

## Protein Design and Variant Prediction Using Autoregressive Generative Models

Jung-Eun Shin<sup>\*1</sup>, Adam J. Riesselman<sup>\*1,2</sup>, Aaron W. Kollasch<sup>\*1</sup>, Conor McMahon<sup>3</sup>, Elana Simon<sup>4,5</sup>, Chris Sander<sup>6</sup>, Aashish Manglik<sup>7,8</sup>, Andrew C. Kruse<sup>§3</sup>, Debora S. Marks<sup>§1,9</sup>

<sup>1</sup> Department of Systems Biology, Harvard Medical School.

<sup>2</sup> Currently at insitro.

<sup>3</sup> Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School.

<sup>4</sup> Harvard College.

<sup>5</sup> Currently at Reverie Labs.

<sup>6</sup> Department of Cell Biology, Harvard Medical School and Department of Data Sciences, Dana-Farber Cancer Institute.

<sup>7</sup> Department of Pharmaceutical Chemistry, University of California San Francisco.

<sup>8</sup> Department of Anesthesia and Perioperative Care, University of California San Francisco.

<sup>9</sup> Broad Institute of Harvard and MIT.

\* These authors contributed equally to this work.

§ Corresponding authors. Email: [Andrew\\_Kruse@hms.harvard.edu](mailto:Andrew_Kruse@hms.harvard.edu) (A.C.K.);

[Debora\\_Marks@hms.harvard.edu](mailto:Debora_Marks@hms.harvard.edu) (D.S.M.)

29        **Abstract**

30        The ability to design functional sequences and predict effects of variation is central to protein  
31        engineering and biotherapeutics. State-of-art computational methods rely on models that  
32        leverage evolutionary information but are inadequate for important applications where multiple  
33        sequence alignments are not robust. Such applications include the prediction of variant effects of  
34        indels, disordered proteins, and the design of proteins such as antibodies due to the highly  
35        variable complementarity determining regions. We introduce a deep generative model adapted  
36        from natural language processing for prediction and design of diverse functional sequences  
37        without the need for alignments. The model performs state-of-art prediction of missense and  
38        indel effects and we successfully design and test a diverse  $10^5$ -nanobody library that shows better  
39        expression than a 1000-fold larger synthetic library. Our results demonstrate the power of the  
40        ‘alignment-free’ autoregressive model in generalizing to regions of sequence space traditionally  
41        considered beyond the reach of prediction and design.

## 42 **Introduction**

43 Over the past twenty years, success in protein engineering has emerged from two distinct  
44 approaches, directed evolution<sup>1,2</sup> and knowledge-based force-field modeling<sup>3,4</sup>. Designing and  
45 generating biomolecules with known function is now a major goal of biotechnology and  
46 biomedicine, propelled by our ability to synthesize and sequence DNA at increasingly low costs.  
47 However, since the space of possible protein sequences is so large (for a protein of length 100  
48 this is  $10^{130}$ ), deep mutational scans<sup>5</sup> and even very large libraries (e.g.  $>10^{10}$  variants) barely  
49 scratch the surface of the possibilities. As the vast majority of possible sequences will be non-  
50 functional proteins, it is crucial to minimize or eliminate these sequences from libraries.  
51 Therefore, the open challenge is to develop computational methods that can accelerate this  
52 search and bias the search space for protein sequences that are likely to be functional. This will  
53 enable design of libraries for tractable high-throughput experiments that are optimized for  
54 functional sequences and variants that are distant in sequence.

55 Antibody design is a particularly challenging problem in the area of statistical modeling of  
56 sequences for the purposes of prediction and design. Antibodies are valuable tools for molecular  
57 biology and therapeutics because they can detect low concentrations of target antigens with high  
58 sensitivity and specificity<sup>6</sup>. Single-domain antibodies, or nanobodies, are composed solely of the  
59 variable domain of the canonical antibody heavy chain. The increasing demand for and success  
60 with rapid and efficient discovery of novel nanobodies using phage and yeast display methods<sup>7-10</sup>  
61 have spurred interest in the design of optimal starting libraries. Previous statistical and structural  
62 modeling of antibody repertoires<sup>11-18</sup> have addressed the characterization of sequences of natural  
63 antibodies or predicted higher affinity sequences from immunization or selection experiments.  
64 One of the biggest challenges is to design libraries diverse enough to target many antigens but  
65 also be well-expressed, stable, and non poly-reactive. In fact, a large, state-of-art synthetic  
66 library contains a substantial fraction of non-functional proteins<sup>8</sup> because library construction  
67 methods lack higher-order sequence constraints. Eliminating these non-functional proteins  
68 requires multiple rounds of selection and poses the single highest barrier to identifying high-  
69 affinity antibodies. In order to circumvent these limitations, there has been emphasis on very  
70 large libraries ( $\sim 10^9$ - $10^{10}$ ) to achieve these desired features<sup>19,20</sup>.  
71 Instead of experimentally producing unnecessarily massive, largely non-functional libraries, we  
72 can design smart libraries of fit and diverse nanobodies for the development of highly specific

73 and possibly therapeutic nanobodies. One way to approach this is to leverage the information in  
74 natural sequences to learn constraints on specific amino acids in individual positions in a way  
75 that captures their dependency on amino acids in other positions. The sequences of these variants  
76 contain rich information about what contributes to a stable, functional protein, and in recent  
77 years generative models of these natural protein sequences have been powerful tools for the  
78 prediction of the first 3D fold from sequences alone<sup>21,22</sup>, to generally more 3D structures and  
79 conformational plasticity<sup>23,24</sup>, protein interactions<sup>25-28</sup>, and most recently, mutation effects<sup>29-34</sup>.  
80 However, these state-of-art methods and established methods<sup>35-38</sup> rely on sequence families and  
81 alignments, and alignment-based methods are inherently unsuitable for the statistical description  
82 of the variable length, hypermutated complementarity determining regions (CDRs) of antibody  
83 sequences, which encode the diverse specific of binding to antigens. While antibody numbering  
84 schemes such as IMGT provide consistent alignments of framework residues, alignments of the  
85 CDRs rely on symmetrical deletions<sup>39</sup>. Alignment-based models are also unreliable for low-  
86 complexity or disordered proteins<sup>40</sup> and cannot handle variants that are insertions and deletions.  
87 Indels make up 15-21% of human polymorphisms<sup>41-43</sup>, 44% of human proteins contain  
88 disordered regions longer than 30 amino acids<sup>40,44</sup>, and both are enriched in association with  
89 human diseases such as cystic fibrosis, many cancers<sup>45,46</sup>, cardiovascular and neurodegenerative  
90 diseases, and diabetes<sup>47,48</sup>.  
91 By contrast, the deep models that have transformed our ability to generate realistic speech such  
92 as text-to-speech<sup>49,50</sup> and translation<sup>51,52</sup> use generative models that do not require “word  
93 alignment”, e.g., between equisemantic sentences, but instead employ an autoregressive  
94 likelihood to tackle context-dependent language prediction and generation. Using this process, an  
95 audio clip is decomposed into discrete time steps, a sentence into words, and a protein sequence  
96 into amino acid residues. Models that decompose high-dimensional data into a series of steps  
97 predicted sequentially are termed autoregressive models, and they are well suited to variable-  
98 length data that have not been forced into a defined structure such as a multiple-sequence  
99 alignment. Autoregressive generative models are uniquely suited for modeling and designing the  
100 complex, highly diverse CDRs of antibodies. Here, we develop and apply a new autoregressive  
101 generative model that aims to capture key statistical properties of sets of sequences of variable  
102 lengths.

103 We first test our method on the problem of prediction of mutation effects, which are typically  
104 analyzed using alignment based statistical methods. The new method performs on par with the  
105 DeepSequence machine-learning VAE-based method<sup>30</sup>, which does require aligned sequences  
106 and which in an independent evaluation, testing against experimental data, was reported to  
107 outperform all currently available methods<sup>34</sup>. In addition to this state-of-the-art performance, our  
108 new alignment-free method is inherently more general. It can deal with a much larger class of  
109 sequences and take into account variable length effects. Another recently developed method<sup>53</sup>  
110 does aim to quantify the of mutation effects without the need for alignments. However, 80% of  
111 the mutational data labelled with experimental outcomes from the same experiments it is tested  
112 on as well as fine-tuning with specific families as input. Previous neural language models<sup>54-56</sup> are  
113 so far not suitable for mutation effect prediction for sequences without extensive experimental  
114 data or sequences with high variability, such as the complementarity-determining regions  
115 (CDRs) of antibody variable domains. By contrast, a fully unsupervised, alignment-free  
116 generative model of functional sequences is therefore desirable for the design of efficient  
117 nanobody libraries.

118 We then trained our validated statistical method on naïve nanobody repertoires<sup>57</sup> as naïve  
119 antibody repertoires have been shown to have functional sequences with capacity to target  
120 diverse antigens<sup>58</sup> and used it to generate probable sequences. In this manner we designed a  
121 sequence library that is 1000-fold smaller than state-of-art synthetic libraries but has an almost  
122 two-fold higher expression level, from which we identified a candidate binder for affinity  
123 maturation. A well designed library can also be used in continuously evolving systems<sup>59</sup> to  
124 combine the hypermutation and affinity maturation processes of living organisms in a single  
125 experiment. Smart library design opens doors to more efficient search methods of nanobody  
126 sequence space for rapid discovery of stable and functional nanobodies.

