

# UCSF

## UC San Francisco Previously Published Works

### Title

Protein design and variant prediction using autoregressive generative models.

### Permalink

<https://escholarship.org/uc/item/3gt7q4tp>

### Journal

Nature communications, 12(1)

### ISSN

2041-1723

### Authors

Shin, Jung-Eun  
Riesselman, Adam J  
Kollasch, Aaron W  
et al.

### Publication Date





2021-04-01

### DOI

10.1038/s41467-021-22732-w

Peer reviewed

# Protein design and variant prediction using autoregressive generative models

Jung-Eun Shin<sup>1,12</sup>, Adam J. Riesselman<sup>1,9,12</sup>, Aaron W. Kollasch<sup>1,12</sup> , Conor McMahon<sup>2,10</sup>, Elana Simon<sup>3,11</sup>, Chris Sander<sup>4,5</sup>, Aashish Manglik<sup>6,7</sup> , Andrew C. Kruse<sup>2,13</sup>  & Debora S. Marks<sup>1,8,13</sup> 

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

<sup>1</sup>Department of Systems Biology, Harvard Medical School, Boston, MA, USA. <sup>2</sup>Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, Boston, MA, USA. <sup>3</sup>Harvard College, Cambridge, MA, USA. <sup>4</sup>Department of Cell Biology, Harvard Medical School, Boston, MA, USA. <sup>5</sup>Department of Data Sciences, Dana-Farber Cancer Institute, Boston, MA, USA. <sup>6</sup>Department of Pharmaceutical Chemistry, University of California San Francisco, San Francisco, CA, USA. <sup>7</sup>Department of Anesthesia and Perioperative Care, University of California San Francisco, San Francisco, CA, USA. <sup>8</sup>Broad Institute of Harvard and MIT, Cambridge, MA, USA. <sup>9</sup>Present address: insitro, South San Francisco, CA, USA. <sup>10</sup>Present address: Vertex Pharmaceuticals, Boston, MA, USA. <sup>11</sup>Present address: Reverie Labs, Cambridge, MA, USA. <sup>12</sup>These authors contributed equally: Jung-Eun Shin, Adam J. Riesselman, Aaron W. Kollasch. <sup>13</sup>These authors jointly supervised this work: Andrew C. Kruse, Debora S. Marks. ✉email: [Andrew\\_Kruse@hms.harvard.edu](mailto:Andrew_Kruse@hms.harvard.edu); [Debora\\_Marks@hms.harvard.edu](mailto:Debora_Marks@hms.harvard.edu)

Over the past 20 years, success in protein engineering has emerged from two distinct approaches, directed evolution<sup>1,2</sup> and knowledge-based force-field modeling<sup>3,4</sup>. Designing and generating biomolecules with known functions is now a major goal of biotechnology and biomedicine, propelled by our ability to synthesize and sequence DNA at increasingly low costs. However, since the space of possible protein sequences is so large (for a protein of length 100 this is  $10^{130}$ ), deep mutational scans<sup>5</sup> and even very large libraries (e.g.,  $>10^{10}$  variants) barely scratch the surface of the possibilities. As the vast majority of possible sequences will be non-functional proteins, it is crucial to minimize or eliminate these sequences from libraries. Therefore, the open challenge is to develop computational methods that can accelerate this search and bias the search space for protein sequences that are likely to be functional. This will enable the design of libraries for tractable high-throughput experiments that are optimized for functional sequences and variants that are distant in sequence.

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

Instead of experimentally producing unnecessarily massive, largely non-functional libraries, we can design smart libraries of fit and diverse nanobodies for the development of highly specific and possibly therapeutic nanobodies. One way to approach this is to leverage the information in natural sequences to learn constraints on specific amino acids in individual positions in a way that captures their dependency on amino acids in other positions. The sequences of these variants contain rich information about what contributes to a stable, functional protein, and in recent years generative models of these natural protein sequences have been powerful tools for the prediction of the first 3D fold from sequences alone<sup>21,22</sup>, to generally more 3D structures and conformational plasticity<sup>23,24</sup>, protein interactions<sup>25–28</sup>, and most recently, mutation effects<sup>29–34</sup>. However, these state-of-art methods and established methods<sup>35–38</sup> rely on sequence families and alignments, and alignment-based methods are inherently unsuitable for the statistical description of the variable length, hypermutated complementarity determining regions (CDRs) of antibody sequences, which encode the diverse specificity of binding to antigens. While antibody numbering schemes such as IMGT provide consistent alignments of framework residues, alignments of the CDRs rely on symmetrical deletions<sup>39</sup>. Alignment-based models are also unreliable for low-complexity or disordered proteins<sup>40</sup> and cannot handle variants that are

