). *European Journal of Can-*
584 *cer* **37**, 18–24 (2001). URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0959804900004044)
585 [S0959804900004044](https://www.sciencedirect.com/science/article/pii/S0959804900004044).
- 586 [23] Komarova, T. V. *et al.* The biological activity of bispecific trastuzumab/pertuzumab plant
587 biosimilars may be drastically boosted by disulfiram increasing formaldehyde accumula-
588 tion in cancer cells. *Scientific Reports* **9**, 16168 (2019). URL [https://doi.org/10.1038/](https://doi.org/10.1038/s41598-019-52507-9)
589 [s41598-019-52507-9](https://doi.org/10.1038/s41598-019-52507-9).
- 590 [24] Uçar, T. & Sormanni, P. Blosum is all you learn — generative antibody models reflect
591 evolutionary priors. *bioRxiv* (2025). URL [https://www.biorxiv.org/content/early/2025/](https://www.biorxiv.org/content/early/2025/10/27/2025.10.26.684652.1)
592 [10/27/2025.10.26.684652.1](https://www.biorxiv.org/content/early/2025/10/27/2025.10.26.684652.1).
- 593 [25] Ferrara, F. *et al.* Using phage and yeast display to select hundreds of monoclonal antibodies:
594 Application to antigen 85, a tuberculosis biomarker. *PLOS ONE* **7**, e49535 (2012). URL
595 <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0049535>.
- 596 [26] Chao, G. *et al.* Isolating and engineering human antibodies using yeast surface display. *Nature*
597 *Protocols* 2006 1:2 **1**, 755–768 (2006). URL [https://www.nature.com/articles/nprot.](https://www.nature.com/articles/nprot.2006.94)
598 [2006.94](https://www.nature.com/articles/nprot.2006.94).
- 599 [27] Eddy, S. R. Where did the blosum62 alignment score matrix come from? *Nature Biotechnology*
600 *2004 22:8* **22**, 1035–1036 (2004). URL <https://www.nature.com/articles/nbt0804-1035>.
- 601 [28] Dunbar, J. *et al.* Sabdab: the structural antibody database URL <http://opig>.

- 602 [29] Schneider, C., Raybould, M. I. & Deane, C. M. Sabdab in the age of biotherapeutics: updates
603 including sabdab-nano, the nanobody structure tracker. *Nucleic Acids Research* **50**, D1368–
604 D1372 (2022). URL <https://academic.oup.com/nar/article/50/D1/D1368/6431822>.
- 605 [30] Kovaltsuk, A. *et al.* Observed antibody space: A resource for data mining next-
606 generation sequencing of antibody repertoires. *The Journal of Immunology* **201**, 2502–
607 2509 (2018). URL [https://journals.aai.org/jimmunol/article/201/8/2502/107069/
608 Observed-Antibody-Space-A-Resource-for-Data-Mining](https://journals.aai.org/jimmunol/article/201/8/2502/107069/Observed-Antibody-Space-A-Resource-for-Data-Mining).
- 609 [31] Olsen, T. H., Boyles, F. & Deane, C. M. Observed antibody space: A diverse database of
610 cleaned, annotated, and translated unpaired and paired antibody sequences. *Protein Science*
611 : *A Publication of the Protein Society* **31**, 141 (2022). URL [https://www.ncbi.nlm.nih.
612 gov/pmc/articles/PMC8740823/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8740823/).
- 613 [32] Yuan, T. Z. *et al.* Rapid discovery of diverse neutralizing sars-cov-2 antibodies from large-
614 scale synthetic phage libraries. *mAbs* **14**, 2002236 (2022). URL [https://doi.org/10.1080/
615 19420862.2021.2002236](https://doi.org/10.1080/19420862.2021.2002236).