## 127 **Results**

### 128 **An autoregressive generative model of biological sequences**

129 Protein sequences observed in organisms today result from mutation and selection for functional,  
130 folded proteins over time scales of a few days to a billion years. Generative models can be used  
131 to parameterize this view of evolution. Namely, they express the probability that a sequence  $\mathbf{x}$   
132 would be generated by evolution as  $p(\mathbf{x}|\boldsymbol{\theta})$ , where parameters  $\boldsymbol{\theta}$  capture the constraints essential

133 to functional sequences. An autoregressive model is one that makes a prediction in a time series  
134 (or sequence) using the previous observations. In our context, this means predicting the amino  
135 acid in a sequence using all of the amino acids that come before it. With the autoregressive  
136 model, the probability distribution  $p(\mathbf{x}|\boldsymbol{\theta})$  can be decomposed into the product of conditional  
137 probabilities on previous characters along a sequence of length  $L$  (**Supplementary Fig. 1**) via an  
138 autoregressive likelihood:

$$139 \quad p(\mathbf{x}|\boldsymbol{\theta}) = p(x_1|\boldsymbol{\theta}) \prod_{i=2}^L p(x_i|x_1, \dots, x_{i-1}; \boldsymbol{\theta})$$

140 Many different neural network architectures can model an autoregressive likelihood, including  
141 attention-based models<sup>60</sup> and recurrent neural networks<sup>61</sup>. However, we encountered exploding  
142 gradients<sup>62</sup> during training on long sequence families with LSTM<sup>63</sup> or GRU<sup>64</sup> architectures.  
143 Instead, we parameterize this process with dilated convolutional neural networks  
144 (**Supplementary Fig. 1**), which are feed-forward deep neural networks that aggregate long-  
145 range dependencies in sequences over an exponentially large receptive field<sup>65-67</sup> (See Methods).  
146 The model is tasked with predicting an amino acid at some position in the sequence given all the  
147 previous amino acids in the sequence, i.e. forward language modeling. The causal structure of  
148 the model allows for efficient training to a set of sequences, inference of mutation effects, and  
149 sampling of new sequences. By learning these sequential constraints, the model can be directly  
150 applied to generating novel, fit proteins, one residue at a time. The autoregressive nature of this  
151 model obviates the need for a structural alignment and opens doors for application to modeling  
152 and design of previously challenging sequences such as non-coding regions, antibodies, and  
153 disordered proteins.

### 154 **The autoregressive model predicts experimental phenotype effects from sequences**

155 In order to gain confidence in the new model for generating designed sequences, we first tested  
156 the ability of our new model to capture the dependencies between positions by testing the  
157 accuracy of mutation effect prediction. Somewhat surprisingly, unsupervised, generative models  
158 trained only on evolutionary sequences are proving the most accurate for predicting the effect of  
159 mutations when compared to large datasets of experimentally measured mutation effects<sup>30, 34</sup>, and  
160 they avoid the risk of overfitting that can occur as a result of circularity in supervised methods<sup>68</sup>.  
161 We compared the accuracy of this new, non-alignment-based model to state-of-art methods for a

162 benchmark set of 40 deep mutational scans across 33 different proteins, totaling 690,257  
163 individual sequences (**Supplementary Table 1**).

164 The autoregressive model was first fit to each family of protein sequences and then we used the  
165 log-ratio of likelihoods of individual sequences to predict mutation effects:

$$166 \log \frac{p(\mathbf{x}^{Mutant} | \boldsymbol{\theta})}{p(\mathbf{x}^{Wild-type} | \boldsymbol{\theta})}$$

167 which estimates the plausibility of mutant sequence  $\mathbf{x}^{Mutant}$  relative to its wild-type, un-mutated  
168 counterpart,  $\mathbf{x}^{Wild-type}$ . This log-ratio has been shown to be predictive of mutation effects<sup>29, 30</sup>.

169 Importantly, this approach is fully unsupervised: rather than learning from experimental mutation  
170 effects, we can learn evolutionary constraints using only the space of natural sequences. We  
171 benchmark the model predictions against the deep mutational scan experiments and compare the  
172 Spearman's rank correlation to state-of-art models trained on alignments of the same sequences.

173 The autoregressive model is able to consistently match or outperform a model with only site-  
174 independent terms (30/40 datasets) and the EVmutation model<sup>29</sup> that includes dependencies  
175 between pairs of sites (30/40 datasets); it performs on par with the state-of-the-art results of  
176 DeepSequence<sup>30</sup> (19/40 datasets, average difference in rank correlation is only 0.09); and it  
177 outperforms the supervised Envision model<sup>31</sup> for 6/9 of the datasets tested (**Fig. 2a**;

178 **Supplementary Figs. 2, 3**). Previously published benchmarks<sup>29</sup> demonstrate the higher  
179 accuracy of the probabilistic models, EVmutation compared to SIFT and PolyPhen, and recent  
180 work demonstrates that DeepSequence outperforms all currently available methods when  
181 measured against experimental mutation scans<sup>34</sup>. These benchmarks, taken together with our  
182 previous benchmarks<sup>29</sup> and evidence from independent assessments<sup>34</sup>, show that our  
183 autoregressive model outperforms all methods including supervised and performs on par with  
184 our own state-of-art alignment-based method<sup>30</sup> for single mutation effect prediction, providing us  
185 with the confidence to use the model for sequence design.

186 As with previous models that use evolutionary sequences, the accuracy of mutation effect  
187 prediction increases with increasing numbers of non-redundant sequences, as long as there is  
188 coverage of the length, tested here across eight of the protein families for four sequence depths  
189 (**Supplementary Fig. 4, Supplementary Table 2**). Interestingly, the accuracy of effect  
190 predictions against the aliphatic amidase mutation scan are remarkably robust even with a low



191 number of training sequences—123 non-redundant sequences provide the same accuracy as  
192 36,000—suggesting that there is more to learn about the relationship between evolutionary  
193 sampling and model learning. For now, we advise a conservative  $M_{eff}/L$  (number of effective  
194 sequences normalized by length) requirement of 5 in order to sample enough diversity.

195 Because the autoregressive model is not dependent on alignments, we can now learn mappings  
196 of sequences of high variability and diverse lengths for which meaningful alignments are  
197 difficult or non-sensical to construct, such as antibody and nanobody sequences. The  
198 autoregressive model was thus also validated on nanobody thermostability measurements to test  
199 whether we could learn the sequence constraints of fit nanobodies, including the highly variable  
200 regions. To do so, we fit the autoregressive model to a set of ~1.2 million natural llama  
201 nanobody sequences<sup>57</sup>. Sequence likelihoods from this trained model are expected to reflect  
202 nanobody fitness, i.e., the multiple convolved aspects that nanobodies are selected for *in vivo*,  
203 including thermostability, expression, and potentially low polyreactivity. Using this model, we  
204 find that the log-probability fitness calculations predict the thermostability of unseen llama  
205 nanobody sequences from four different stability experiments<sup>69-72</sup> (**Fig. 2b, Supplementary Fig.**  
206 **5, Supplementary Table 3**). These experiments span a wide range of mutation types, lengths,  
207 and sequence diversity. The autoregressive model consistently outperforms a hidden Markov  
208 model (HMM, hmmer3)<sup>73, 74</sup> in predicting the relationship between sequence and thermostability  
209 of nanobodies.

210 Previous alignment-dependent generative models are constrained to predicting the effects of  
211 missense mutations. However, in-frame insertions and deletions can also have large phenotypic  
212 consequences for protein function, yet these changes have proved difficult to model. We  
213 compare the fitness predictions calculated as log probabilities by the autoregressive model to  
214 experimental assays for the fitness of mutated biomolecules, using rank correlation ( $\rho$ ) for  
215 quantitative measurements and area under the receiver-operator curve (AUC) for binary fitness  
216 categorization, identifying the two groups with a two-component Gaussian mixture model. The  
217 model is able to capture the effects of single amino acid deletions on PTEN phosphatase<sup>75</sup>  
218 ( $\rho=0.69$ ,  $N=340$ , HMM  $\rho=0.75$ ; PROVEAN  $\rho=0.7$ ; **Fig. 2c**) and multiple amino acid insertions  
219 and deletions in imidazoleglycerol-phosphate (IGP) dehydratase<sup>76</sup> (AUC=0.90,  $N=6102$ , HMM  
220 AUC=0.88; **Fig. 2d, Supplementary Table 4**). Here we use the AUROC metric for IGP



221 dehydratase as the experimental data are bimodal with a large fraction at zero fitness. While  
222 PROVEAN<sup>77</sup> predicted the effect of single PTEN deletions comparably to our model, it fails to  
223 predict the effect of multiple insertions, deletions, and substitutions as were tested in IGP  
224 dehydratase and it cannot generate new sequences. Three additional insertion and deletion  
225 mutation scan fitness predictions are included in the supplement: yeast snoRNA ( $\rho=0.49$ ), beta  
226 lactamase ( $\rho=0.45$ ), and p53 ( $\rho=0.035$ ; **Supplementary Fig. 6**). Predicting the effects of indels  
227 also has clinical significance: the four different single amino acid deletions annotated as  
228 pathogenic by Clinvar<sup>78</sup> in two cancer genes, BRCA1 and P53, and one Alzheimer's-linked gene,  
229 APOE, are in the bottom 25<sup>th</sup> percentile of predicted deletion effect distributions  
230 (**Supplementary Fig. 7**). Other indels that are predicted to be highly deleterious by the  
231 autoregressive model may be of clinical interest for experimental study of pathogenicity. We  
232 expect that the autoregressive model can predict mutation effects in disordered and low-  
233 complexity sequences. As a proof-of-concept, we have provided an *in silico* mutation scan of the  
234 human tau protein, which contains regions of low complexity and is strongly associated with  
235 neurodegenerative diseases, (**Supplementary Fig. 8**). Our mutation effect prediction  
236 distinguishes between 40 pathogenic and 10 non-pathogenic mutations (two-tailed independent  
237  $t=-4.1$ ,  $P=0.001$ ,  $AUC=0.86$ ) that were collected from the Alzforum database<sup>79</sup>.