insertions and deletions. Indels make up 15–21% of human polymorphisms<sup>41–43</sup>, 44% of human proteins contain disordered regions longer than 30 amino acids<sup>40,44</sup>, and both are enriched in association with human diseases such as cystic fibrosis, many cancers<sup>45,46</sup>, cardiovascular and neurodegenerative diseases, and diabetes<sup>47,48</sup>.

By contrast, the deep models that have transformed our ability to generate realistic speech such as text-to-speech<sup>49,50</sup> and translation<sup>51,52</sup> use generative models that do not require “word alignment”, e.g., between equisemantic sentences, but instead employ an autoregressive likelihood to tackle context-dependent language prediction and generation. Using this process, an audio clip is decomposed into discrete time steps, a sentence into words, and a protein sequence into amino acid residues. Models that decompose high-dimensional data into a series of steps predicted sequentially are termed autoregressive models, and they are well suited to variable-length data that have not been forced into a defined structure such as a multiple-sequence alignment. Autoregressive generative models are uniquely suited for modeling and designing the complex, highly diverse CDRs of antibodies. Here, we develop and apply a new autoregressive generative model that aims to capture key statistical properties of sets of sequences of variable lengths.

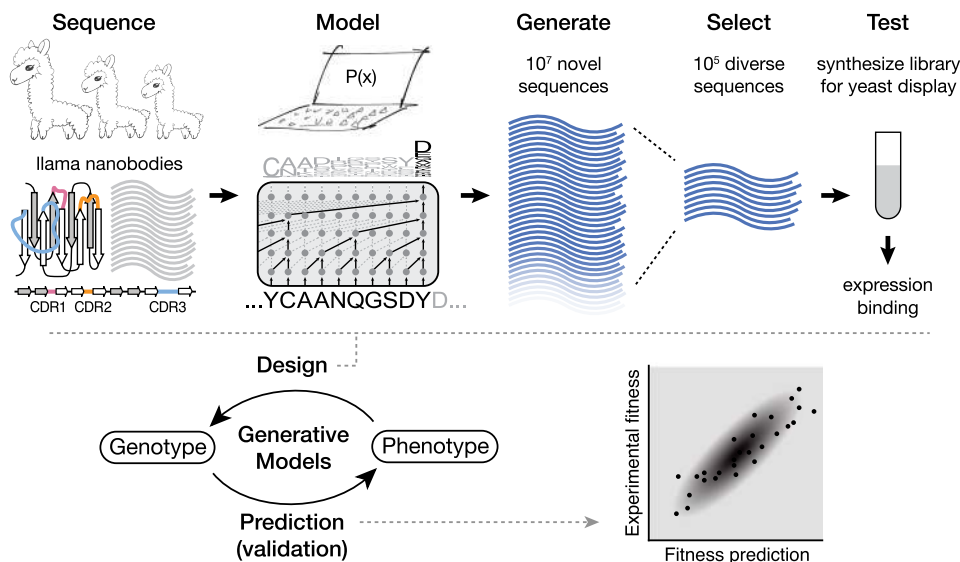
We first test our method on the problem of prediction of mutation effects, which are typically analyzed using alignment-based statistical methods. The new method performs on par with the DeepSequence machine-learning VAE-based method<sup>30</sup>, which does require aligned sequences and is an independent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

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

## Results

### An autoregressive generative model of biological sequences.

Protein sequences observed in organisms today result from mutation and selection for functional, folded proteins over time scales of a few days to a billion years. Generative models can be used to parameterize this view of evolution. Namely, they express the probability that a sequence  $x$  would be generated by evolution



**Fig. 1 Autoregressive models of biological sequences can learn the genotype-phenotype map for both prediction and design.** From natural sequences (gray) 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. Nanobodies have three highly variable complementarity determining regions (CDR1, CDR2, and CDR3). We then use these constraints to generate millions of novel nanobody sequences (blue)—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.