### 238 **Generating an efficient library of functional nanobodies**

239 Screening large, high-throughput libraries of antibodies and nanobodies *in vitro* has become  
240 increasingly prevalent because it can allow for rapid identification of diverse monoclonal binders  
241 to target antigens. However, these synthetic libraries contain a large fraction of non-functional  
242 nanobody sequences. Natural nanobody sequences are selected against unfavorable biochemical  
243 properties such as instability, poly-reactivity, and aggregation during affinity maturation<sup>6</sup>.  
244 Similarly to nanobody thermostability prediction, we sought to learn the constraints that  
245 characterize functional nanobodies by fitting the autoregressive model to a set of ~1.2 million  
246 nanobody sequences from the immune repertoires of seven different naïve llamas<sup>57</sup>. Using this  
247 trained model and conditioning on the germline framework-CDR1-CDR2 nanobody sequence,  
248 we then generate over  $10^7$  fit sequences, generating one amino acid at a time based on the  
249 learned sequential constraints. As nanobody CDR3s often contact the framework in 3D,  
250 conditioning in this way allows the model to learn any resulting constraints on the CDR3

251 sequence and incorporate them during generation. We remove sequences that do not end with the  
252 final beta strand of our nanobody template, duplicate sequences, and CDR3s likely to suffer post-  
253 translational modification to obtain ~3.7 million sequences (**Supplementary Table 5**). From  
254 these, we select 185,836 highly diverse CDR3 sequences for inclusion in our designed library.  
255 We compare our designed library to a state-of-art synthetic library<sup>8</sup>, which was constructed  
256 combinatorically based on position specific amino acid frequencies of nanobody sequences with  
257 crystal structures in the PDB database. This library contains CDR3 sequences that have a similar  
258 distribution of biochemical properties as the naïve llama immune repertoire (Methods; **Fig. 3a**).  
259 The distribution of hydrophobicity and isoelectric points are similar to the natural llama  
260 repertoire even though explicit constraints on these properties were never imposed during  
261 generation or selection of sequences for the designed library. The lengths of the CDR3 sequences  
262 in the designed library are shorter than the natural repertoire; this is due to the strategy of  
263 choosing cluster centroids during selection of the 10<sup>5</sup> sequences and can be adjusted by changing  
264 the sampling method. Longer CDR3s may also be attained by allowing interloop disulfide  
265 bridges that stabilize longer CDR3s in some VHH domains<sup>80</sup>; this would require a different  
266 nanobody template and ideally camel or dromedary nanobody repertoires. The sequences in the  
267 designed library are extremely diverse and are more distant from each other than sequences in  
268 the natural repertoire (**Fig. 3b**), while maintaining nearly as much diversity as an equivalent  
269 sample of a combinatorial synthetic library<sup>8</sup> (**Supplementary Fig. 9**). Additionally, we are  
270 exploring new regions of sequence space because the generated sequences in the designed library  
271 are diverse from the naïve repertoire (**Fig. 3c**).

272 Using these designed CDR3 sequences, a nanobody library was constructed using our yeast-  
273 display technology for experimental characterization alongside a combinatorial synthetic  
274 nanobody library<sup>8</sup>. The designed library had more length diversity and a longer CDR3 median  
275 length (13) than the synthetic library (12) (**Supplementary Fig. 9**), while the synthetic library  
276 included designed diversity in specific residues of the CDR1 and CDR2. Individual nanobody  
277 sequences were expressed on the surface of yeast cells, allowing for rapid sorting of nanobody  
278 clones based on expression and/or binding levels. Upon induction, the designed nanobody library  
279 contained 1.5 times higher proportion of cells expressing and displaying nanobodies on their cell  
280 surface than the synthetic nanobody library (**Fig. 4a,b, Supplementary Fig. 10**). In the designed  
281 library, we can also see a clearer separation of cells expressing nanobodies and those that are not.

282 Of cells expressing nanobodies, the mean nanobody display levels from the designed library is  
283 almost twice the level of the previous library (**Fig. 4a,b**). Furthermore, the designed library had  
284 nearly half the fraction of poorly expressed nanobodies (cells with fluorescence below 10,000  
285 AU) as compared to the synthetic library (**Fig. 4a,b**) as well as a significant increase in the  
286 fraction of highly expressed nanobodies as can be seen in the upper limits in the respective  
287 expression distributions (**Fig. 4a, Supplementary Fig. 10**). Expression experiments were  
288 performed with two replicates in addition to a single control experiment of yeast expressing a  
289 single well-behaved nanobody clone (Nb. 174684). These experimental results demonstrate that  
290 with the autoregressive model trained on natural llama nanobody sequences, we successfully  
291 designed a smart library consisting of a higher proportion of stable, well-expressed nanobodies.

292 With this small designed library, we selected nanobody sequences that bound to human serum  
293 albumin (HSA) using fluorescence activated cell sorting (FACS) (**Fig. 4c**), from which we were  
294 even able to identify weak to moderate binders—the strongest binder has a predicted  $K_d$  of 9.8  
295  $\mu M$  (**Fig. 4d**). This experiment is a proof-of-concept that this small library contains antigen-  
296 binding sequences that can be starting points for affinity maturation to identify strong binders.  
297 Though not explicitly designed to minimize poly-reactive nanobody sequences, training on a  
298 naïve llama repertoire, which presumably contain a moderate proportion of poly-reactive  
299 sequences<sup>81-87</sup>, the designed library shows similar levels of poly-reactivity to the synthetic  
300 library, which had been designed according to a small set of highly specific nanobodies  
301 (**Supplementary Fig. 11**). These results indicate that we have successfully designed an efficient  
302 library containing a high proportion of promising diverse, stable, specific, and sensitive  
303 nanobody sequences.

## 304 **Discussion**

305 Here we show how neural network-powered generative autoregressive models can be used to  
306 model sequence constraints independent of alignments and design novel functional sequences for  
307 previously out of reach applications such as nanobodies. The capability of these models is based  
308 on demonstrated state-of-the-art performance and on an extended range of applicability in the  
309 space of sequences. In the particular version in this paper, we validated our model first on deep  
310 mutational scan data, with on par performance with the best currently available model<sup>29-31, 34, 77</sup>,  
311 and demonstrated application to examples for which robust alignments cannot be constructed,

312 such as sequences with multiple insertions, deletions, and substitutions, and cases for which  
313 protein structures and experimental data are not available. As for comparison with a potentially  
314 competing alignment-free model, while we do not discount the utility of semi-supervised  
315 methods (exploiting mutation effect-labeled experimental data), great care must be taken in the  
316 way the split between training and test is conducted to evaluate the true generalizability of the  
317 method. For instance, randomized subsets excluded from training will still be learned from the  
318 labeled data in a way that is not generalizable to required predictions for other proteins<sup>53,88,89</sup>.  
319 Our model is not subject to these limitations as its training is fully unsupervised.

320 Due to their flexibility, deep autoregressive models could also open the door to new  
321 opportunities in biological sequence analysis and design. Unlike alignment-based techniques,  
322 since no homology between sequences is explicitly required, generative models with  
323 autoregressive likelihoods can be applied to variants with insertions and deletions, disordered  
324 proteins, multiple protein families, promoters and enhancers, or even entire genomes.  
325 Specifically, prediction of insertions and deletions and mutation effects in disordered regions has  
326 been a difficult research area, despite their prevalence in human genomes. Disordered regions are  
327 enriched in disease-associated proteins, so understanding variant effects will be important in  
328 understanding the biology and mechanism of genes indicated in cardiovascular, cancer, and  
329 neurodegenerative diseases. For example, classical tumor suppressor genes, such as p53,  
330 BRCA1, and VHL, and proteins indicated in Alzheimer's disease, such as Tau, have long  
331 disordered regions where these models may prove particularly useful.

332 With this model, we designed a smart, diverse, and efficient library of fit nanobody sequences  
333 for experimental screening against target antigens. Designing individual hypervariable CDR  
334 sequences that make up a library of diverse, functional, and developable nanobodies allows for  
335 much faster and cheaper discovery of new therapeutics, minimizing both library waste and  
336 necessary experimental steps. Our streamlined library (1000-fold smaller than combinatorial  
337 synthetic libraries) enables rapid, efficient discovery of candidate nanobodies, quickly providing  
338 a starting point for affinity maturation to enhance binding affinity. In combination with a  
339 continuous evolution system, candidate binders from the designed library have been identified  
340 and affinity matured after only a few rounds of selection with a single experiment<sup>90</sup>. As the cost  
341 to synthesize sequences decreases, the demand for methods that can design highly optimized and

342 diverse sequences will increase as compared to constructing libraries via random or semi-random  
343 generation strategies.

344 A challenge of using synthetic libraries is the poly-reactivity of many sequences that *in vivo*,  
345 would be cleared by an organism's immune system. Naïve llama repertoires also contain poly-  
346 specific sequences, so training a model on sequences from mature or memory B cell repertoires  
347 may provide information on how to improve library design in the future and minimize the poly-  
348 reactivity of the designed library sequences. Multi-chain proteins such as antibodies present an  
349 additional challenge that multiple domains must be designed together. Models incorporating  
350 direct long-range interactions such as dilated convolutions or attention may identify the relevant  
351 dependencies between domains, even when the domains simply concatenated and generated  
352 sequentially. Paired antibody chains are more challenging to sequence than nanobodies, but more  
353 repertoires are becoming available<sup>91</sup>. Beyond antibody and antibody fragment libraries, this  
354 method is translatable to library design for any biomolecule of interest, including disordered  
355 proteins.

356 Our model is the first alignment-free method demonstrating state-of-art mutation effect  
357 prediction without experimental data and applied to at scale to design of protein sequences. New  
358 developments in machine learning will enhance the power of such autoregressive models and  
359 incorporating protein structural information may further improve the capacity to capture long-  
360 range dependencies<sup>92</sup> for these applications. The addition of latent variables could also allow for  
361 targeted design of high affinity and specificity sequences to a desired target antigen<sup>56, 93-95</sup>.  
362 Conversely, we also anticipate better exploration of broader spans of sequence space for  
363 generation, either by exploiting variance explained by latent variables<sup>96</sup> or diverse beam search  
364 strategies<sup>97</sup>. With the increased number of available sequences and growth in both computing  
365 power and new machine learning algorithms, autoregressive sequence models may enable  
366 exploration into previously inaccessible pockets of sequence space.

367

## 368 **Methods**

### 369 **Model**

370 Sequences are represented by a 21-letter alphabet for proteins or 5-letter alphabet for RNAs, one  
371 for each residue type and a ‘start/stop’ character. Training sequences are weighted inversely to  
372 the number of neighbors for each sequence at a minimum identity of 80%, except for viral  
373 families, where a 99% identity threshold was used, as was done previously<sup>30</sup>. Sequence sets are  
374 derived from alignments by extracting full sequences for each aligned region; sequence  
375 identities, boundaries, and weights are the only information provided to the model by alignments.  
376 The log-likelihood for a sequence is the sum of the cross-entropy between the true residue at  
377 each position and the predicted distribution over possible residues, conditioned on the previous  
378 characters. Since we encountered exploding gradients<sup>62</sup> during training on long sequence  
379 families with LSTM<sup>63</sup> or GRU<sup>64</sup> architectures, we parameterize an autoregressive likelihood with  
380 dilated convolutional neural networks (**Supplementary Fig. 1**). These feed-forward deep neural  
381 networks aggregate long-range dependencies in sequences over an exponentially large receptive  
382 field<sup>65-67</sup>. Specifically, we use a residual causal dilated convolutional neural network architecture  
383 with 6 blocks of 9 dilated convolutional layers and both weight normalization<sup>98</sup> and layer  
384 normalization<sup>99</sup>, where the number of blocks and layers were chosen to cover protein sequences  
385 of any length. To help prevent overfitting, we use L2 regularization on the weights and place  
386 Dropout layers ( $p = 0.5$ ) immediately after each of the 6 residual blocks<sup>100</sup>. We use a batch size  
387 of 30 for all sequence families tested. Channel sizes of 24 and 48 were tested for all protein  
388 families, and channel size 48 was chosen for further use. Six models are built for each family:  
389 three replicates in both the N-to-C and C-to-N directions, respectively. Each model is trained for  
390 250,000 updates using Adam with default parameters<sup>101</sup> at which point the loss had visibly  
391 converged, and the gradient norm is clipped<sup>62</sup> to 100.

### 392 **Data collection**

393 40 datasets which include experimental mutation effects, the sequence families, and effect  
394 predictions were taken from our previous publication<sup>30</sup> and 5 datasets that include indels and  
395 nanobody thermostability data were added for this work (references and data in **Supplementary**  
396 **Table 4 and Extended Data**). For new mutation effect predictions such as the indel mutation  
397 scans, sequence families were collected from the UniProt database in the same procedure as  
398 described in previous published work<sup>30</sup>. Pathogenic mutations for the Tau protein were



399 downloaded from the Alzforum database<sup>79</sup>. The naïve llama immune repertoire was acquired  
400 from<sup>57</sup>. Due to the large number of sequences in the llama immune repertoire, sequence weights  
401 were approximated using Linclust<sup>102</sup> by clustering sequences at both 80% and 90% sequence  
402 identity thresholds.

### 403 **Nanobody library generation**

404 Using the N-to-C terminus model trained on llama nanobody sequences, we generated  
405 33,047,639 CDR3 sequences by ancestral sampling<sup>61</sup>, conditioned on the germline framework-  
406 CDR1-CDR2 sequence and continued until generation of the stop character. Duplicates of the  
407 training set or generated sequences and those not matching the final beta strand of our nanobody  
408 template were excluded. CDR3 sequences were also removed if they contained glycosylation  
409 (NxS and NxT) sites, asparagine deamination (NG) motifs, or sulfur-containing amino acids  
410 (cysteine and methionine), resulting in 3,690,554 sequences.

411 From this large number of sequences, we then sought to choose roughly 200,000 CDR3  
412 sequences that are both deemed fit by the model and as diverse from one another as possible to  
413 cover the largest amount of sequence space. First, we featurized these sequences into fixed  
414 length, L2 normalized k-mer vectors with k-mers of size 1, 2, and 3. We then used BIRCH  
415 clustering<sup>103</sup> to find diverse members of the dataset in O(n) time. We used a diameter threshold  
416 of 0.575, resulting in 382,675 clusters. K-mer size and BIRCH diameter threshold were chosen  
417 to maximize the number of clusters within a memory constraint of 70 GB. From the cluster  
418 centroids, we chose the 185,836 most probable sequences for final library construction.

### 419 **Construction of nanobody library**

420 FragmentGENE\_NbCM coding for the nanobody template was amplified with oligonucleotides  
421 NbCM\_pydsF2.0 and NbCM\_pydsR and then cloned into the pYDS649 yeast-display plasmid<sup>8</sup>  
422 using HiFi Mastermix (New England Biolabs). The original NotI site in pYDS649 was then  
423 removed by amplification with primers NotI\_removal\_1F and Pyds\_NbCM\_cloning\_R followed  
424 by cloning again into pYDS649 to generate the pYDS\_NbCM display plasmid for the nanobody  
425 template.

426 An oligonucleotide library was synthesized (Agilent) with the following design ACTCTGT  
427 [CDR3] ATCGT where CDR3 is a sequence for one of the computationally designed clones.  
428 Two-hundred picomoles of the library was PCR amplified over 15 cycles with oligonucleotides



429 Oligo\_library\_F and Oligo\_library\_R using Q5 polymerase (New England Biolabs). Amplified  
430 DNA was PCR purified (Qiagen) and ethanol precipitated in preparation for yeast  
431 transformation.  $4.8 \times 10^8$  BJ5465 (MAT $\alpha$  ura352 trp1 leu2 $\Delta$ 1 his3 $\Delta$ 200 pep4::HIS3 prb1 $\Delta$ 1.6 R  
432 can1 GAL) yeast cells, grown to OD600 1.6, were transformed, using an ECM 830  
433 Electroporator (BTX-Harvard Apparatus), with 2.4  $\mu$ g of NotI digested pYDS\_NbCM vector and  
434 9.9  $\mu$ g of CDR3 library PCR product yielding  $2.7 \times 10^6$  transformants. Library aliquots of  $2.4 \times$   
435  $10^8$  cells per vial were frozen in tryptophan dropout media containing 10% DMSO.

#### 436 **Characterization of nanobody library**

437 Yeast displaying the computationally designed or combinatorial synthetic nanobody library<sup>8</sup>  
438 were grown in tryptophan dropout media with glucose as the sugar source for one day at 30 °C  
439 and then passaged into media with galactose as the sole sugar source to induce expression of  
440 nanobodies at 25 °C. After two days of induction, one million cells from each library were  
441 stained with a 1:25 dilution of anti-HA AlexaFluor647 conjugated antibody (Cell Signaling  
442 Technology) in Buffer A (20 mM HEPES pH 7.5, 150 mM NaCl, 0.1% BSA, 0.2% maltose) for  
443 30 minutes at 4 °C. After staining, cells were centrifuged, the supernatant was removed, and cells  
444 were resuspended in Buffer A for flow analysis with an Accuri C6 (BD Biosciences,

#### 445 **Supplementary Fig. 12).**

446 To find nanobody binders to human serum albumin (HSA) one round of magnetic-activated cell  
447 sorting (MACS) followed by two rounds of fluorescence-activated cell sorting (FACS) were  
448 performed on our yeast-displayed library of nanobodies. For MACS,  $4 \times 10^7$  induced cells were  
449 resuspended in binding buffer (20 mM HEPES pH 7.5, 150 mM NaCl, 0.1% ovalbumin) along  
450 with anti-fluorescein isothiocyanate (FITC) microbeads (Miltenyi) and FITC-labeled streptavidin  
451 for 35 min at 4°C and then passed through an LD column (Miltenyi) to remove binders to  
452 microbeads and streptavidin. Remaining yeast were centrifuged and resuspended in binding  
453 buffer and incubated with 500 nM streptavidin-FITC and 2  $\mu$ M of biotinylated HSA for one hour  
454 at 4°C. Yeast were then centrifuged and resuspended in binding buffer containing anti-FITC  
455 microbeads for 15 min at 4°C before passing them into an LS column and eluting and collecting  
456 the bound yeast. For the first round of FACS, induced yeast were first stained with 1  $\mu$ M of  
457 biotinylated HSA for 45 min at 4°C and then briefly stained with 500 nM of streptavidin tetramer  
458 along with antiHA-488 to assess expression levels. Both yeast stainings were performed in

459 FACS buffer (20 mM HEPES pH 7.5, 150 mM NaCl, 0.1% ovalbumin, 0.2% maltose).  $5 \times 10^6$   
460 yeast were sorted and 28,000 were collected and expanded for a second round of FACS. The  
461 second round of FACS was performed under the same conditions as the first and from  $3.8 \times 10^6$   
462 sorted yeast 21,455 were collected. Nanobody Nb174684 was isolated from a screen of 36 clones  
463 for binding to HSA using a flow cytometer and then sequenced. In order to characterize binding  
464 of Nb174684, yeast displaying Nb174684 were stained with varying amounts of AlexaFluor 488  
465 labeled HSA and fluorescence was analyzed with a flow cytometer.