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

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

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

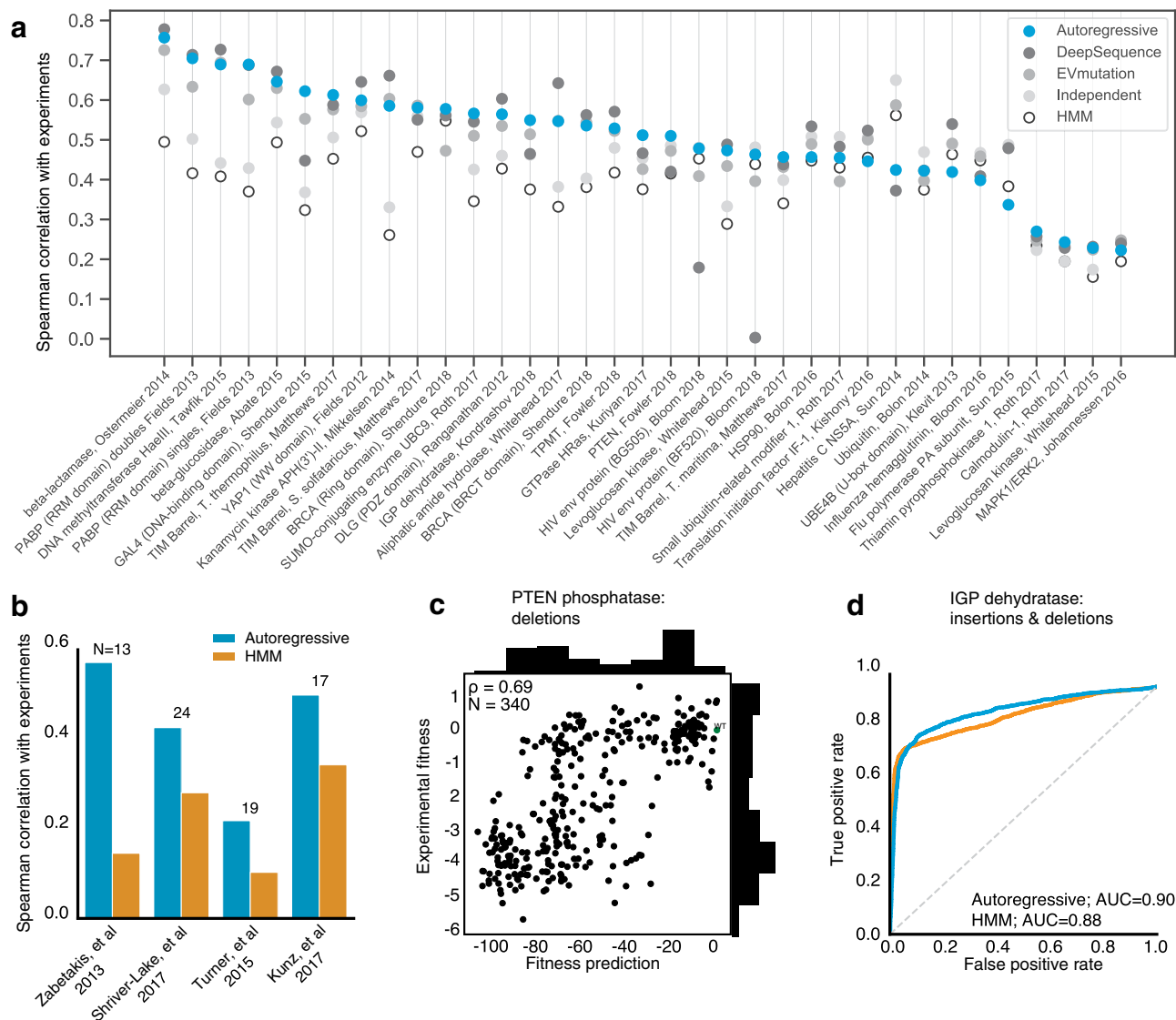
**The autoregressive model predicts experimental phenotype effects from sequences.** In order to gain confidence in the new