466 Oligonucleotides:

467 FragmentGENE\_NbCM:

468 GCTGCCAGCCGGCGATGGCCAGGTCCAACCTCAAGAATCAGGCGGGGGCCTGGT  
469 ACAGGCAGGCGGTTCTCTTCGGCTGTCGTGTGCGGCAAGCGGATTTACATTCAGTAG  
470 CTACGCTATGGGCTGGTACCGTCAGGCACCGGGGAAAGAACGGGAATTTGTTGCTG  
471 CAATCTCTTGAGCGGTGGGAGCACATATTATGCAGATTCCGTTAAAGGCAGATTCA  
472 CGATCAGTCGCGATAACGCAAAAAATACAGTGTACTTACAAATGAACTCTTTGAAA  
473 CCCGAAGACACCGCAGTCTATTACTGCGCGGCCGCTACTGGGGACAAGGCACCCAG  
474 GTGACTGTATCATCCCACCACCACCACCACCTGA

475 NbCM\_pydsF2.0:

476 GGTGTTCAATTGGACAAGAGAGAAGCTGACGCAGAAGTCCAACCTTGTCGAATCAGG  
477 CGGGGGCCTGGTACAG

478 NbCM\_pydsR:

479 CGTAATCTGGAACATCGTATGGGTAGGATCCGGATGATACAGTCACCTGGGT

480 NotI\_removal\_1F:

481 CAACCCTCACTAAAGGGCGTTCGCCATGAGATTCCCATCTATCTTCA

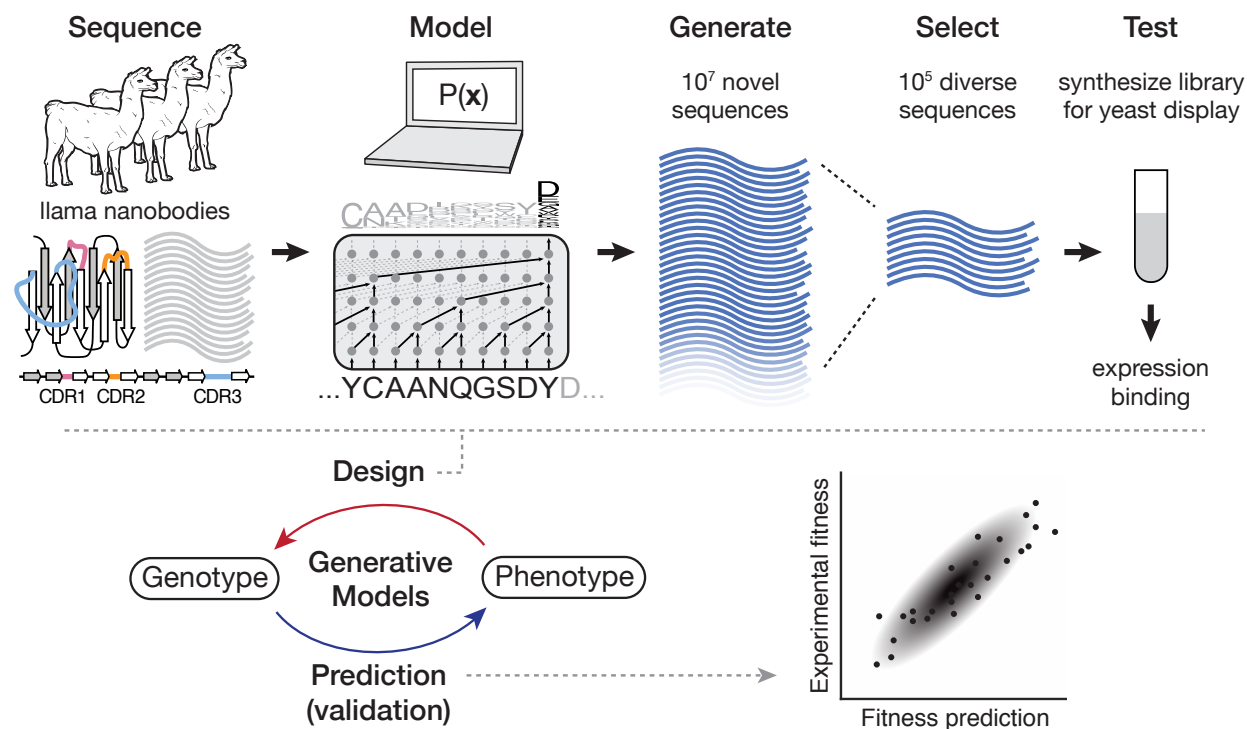
482 Pyds\_NbCM\_cloning\_R:

483 CACCTGGGTGCCTTGTCCCCAGTA

484

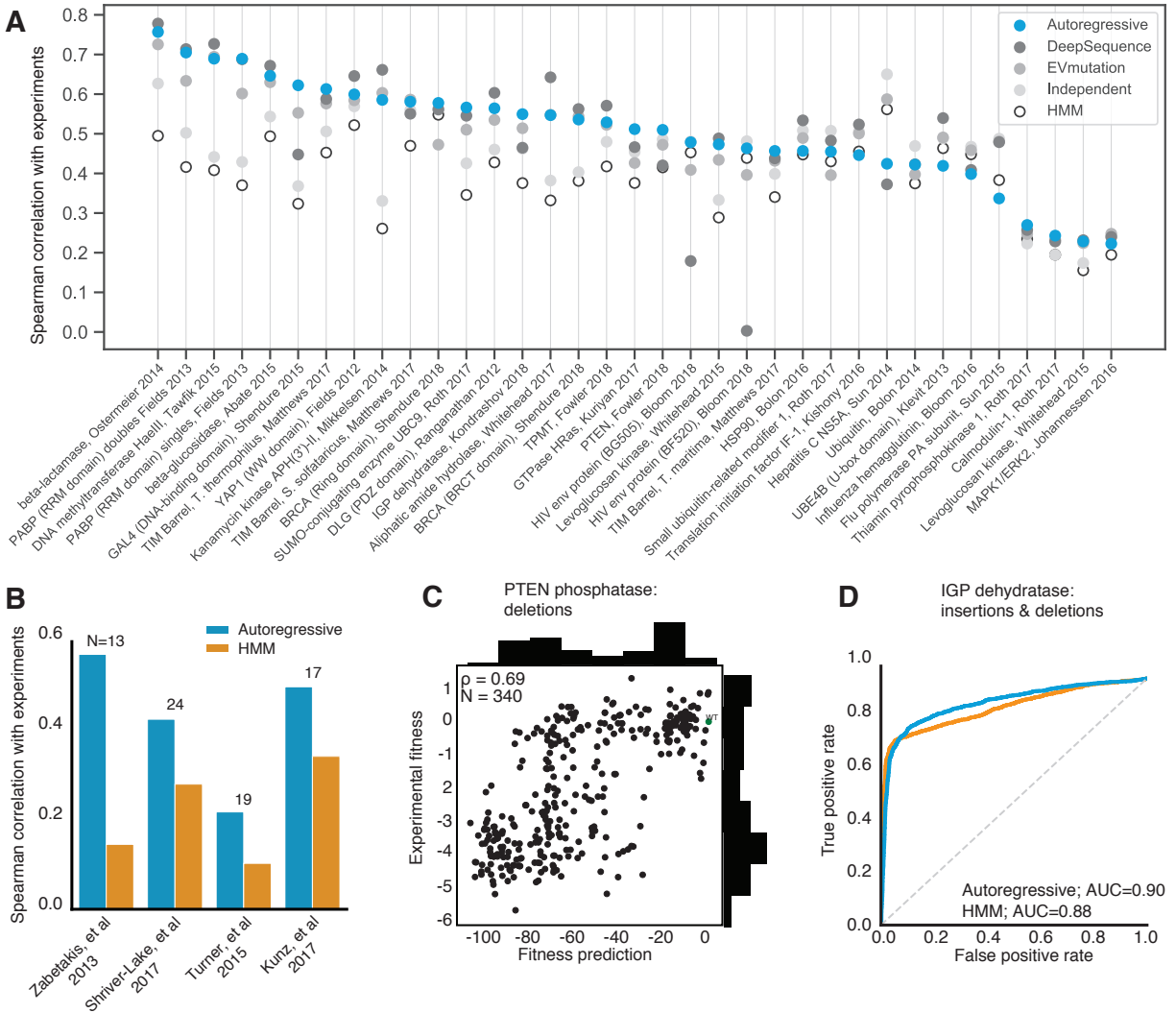
485  
486

## Figures



487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497

**Figure 1.** Autoregressive models of biological sequences can learn the genotype-phenotype map for both prediction and design. From natural sequences in a naïve llama repertoire<sup>57</sup>, the autoregressive model can learn functional constraints by predicting the likelihood of each residue in the sequence conditioned on preceding residues. We then use these constraints to generate millions of novel nanobody sequences—as many can be generated as desired. Of these designed sequences we select hundreds of thousands of diverse sequences, synthesize a library, and screen for expression and binding. We also validate the model on mutation effect prediction tasks of deep mutational scans including the effects of multiple insertions and deletions, and the thermostabilities of highly variable nanobody sequences.