model for generating designed sequences, we first tested the ability of our new model to capture the dependencies between positions by testing the accuracy of mutation effect prediction. Somewhat surprisingly, unsupervised, generative models trained only on evolutionary sequences are proving the most accurate for predicting the effect of mutations when compared to large datasets of experimentally measured mutation effects<sup>30,34</sup>, and they avoid the risk of overfitting that can occur as a result of circularity in supervised methods<sup>68</sup>. We compared the accuracy of this new, non-alignment-based model to state-of-art methods for a benchmark set of 40 deep mutational scans across 33 different proteins, totaling 690,257 individual sequences (Supplementary Table 1).

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

$$\log \frac{p(\mathbf{x}^{\text{Mutant}}|\boldsymbol{\theta})}{p(\mathbf{x}^{\text{Wild-type}}|\boldsymbol{\theta})} \quad (2)$$

which estimates the plausibility of mutant sequence  $\mathbf{x}^{\text{Mutant}}$  relative to its wild-type, un-mutated counterpart,  $\mathbf{x}^{\text{Wild-type}}$ . This log-ratio has been shown to be predictive of mutation effects<sup>29,30</sup>. Importantly, this approach is fully unsupervised: rather than learning from experimental mutation effects, we can learn evolutionary constraints using only the space of natural sequences. We benchmark the model predictions against the deep mutational scan experiments and compare the Spearman’s rank correlation to state-of-art models trained on alignments of the same sequences. The autoregressive model is able to consistently match or outperform a model with only site-independent terms (30/40 datasets) and the EVmutation model<sup>29</sup> that includes dependencies between pairs of sites (30/40 datasets); it performs on par with the state-of-the-art results of DeepSequence<sup>30</sup> (19/40 datasets, average difference in rank correlation is only 0.09); and it outperforms the supervised Envision model<sup>31</sup> for 6/9 of the datasets tested (Fig. 2a;



**Fig. 2 Validation of the autoregressive model in learning the genotype to phenotype map.** **a** Even without using alignments, the autoregressive model (blue) can competitively match mutation effect prediction accuracies of state-of-art alignment-dependent models, such as conservation (light gray), evolutionary couplings (gray), and DeepSequence (dark gray)<sup>30</sup>. In addition, the mutation effect prediction accuracies improve upon hidden Markov model<sup>74</sup> (HMM, white) accuracies. Without using alignments, the autoregressive model matches alignment-dependent state-of-art missense mutation effect prediction (DeepSequence) 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 (blue) 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 HMMs (orange). 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 experimental fitness is accurate, with a Spearman correlation of 0.69. **d** Accurate prediction of binary fitness for IGP dehydratase with a range of insertions, deletions, and missense mutations of the autoregressive model (blue), higher than HMM (orange).

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

As with previous models that use evolutionary sequences, the accuracy of mutation effect prediction increases with increasing numbers of non-redundant sequences, as long as there is coverage of the length, tested here across eight of the protein families for four sequence depths (Supplementary Fig. 4 and Supplementary Table 2). Interestingly, the accuracy of effect predictions against the aliphatic amidase mutation scan is remarkably robust even with a low number of training sequences—123 non-redundant sequences provide the same accuracy as 36,000—suggesting that there is more to learn about the relationship between evolutionary sampling and model learning. For now, we suggest a  $M_{eff}/L > 5$  (number of effective sequences normalized by length) in order to sample enough diversity.