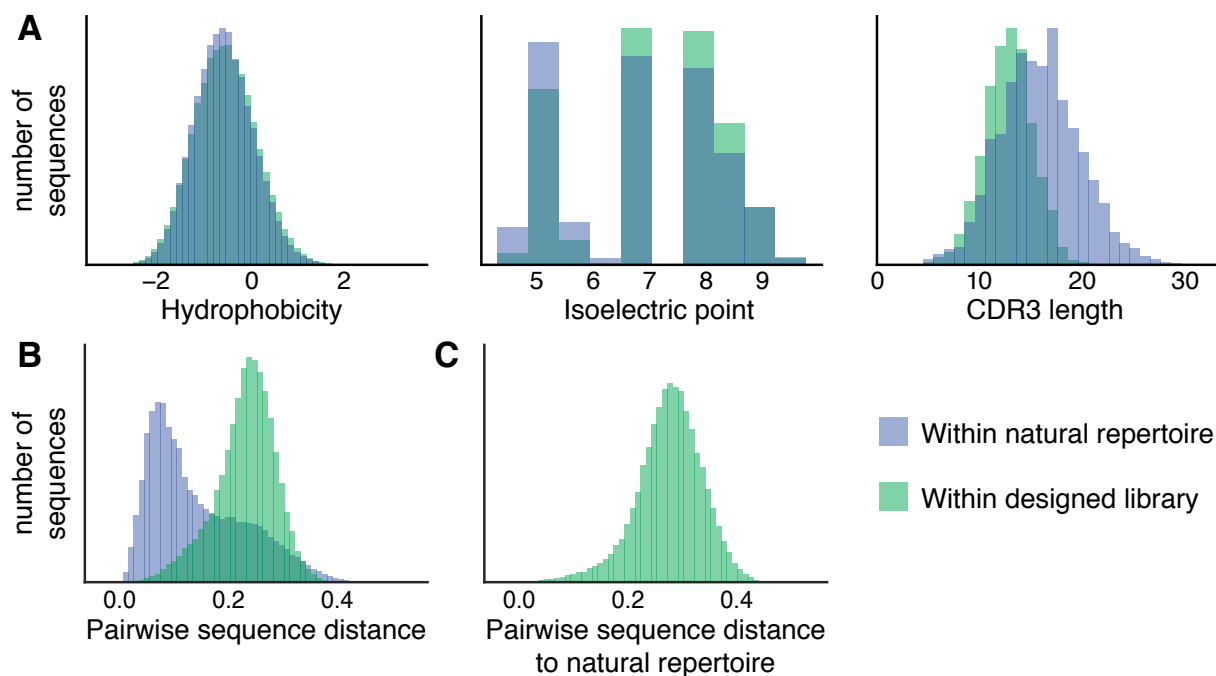
498

499

**Figure 2.** Validation of the autoregressive model in learning the genotype to phenotype map. The model accurately predicts fitness of biological sequences of various lengths. **a.** Even without using alignments, the autoregressive model can competitively match mutation effect prediction accuracies of state-of-art alignment-dependent models, such as conservation, evolutionary couplings, and DeepSequence. Additionally, the mutation effect prediction accuracies improves upon HMM model accuracies. Without using alignments, the autoregressive model matches alignment-dependent state-of-art missense mutation effect prediction (DeepSequence<sup>30</sup>) for 40 different deep mutational scan experiments. Three datasets show significant improvement with the autoregressive model: HIV env (BF520), HIV env (BG505), and GAL4 DNA-binding domain. **b.** The autoregressive model can learn from natural sequence repertoires of llama nanobodies to predict the thermostability of llama nanobody sequences with variation in the framework and complementarity determining regions with greater accuracy than hidden Markov models<sup>74</sup>. The number of llama nanobody sequences from each study is shown above each pair of bars. **c.** Fitness predictions for single deletions in PTEN phosphatase compared with measured

512

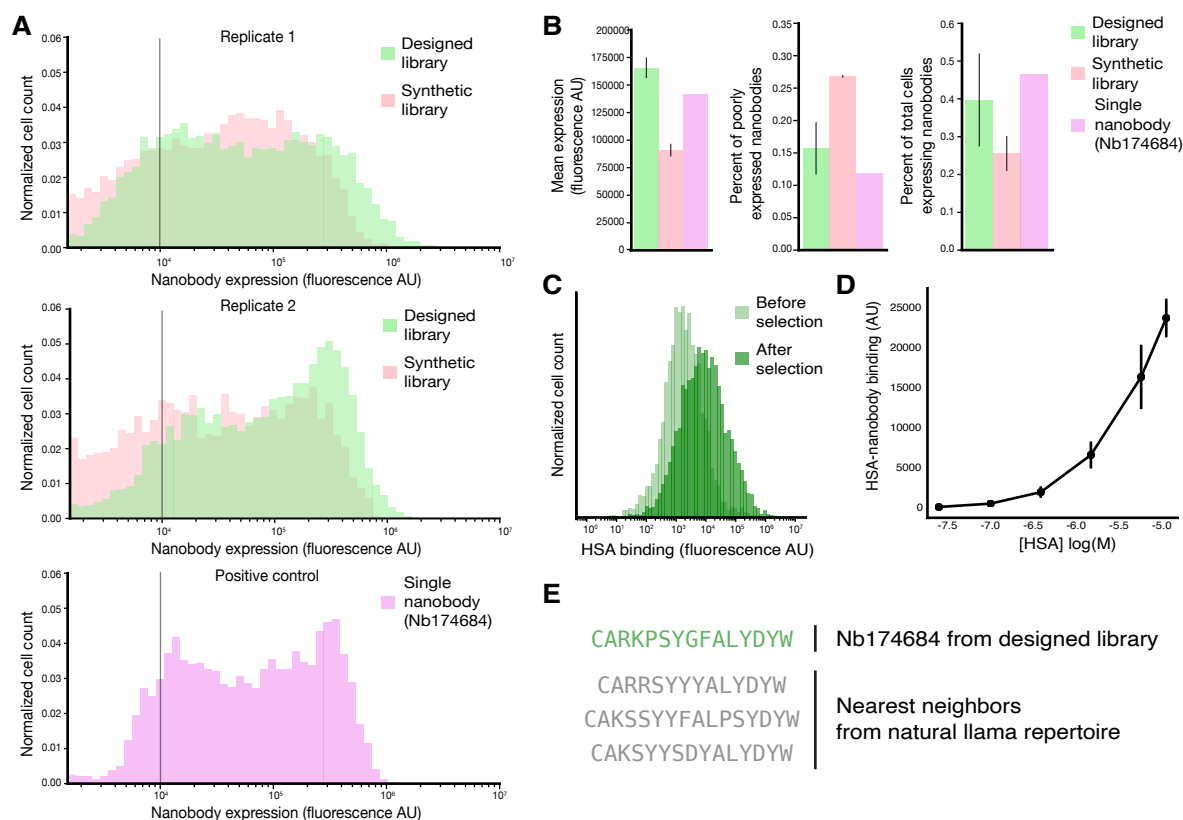
513 experimental fitness is accurate, with a Spearman correlation of 0.69. **d.** Accurate prediction of  
514 binary fitness for IGP dehydratase with a range of insertions, deletions, and missense mutations.  
515



516  
517

518 **Figure 3.** The designed library has comparable biochemical property distributions and improved  
519 diversity to the natural llama repertoire. **a.** Conditioned on the framework-CDR1-CDR2  
520 sequence, a diverse set of CDR3 sequences are generated and selected. These CDR3 sequences  
521 are similar to the natural repertoire in their distributions of hydrophobicity<sup>104</sup> and isoelectric  
522 point<sup>105, 106</sup>, while having shorter length distributions due to selection strategies in the final  
523 library construction. **b.** The designed library contains more diversity in sequences than the  
524 natural repertoire as evidenced by the larger cosine distance to its nearest neighbor. **c.** Each  
525 sequence in the designed library is diverse from any sequence seen in the natural repertoire,  
526 indicating that we have learned fit sequence constraints but are traversing previously unexplored  
527 regions of sequence space.

528



529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

**Figure 4.** The designed library contains stable and functional nanobody sequences that are well expressed and can bind target antigens. **a.** Fluorescence distributions of cells expressing nanobodies comparing the synthetic combinatorial library and our designed library in two biological replicate experiments as well as a control experiment of a single, well-expressed nanobody clone (Nb174684). The distributions of the designed library are consistently right-shifted compared to the combinatorial library and resemble the control nanobody. **b.** Compared to the combinatorial library, the designed library has almost double the mean expression level (left panel, 166,193 AU compared to 92,183 AU), nearly half the fraction of poorly expressed nanobodies (of cells expressing nanobodies) (middle panel, 15.4% compared to 25.7% of clones with less than 10,000 AU indicated as a grey bar in panel **a**), and one and a half times the fraction of total cells that express nanobodies (right panel, 39.6% compared to 25.1%). The thresholds for determining the proportion of total cells expressing nanobodies were found by identifying the local minima on the distributions and are displayed in **Supplementary Fig. 10**. Values displayed on the bar graphs are means of the two replicates and the standard deviations are shown as error bars. There is only one replicate for the control experiment of the single nanobody clone. **c.** Fluorescence distributions of nanobodies bound to HSA shows a rightward shift after screening and selection, indicating a successful enrichment of binders to the target antigen. **d.** On-yeast binding assay of Nb.174684, an HSA binder identified from the designed library with moderate binding affinity. Error bars represent standard deviations in measurements at each concentration of HSA. **e.** CDR3 sequence of binder Nb.174684 and the sequences of the



551 nearest neighbors from the natural llama repertoire that was used to train the autoregressive  
552 model.