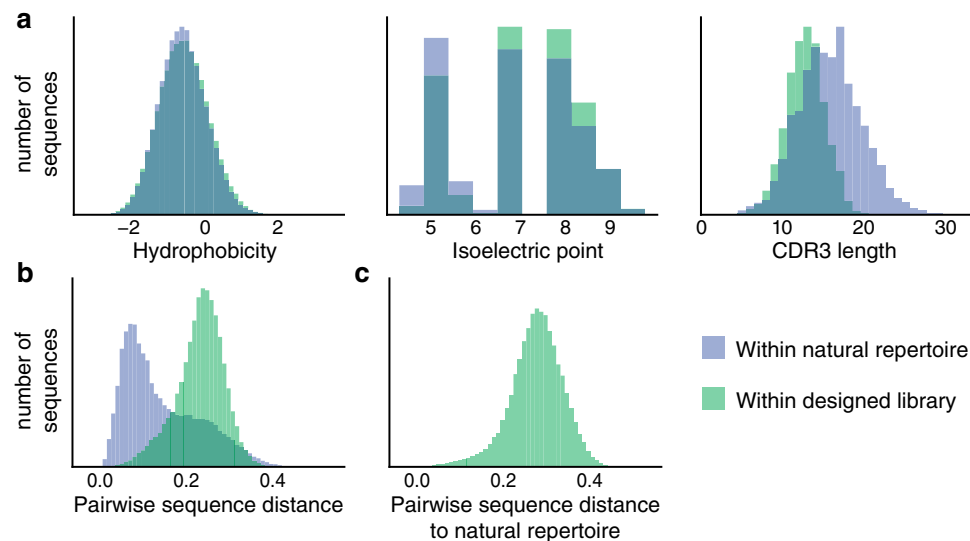
Because the autoregressive model is not dependent on alignments, we can now learn mappings of sequences of high variability and diverse lengths for which meaningful alignments are difficult or nonsensical to construct, such as antibody and nanobody sequences. The autoregressive model was thus also validated on nanobody thermostability measurements to test whether we could learn the sequence constraints of fit nanobodies, including the highly variable regions. To do so, we fit the autoregressive model to a set of ~1.2 million natural llama nanobody sequences<sup>57</sup>. Sequence likelihoods from this trained model are expected to reflect nanobody fitness, i.e., the multiple convolved aspects that nanobodies are selected for in vivo, including thermostability, expression, and potentially low poly-reactivity. Using this model, we find that the log-probability fitness calculations predict the thermostability of unseen llama nanobody sequences from four different stability experiments<sup>69–72</sup> (Fig. 2b, Supplementary Fig. 5, Supplementary Table 3, Supplementary Data 1). These experiments span a wide range of mutation types, lengths, and sequence diversity. The autoregressive model consistently outperforms a hidden Markov model (HMM, hmmer3)<sup>73,74</sup> in predicting the relationship between sequence and thermostability of nanobodies.

Previous alignment-dependent generative models are constrained to predicting the effects of missense mutations. However, in-frame insertions and deletions can also have large phenotypic consequences for protein function, yet these changes have proved difficult to model. We compare the fitness predictions calculated as log probabilities by the autoregressive model to experimental assays for the fitness of mutated biomolecules, using rank correlation ( $\rho$ ) for quantitative measurements and area under the receiver-operator curve (AUC) for binary fitness categorization, identifying the two groups with a two-component Gaussian mixture model. The model is able to capture the effects of single amino acid deletions on PTEN phosphatase<sup>75</sup> ( $\rho = 0.69$ ,  $N = 340$ , HMM  $\rho = 0.75$ ; PROVEAN  $\rho = 0.7$ ; Fig. 2c, Supplementary Data 2) and multiple amino acid insertions and deletions in imidazoleglycerol-phosphate (IGP) dehydratase<sup>76</sup> (AUC = 0.90,  $N = 6102$ , HMM AUC = 0.88; Fig. 2d, Supplementary Table 4, Supplementary Data 3). Here we use the AUROC metric for IGP dehydratase as the experimental data are bimodal with a large fraction at zero fitness. While PROVEAN<sup>77</sup> predicted the effect of single PTEN deletions comparably to our model, it fails to predict the effect of multiple insertions, deletions, and substitutions as were tested in IGP dehydratase and it cannot generate new sequences. Three additional insertion and deletion mutation scan fitness predictions are included in the supplement: yeast snoRNA ( $\rho = 0.49$ ; Supplementary Data 4), beta-lactamase ( $\rho = 0.45$ ; Supplementary Data 5), and p53 ( $\rho = 0.035$ ; Supplementary Data 6) (see Supplementary Fig. 6). Predicting the effects of indels can be important for disease-related genes: the four different single amino acid deletions annotated as pathogenic by Clinvar<sup>78</sup> in two cancer genes, *BRCA1* and *P53*, and one Alzheimer's-linked gene, *APOE*, is in the bottom 25th percentile of predicted deletion effect distributions (Supplementary Fig. 7). Other indels that are predicted to be highly deleterious by the autoregressive model may be of clinical interest for the experimental study of pathogenicity. We expect that the autoregressive model can predict mutation effects in disordered and low-complexity sequences. As a proof-of-concept, we have provided an in silico mutation scan of the human tau protein, which contains regions of low complexity and is strongly associated with neurodegenerative diseases (Supplementary Fig. 8; Supplementary Data 7). Our mutation effect prediction distinguishes between 40 pathogenic and 10 non-pathogenic mutations (two-tailed independent  $t = -4.1$ ,  $p = 0.001$ , AUC = 0.86; Supplementary Data 8) that were collected from the Alzforum database<sup>79</sup>.