553

## References and Notes:

- 554 1. Romero, P.A. & Arnold, F.H. Exploring protein fitness landscapes by directed evolution.  
555 *Nature Reviews Molecular Cell Biology* **10**, 866-876 (2009).
- 556 2. Dougherty, M.J. & Arnold, F.H. Directed evolution: new parts and optimized function.  
557 *Current Opinion in Biotechnology* **20**, 486-491 (2009).
- 558 3. Baker, D. An exciting but challenging road ahead for computational enzyme design.  
559 *Protein Science* **19**, 1817-1819 (2010).
- 560 4. Huang, P.-S., Boyken, S.E. & Baker, D. The coming of age of de novo protein design.  
561 *Nature* **537**, 320-327 (2016).
- 562 5. Fowler, D.M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nat*  
563 *Methods* **11**, 801-807 (2014).
- 564 6. Muyldermans, S. Nanobodies: natural single-domain antibodies. *Annu Rev Biochem* **82**,  
565 775-797 (2013).
- 566 7. Sall, A. et al. Generation and analyses of human synthetic antibody libraries and their  
567 application for protein microarrays. *Protein Eng Des Sel* **29**, 427-437 (2016).
- 568 8. McMahan, C. et al. Yeast surface display platform for rapid discovery of  
569 conformationally selective nanobodies. *Nat Struct Mol Biol* **25**, 289-296 (2018).
- 570 9. Bradbury, A.R., Sidhu, S., Dubel, S. & McCafferty, J. Beyond natural antibodies: the  
571 power of in vitro display technologies. *Nat Biotechnol* **29**, 245-254 (2011).
- 572 10. Schoof, M. et al. An ultra-potent synthetic nanobody neutralizes SARS-CoV-2 by locking  
573 Spike into an inactive conformation. *bioRxiv*, 2020.2008.2008.238469 (2020).
- 574 11. Miho, E., Roskar, R., Greiff, V. & Reddy, S.T. Large-scale network analysis reveals the  
575 sequence space architecture of antibody repertoires. *Nat Commun* **10**, 1321 (2019).
- 576 12. Jain, T. et al. Biophysical properties of the clinical-stage antibody landscape. *Proc Natl*  
577 *Acad Sci U S A* **114**, 944-949 (2017).
- 578 13. Marks, C. & Deane, C.M. How repertoire data are changing antibody science. *J Biol*  
579 *Chem* **295**, 9823-9837 (2020).
- 580 14. Asti, L., Uguzzoni, G., Marcatili, P. & Pagnani, A. Maximum-Entropy Models of  
581 Sequenced Immune Repertoires Predict Antigen-Antibody Affinity. *PLoS Comput Biol*  
582 **12**, e1004870 (2016).
- 583 15. Mora, T., Walczak, A.M., Bialek, W. & Callan, C.G., Jr. Maximum entropy models for  
584 antibody diversity. *Proc Natl Acad Sci U S A* **107**, 5405-5410 (2010).
- 585 16. Marcou, Q., Mora, T. & Walczak, A.M. High-throughput immune repertoire analysis  
586 with IGoR. *Nat Commun* **9**, 561 (2018).
- 587 17. Liu, G. et al. Antibody complementarity determining region design using high-capacity  
588 machine learning. *Bioinformatics* **36**, 2126-2133 (2020).
- 589 18. DeKosky, B.J. et al. Large-scale sequence and structural comparisons of human naive  
590 and antigen-experienced antibody repertoires. *Proc Natl Acad Sci U S A* **113**, E2636-  
591 2645 (2016).
- 592 19. Muyldermans, S. A guide to: generation and design of nanobodies. *FEBS J* (2020).
- 593 20. Zimmermann, I. et al. Synthetic single domain antibodies for the conformational trapping  
594 of membrane proteins. *Elife* **7** (2018).
- 595 21. Marks, D.S. et al. Protein 3D structure computed from evolutionary sequence variation.  
596 *PLoS One* **6**, e28766 (2011).
- 597 22. Hopf, T.A. et al. Three-dimensional structures of membrane proteins from genomic  
598 sequencing. *Cell* **149**, 1607-1621 (2012).

- 599 23. Marks, D.S., Hopf, T.A. & Sander, C. Protein structure prediction from sequence  
600 variation. *Nat Biotechnol* **30**, 1072-1080 (2012).
- 601 24. Kamisetty, H., Ovchinnikov, S. & Baker, D. Assessing the utility of coevolution-based  
602 residue-residue contact predictions in a sequence- and structure-rich era. *Proc Natl Acad*  
603 *Sci U S A* **110**, 15674-15679 (2013).
- 604 25. Ovchinnikov, S., Kamisetty, H. & Baker, D. Robust and accurate prediction of residue-  
605 residue interactions across protein interfaces using evolutionary information. *Elife* **3**,  
606 e02030 (2014).
- 607 26. Hopf, T.A. et al. Sequence co-evolution gives 3D contacts and structures of protein  
608 complexes. *eLife* **3**, e03430 (2014).
- 609 27. Cong, Q., Anishchenko, I., Ovchinnikov, S. & Baker, D. Protein interaction networks  
610 revealed by proteome coevolution. *Science* **365**, 185 (2019).
- 611 28. Green, A.G. et al. Proteome-scale discovery of protein interactions with residue-level  
612 resolution using sequence coevolution. *bioRxiv*, 791293 (2019).
- 613 29. Hopf, T.A. et al. Mutation effects predicted from sequence co-variation. *Nat Biotechnol*  
614 **35**, 128-135 (2017).
- 615 30. Riesselman, A.J., Ingraham, J.B. & Marks, D.S. Deep generative models of genetic  
616 variation capture the effects of mutations. *Nat Methods* **15**, 816-822 (2018).
- 617 31. Gray, V.E., Hause, R.J., Luebeck, J., Shendure, J. & Fowler, D.M. Quantitative Missense  
618 Variant Effect Prediction Using Large-Scale Mutagenesis Data. *Cell Syst* **6**, 116-124  
619 e113 (2018).
- 620 32. Mann, J.K. et al. The fitness landscape of HIV-1 gag: advanced modeling approaches and  
621 validation of model predictions by in vitro testing. *PLoS Comput Biol* **10**, e1003776  
622 (2014).
- 623 33. Figliuzzi, M., Jacquier, H., Schug, A., Tenailon, O. & Weigt, M. Coevolutionary  
624 Landscape Inference and the Context-Dependence of Mutations in Beta-Lactamase TEM-  
625 1. *Molecular Biology and Evolution* **33**, 268-280 (2015).
- 626 34. Livesey, B.J. & Marsh, J.A. Using deep mutational scanning to benchmark variant effect  
627 predictors and identify disease mutations. *Mol Syst Biol* **16**, e9380 (2020).
- 628 35. Sim, N.L. et al. SIFT web server: predicting effects of amino acid substitutions on  
629 proteins. *Nucleic Acids Res* **40**, W452-457 (2012).
- 630 36. Adzhubei, I., Jordan, D.M. & Sunyaev, S.R. Predicting functional effect of human  
631 missense mutations using PolyPhen-2. *Curr Protoc Hum Genet* **Chapter 7**, Unit7 20  
632 (2013).
- 633 37. Kircher, M. et al. A general framework for estimating the relative pathogenicity of human  
634 genetic variants. *Nat Genet* **46**, 310-315 (2014).
- 635 38. Reva, B., Antipin, Y. & Sander, C. Predicting the functional impact of protein mutations:  
636 application to cancer genomics. *Nucleic Acids Res* **39**, e118 (2011).
- 637 39. Lefranc, M.P. et al. IMGT unique numbering for immunoglobulin and T cell receptor  
638 variable domains and Ig superfamily V-like domains. *Dev Comp Immunol* **27**, 55-77  
639 (2003).
- 640 40. van der Lee, R. et al. Classification of intrinsically disordered regions and proteins. *Chem*  
641 *Rev* **114**, 6589-6631 (2014).
- 642 41. Mullaney, J.M., Mills, R.E., Pittard, W.S. & Devine, S.E. Small insertions and deletions  
643 (INDELS) in human genomes. *Hum Mol Genet* **19**, R131-136 (2010).

- 644 42. Lin, M. et al. Effects of short indels on protein structure and function in human genomes.  
645 *Sci Rep* **7**, 9313 (2017).
- 646 43. Mills, R.E. et al. Natural genetic variation caused by small insertions and deletions in the  
647 human genome. *Genome Res* **21**, 830-839 (2011).
- 648 44. Pentony, M.M. & Jones, D.T. Modularity of intrinsic disorder in the human proteome.  
649 *Proteins* **78**, 212-221 (2010).
- 650 45. Campbell, P.J. et al. Pan-cancer analysis of whole genomes. *Nature* **578**, 82-93 (2020).
- 651 46. Turajlic, S. et al. Insertion-and-deletion-derived tumour-specific neoantigens and the  
652 immunogenic phenotype: a pan-cancer analysis. *The Lancet Oncology* **18**, 1009-1021  
653 (2017).
- 654 47. Deiana, A., Forcelloni, S., Porrello, A. & Giansanti, A. Intrinsically disordered proteins  
655 and structured proteins with intrinsically disordered regions have different functional  
656 roles in the cell. *PLoS One* **14**, e0217889 (2019).
- 657 48. Uversky, V.N. et al. Unfoldomics of human diseases: linking protein intrinsic disorder  
658 with diseases. *BMC Genomics* **10 Suppl 1**, S7 (2009).
- 659 49. Graves, A., Mohamed, A. & Hinton, G. Speech recognition with deep recurrent neural  
660 networks. *2013 IEEE International Conference on Acoustics, Speech and Singal*  
661 *Processing*, 6645-6649 (2013).
- 662 50. Wang, Y. et al. Tacotron: Towards end-to-end speech synthesis. *arXiv* (2017).
- 663 51. Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly learning to  
664 align and translate. *arXiv* (2014).
- 665 52. Sutskever, I., Vinyals, O. & Le, Q.V. Sequence to sequence learning with neural  
666 networks. *Advances in neural information processing systems*, 3104-3112 (2014).
- 667 53. Alley, E.C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G.M. Unified rational  
668 protein engineering with sequence-based deep representation learning. *Nature Methods*  
669 **16**, 1315-1322 (2019).
- 670 54. Linder, J., Bogard, N., Rosenberg, A.B. & Seelig, G. A Generative Neural Network for  
671 Maximizing Fitness and Diversity of Synthetic DNA and Protein Sequences. *Cell Syst* **11**,  
672 49-62 e16 (2020).
- 673 55. Strodthoff, N., Wagner, P., Wenzel, M. & Samek, W. UDSMProt: universal deep  
674 sequence models for protein classification. *Bioinformatics* **36**, 2401-2409 (2020).
- 675 56. Brookes, D.H., Park, H. & Listgarten, J. Conditioning by adaptive sampling for robust  
676 design. *Proceedings of the 36th International Conference on Machine Learning* **97**, 773-  
677 782 (2019).
- 678 57. McCoy, L.E. et al. Molecular evolution of broadly neutralizing Llama antibodies to the  
679 CD4-binding site of HIV-1. *PLoS Pathog* **10**, e1004552 (2014).
- 680 58. Chan, S.K., Rahumatullah, A., Lai, J.Y. & Lim, T.S. Naive Human Antibody Libraries  
681 for Infectious Diseases. *Adv Exp Med Biol* **1053**, 35-59 (2017).
- 682 59. Ravikumar, A., Arzumanyan, G.A., Obadi, M.K.A., Javanpour, A.A. & Liu, C.C.  
683 Scalable, Continuous Evolution of Genes at Mutation Rates above Genomic Error  
684 Thresholds. *Cell* **175**, 1946-1957 e1913 (2018).
- 685 60. Vaswani, A. et al. Attention is all you need. *Advances in neural information processing*  
686 *systems*, 5998-6008 (2017).
- 687 61. Sutskever, I., Martens, J. & Hinton, G. Generating text with recurrent neural networks.  
688 *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*,  
689 1017-1024 (2011).