### Generating an efficient library of functional nanobodies.

Screening large, high-throughput libraries of antibodies and nanobodies in vitro has become increasingly prevalent because it can allow for rapid identification of diverse monoclonal binders to target antigens. However, these synthetic libraries contain a large fraction of non-functional nanobody sequences. Natural nanobody sequences are selected against unfavorable biochemical properties such as instability, poly-reactivity, and aggregation during affinity maturation<sup>6</sup>. Similarly to nanobody thermostability prediction, we sought to learn the constraints that characterize functional nanobodies by fitting the autoregressive model to a set of ~1.2 million nanobody sequences from the immune repertoires of seven different naïve llamas<sup>57</sup>. Using this trained model and conditioning on the germline framework-CDR1-CDR2 nanobody sequence, we then generate over  $10^7$  fit sequences, generating one amino acid at a time based on the learned sequential constraints. As nanobody CDR3s often contact the framework in 3D, conditioning in this way allows the model to learn any resulting constraints on the CDR3 sequence and incorporate them during generation. We remove sequences that do not end with the final beta-strand of our nanobody template, duplicate sequences, and CDR3s likely to suffer post-translational modification to obtain ~3.7 million sequences (Supplementary Table 5). From these, we select 185,836 highly diverse CDR3 sequences for inclusion in our designed library. We compare our designed library to a state-of-art synthetic library<sup>8</sup>, which was constructed combinatorially based on position-specific amino acid frequencies of nanobody sequences with crystal structures in the PDB database. This library contains CDR3 sequences that have a similar distribution of biochemical properties as the naïve llama immune repertoire ("Methods"; Fig. 3a). The distribution of hydrophobicity and isoelectric points are similar to the natural llama repertoire even though explicit constraints on these properties were never imposed during the generation or selection of sequences for the designed library. The lengths of the CDR3 sequences in the designed library are shorter than the natural repertoire; this is due to the strategy of choosing cluster centroids during the selection of the  $10^5$  sequences and can be adjusted by changing the sampling method. Longer CDR3s may also be attained by allowing interloop disulfide bridges that stabilize longer CDR3s in some VHH domains<sup>80</sup>; this would require a different nanobody template and ideally camel or dromedary nanobody repertoires. The sequences in the designed library are diverse; they are more distant from each other than sequences in the natural repertoire (Fig. 3b) while maintaining nearly as much diversity as an equivalent sample of a combinatorial synthetic library<sup>8</sup> (Supplementary Fig. 9). In addition, we are exploring new regions of sequence space because the generated sequences in the designed library are diverse from the naïve repertoire (Fig. 3c).

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



**Fig. 3** The designed library has comparable biochemical property distributions and improved diversity to the natural llama repertoire. **a** Conditioned on the framework-CDR1-CDR2 sequence, a diverse set of CDR3 sequences are generated and selected. These designed CDR3 sequences (green) are similar to the natural repertoire (blue) in their distributions of hydrophobicity<sup>105</sup> and isoelectric point<sup>106, 107</sup>, while having shorter length distributions due to selection strategies in the final library construction. **b** The designed library contains more diversity in sequences than the natural repertoire as evidenced by the larger cosine distance to its nearest neighbor. **c** Each sequence in the designed library is diverse from any sequence seen in the natural repertoire, indicating that we have learned fit sequence constraints but are traversing previously unexplored regions of sequence space.

previous library (Fig. 4a, b). Furthermore, the designed library had nearly half the fraction of poorly expressed nanobodies (cells with fluorescence below 10,000 AU) as compared to the synthetic library (Fig. 4a, b) as well as a significant increase in the fraction of highly expressed nanobodies, as can be seen in the upper limits in the respective expression distributions (Fig. 4a and Supplementary Fig. 10). Expression experiments were performed with two replicates in addition to a single control experiment of yeast expressing a single well-behaved nanobody clone (Nb. 174684). These experimental results demonstrate that with the autoregressive model trained on natural llama nanobody sequences, we successfully designed a smart library consisting of a higher proportion of stable, well-expressed nanobodies.

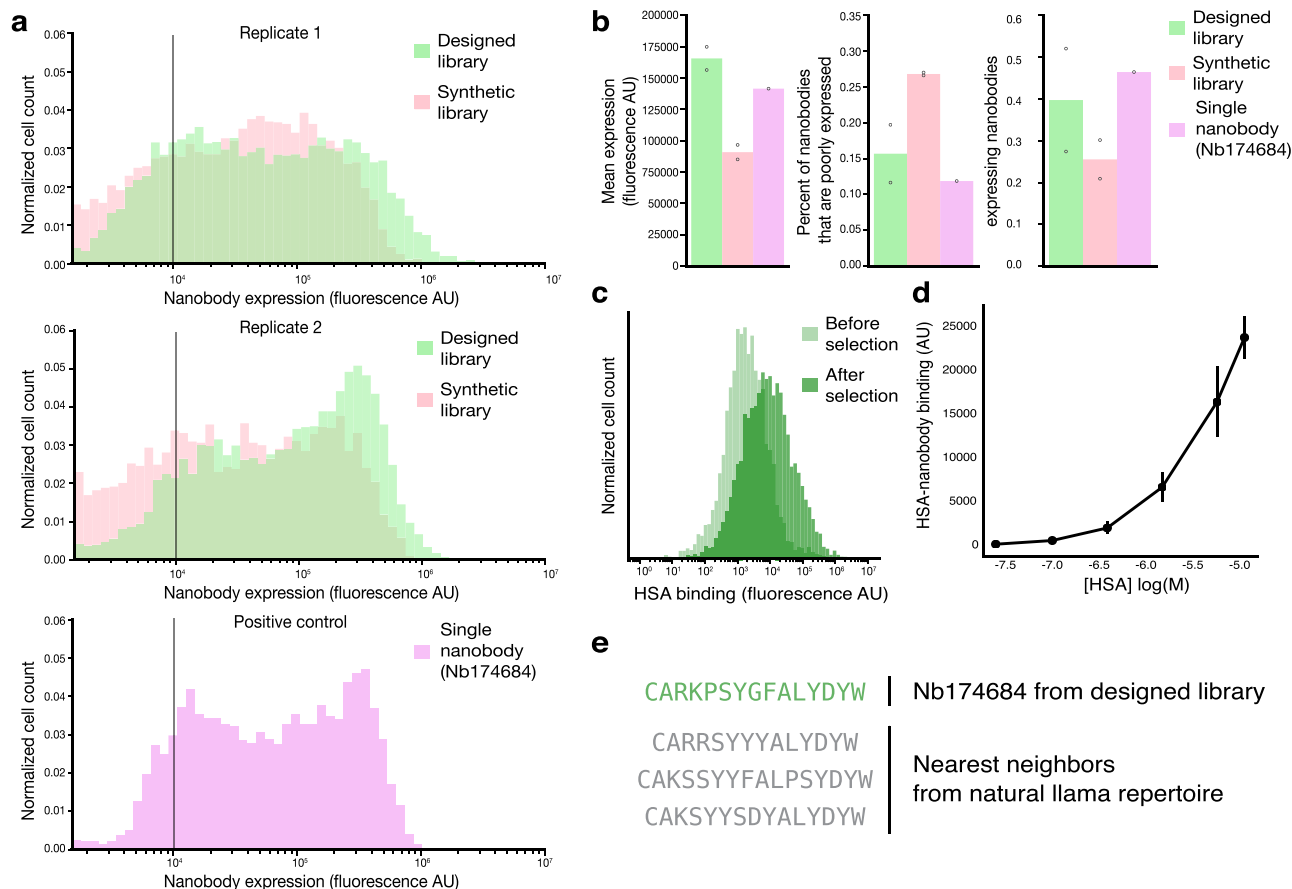
With this small designed library, we selected nanobody sequences that bound to human serum albumin (HSA) using fluorescence-activated cell sorting (FACS) (Fig. 4c), from which we were even able to identify weak to moderate binders—the strongest binder has a predicted  $K_d$  of 9.8  $\mu\text{M}$  (Fig.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