- 690 62. Pascanu, R., Mikolov, T. & Begio, Y. On the difficulty of training recurrent neural  
691 networks. *International Conference on Machine Learning*, 1310-1318 (2013).
- 692 63. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput* **9**, 1735-1780  
693 (1997).
- 694 64. Cho, K. et al. Learning phrase representations using RNN Encoder-Decoder for statistical  
695 machine translation. *Proceedings of the 2014 Conference on Empirical Methods in  
696 Natural Language Processing (EMNLP)*, 1724-1734 (2014).
- 697 65. van den Oord, A. et al. Wavenet: A generative model for raw audio. *arXiv* (2016).
- 698 66. Kalchbrenner, N. et al. Neural machine translation in linear time. *arXiv* (2016).
- 699 67. Gupta, A. & Rush, A. Dilated convolutions for modeling long-distance genomic  
700 dependencies. *arXiv* (2017).
- 701 68. Grimm, D.G. et al. The evaluation of tools used to predict the impact of missense variants  
702 is hindered by two types of circularity. *Hum Mutat* **36**, 513-523 (2015).
- 703 69. Kunz, P. et al. Exploiting sequence and stability information for directing nanobody  
704 stability engineering. *Biochim Biophys Acta Gen Subj* **1861**, 2196-2205 (2017).
- 705 70. Shriver-Lake, L.C., Zabetakis, D., Goldman, E.R. & Anderson, G.P. Evaluation of anti-  
706 botulinum neurotoxin single domain antibodies with additional optimization for  
707 improved production and stability. *Toxicon* **135**, 51-58 (2017).
- 708 71. Turner, K.B. et al. Improving the biophysical properties of anti-ricin single-domain  
709 antibodies. *Biotechnol Rep (Amst)* **6**, 27-35 (2015).
- 710 72. Zabetakis, D., Anderson, G.P., Bayya, N. & Goldman, E.R. Contributions of the  
711 complementarity determining regions to the thermal stability of a single-domain  
712 antibody. *PLoS One* **8**, e77678 (2013).
- 713 73. Durbin, R., Eddy, S.R., Krogh, A. & Mitchison, G. Biological sequence analysis:  
714 probabilistic models of proteins and nucleic acids. (Cambridge university press, 1998).
- 715 74. Eddy, S.R. Accelerated profile HMM searches. *PLoS computational biology* **7**, e1002195  
716 (2011).
- 717 75. Mighell, T.L., Evans-Dutson, S. & O'Roak, B.J. A Saturation Mutagenesis Approach to  
718 Understanding PTEN Lipid Phosphatase Activity and Genotype-Phenotype  
719 Relationships. *Am J Hum Genet* **102**, 943-955 (2018).
- 720 76. Pokusaeva, V.O. et al. An experimental assay of the interactions of amino acids from  
721 orthologous sequences shaping a complex fitness landscape. *PLoS Genet* **15**, e1008079  
722 (2019).
- 723 77. Choi, Y., Sims, G.E., Murphy, S., Miller, J.R. & Chan, A.P. Predicting the functional  
724 effect of amino acid substitutions and indels. *PLoS One* **7**, e46688 (2012).
- 725 78. Landrum, M.J. et al. ClinVar: improving access to variant interpretations and supporting  
726 evidence. *Nucleic Acids Res* **46**, D1062-D1067 (2018).
- 727 79. MAPT | Alzforum. Retrieved August 12, 2020. from  
728 <https://www.alzforum.org/mutations/mapt>.
- 729 80. Harmsen, M.M. et al. Llama heavy-chain V regions consist of at least four distinct  
730 subfamilies revealing novel sequence features. *Mol Immunol* **37**, 579-590 (2000).
- 731 81. Beerli, R.R. & Rader, C. Mining human antibody repertoires. *MAbs* **2**, 365-378 (2010).
- 732 82. Dimitrov, J.D., Pashov, A.D. & Vassilev, T.L. Antibody polyspecificity: what does it  
733 matter? *Adv Exp Med Biol* **750**, 213-226 (2012).
- 734 83. Dimitrov, J.D. et al. Antibody polyreactivity in health and disease: statu variabilis. *J*  
735 *Immunol* **191**, 993-999 (2013).



- 736 84. Kelly, R.L., Zhao, J., Le, D. & Wittrup, K.D. Nonspecificity in a nonimmune human  
737 scFv repertoire. *MAbs* **9**, 1029-1035 (2017).
- 738 85. Lim, C.C., Choong, Y.S. & Lim, T.S. Cognizance of Molecular Methods for the  
739 Generation of Mutagenic Phage Display Antibody Libraries for Affinity Maturation. *Int J*  
740 *Mol Sci* **20** (2019).
- 741 86. Pashova, S., Schneider, C., von Gunten, S. & Pashov, A. Antibody repertoire profiling  
742 with mimotope arrays. *Hum Vaccin Immunother* **13**, 314-322 (2017).
- 743 87. Wardemann, H. et al. Predominant autoantibody production by early human B cell  
744 precursors. *Science* **301**, 1374-1377 (2003).
- 745 88. Rives, A. et al. Biological structure and function emerge from scaling unsupervised  
746 learning to 250 million protein sequences. *bioRxiv* (2019).
- 747 89. Rao, R. et al. Evaluating Protein Transfer Learning with TAPE. *33rd Conference on*  
748 *Neural Information Processing Systems* (2019).
- 749 90. Wellner, A. et al. Rapid generation of potent antibodies by autonomous hypermutation in  
750 yeast. *bioRxiv* (2020).
- 751 91. DeKosky, B.J. et al. High-throughput sequencing of the paired human immunoglobulin  
752 heavy and light chain repertoire. *Nat Biotechnol* **31**, 166-169 (2013).
- 753 92. Ingraham, J.B., Vikas, G.K., Barzilay, R. & Jaakkola, T. Generative models for graph-  
754 based protein design. *33rd Conference on Neural Information Processing Systems*  
755 15794-15805 (2019).
- 756 93. Kim, Y., Wiseman, S., Miller, A.C., Sontag, D. & Rush, A. Semi-amortized variational  
757 autoencoders. *arXiv* (2018).
- 758 94. Yang, Z., Hu, Z., Salakhutdinov, R. & Berg-Kirkpatrick, T. Improved variational  
759 autoencoders for text modeling using dilated convolutions. *arXiv* (2017).
- 760 95. van den Oord, A. & Vinyals, O. Neural discrete representation learning. *Advances in*  
761 *neural information processing systems*, 6306-6315 (2017).
- 762 96. Greener, J.G., Moffat, L. & Jones, D.T. Design of metalloproteins and novel protein folds  
763 using variational autoencoders. *Sci Rep* **8**, 16189 (2018).
- 764 97. Vijayakumar, A.K. et al. Diverse beam search: Decoding diverse solutions from neural  
765 sequence models. *arXiv* (2016).
- 766 98. Salimans, T. & Kingma, D.P. Weight normalization: a simple reparametrization to  
767 accelerate training of deep neural networks. *Advances in neural information processing*  
768 *systems*, 901-909 (2016).
- 769 99. Ba, J.L., Kiros, J.R. & Hinton, G. Layer normalization. *arXiv* (2016).
- 770 100. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a  
771 simple way to prevent neural networks from overfitting. *The Journal of Machine*  
772 *Learning Research* **15**, 1929-1958 (2014).
- 773 101. Kingma, D.P. & Ba, J.L. Adam: a method for stochastic optimization. *arXiv* (2014).
- 774 102. Steinegger, M. & Soding, J. Clustering huge protein sequence sets in linear time. *Nat*  
775 *Commun* **9**, 2542 (2018).
- 776 103. Zhang, T., Ramakrishnan, R. & Livny, M. BIRCH: an efficient data clustering method  
777 for very large databases. *ACM SIGMOD Record* **25**, 103-114 (1996).
- 778 104. Kyte, J. & Doolittle, R.F. A simple method for displaying the hydropathic character of a  
779 protein. *Journal of Molecular Biology* **157**, 105-132 (1982).

- 780 105. Bjellqvist, B. et al. The focusing positions of polypeptides in immobilized pH gradients  
781 can be predicted from their amino acid sequences. *ELECTROPHORESIS* **14**, 1023-1031  
782 (1993).
- 783 106. Cock, P.J.A. et al. Biopython: freely available Python tools for computational molecular  
784 biology and bioinformatics. *Bioinformatics* **25**, 1422-1423 (2009).  
785



786 **Acknowledgments:** We would like to thank John Ingraham, members of the Marks and Sander  
787 labs, and Harvard Research Computing for their insight and feedback to our research. **Author**  
788 **contributions:** D.S.M., A.C.K., and A.J.R. conceived the project; A.J.R. constructed the model;  
789 A.J.R, J.-E.S., and A.W.K. designed and evaluated computational experiments for validation,  
790 prediction, and generation of sequences; A.M. compiled natural nanobody sequence data; C.M.  
791 constructed the library and performed experiments; J.-E.S., A.W.K., and C.M. analyzed the  
792 library experimental data; A.J.R., J.-E.S., A.W.K., C.M., A.C.K., and D.S.M. wrote the  
793 manuscript. **Competing interests:** Authors declare no competing interests. **Data and code**  
794 **availability:** All data generated and analyzed during the study are available in this published  
795 article, its supplementary information files and on the github repository  
796 (<https://github.com/debbiemarkslab/SeqDesign>). All code used for model training, mutation  
797 effect prediction, sequence generation, and library generation is also available on the github  
798 repository.

799

#### 800 **Supplementary Information:**

801 Supplementary Figures 1-12

802 Supplementary Tables 1-5

803 Extended Data 1-6

804