## Discussion

Here we show how neural network-powered generative autoregressive models can be used to model sequence constraints independent of alignments and design novel functional sequences for previously out-of-reach applications such as nanobodies. The

capability of these models is based on demonstrated state-of-the-art performance and on an extended range of applicability in the space of sequences. In the particular version in this paper, we validated our model first on deep mutational scan data, with on-par performance with the best currently available model<sup>29–31,34,77</sup>, and demonstrated application to examples for which robust alignments cannot be constructed, such as sequences with multiple insertions, deletions, and substitutions, and cases for which protein structures and experimental data are not available. As for the comparison with a potentially competing alignment-free model, while we do not discount the utility of semi-supervised methods (exploiting mutation effect-labeled experimental data), great care must be taken in the way the split between training and test is conducted to evaluate the true generalizability of the method. For instance, randomized subsets excluded from training will still be learned from the labeled data in a way that is not generalizable to required predictions for other proteins<sup>53,88,89</sup>. Our model is not subject to these limitations as its training is fully unsupervised.

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



**Fig. 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 (pink) and our designed library (green) in two biological replicate experiments as well as a control experiment of a single, well-expressed nanobody clone (Nb174684, purple). 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 (166,193 AU compared to 92,183 AU), nearly half the fraction of poorly expressed nanobodies (of cells expressing nanobodies) (15.4% compared to 25.7% of clones with less than 10,000 AU indicated as a gray bar in panel (a)), and one and a half times the fraction of total cells that express nanobodies (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 two biological replicates for the two libraries and one replicate for the control experiment of the single nanobody clone. Replicate measurements are displayed as dots for the two library experiments. **c** Fluorescence distributions of nanobodies bound to HSA show 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. The means of HSA binding (AU) of three replicates are shown and error bars represent standard deviations in measurements at each concentration of HSA. **e** CDR3 sequence of binder Nb.174684 and the sequences of the nearest neighbors from the natural llama repertoire that was used to train the autoregressive model.

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

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



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

## Methods

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

**Data collection.** Forty datasets which include experimental mutation effects, the sequence families, and effect predictions were taken from our previous publication<sup>30</sup> and five datasets that include indels and nanobody thermostability data were added for this work (references and data in Supplementary Table 4 and Supplementary Data 1–4). For new mutation effect predictions such as the indel mutation scans, sequence families were collected from the UniProt database in the same procedure as described in previous published work<sup>30</sup>, using jackhammer<sup>74</sup> and a default bit score of 0.5 bits/per residue for inclusion of sequence unless there was low coverage of the target sequence or not enough sequences. Pathogenic mutations for the Tau protein were downloaded from the Alzforum database<sup>79</sup>. The naïve llama immune repertoire was acquired from<sup>57</sup>. Due to a large number of sequences in the llama immune repertoire, sequence weights were approximated using Linclust<sup>102</sup> by clustering sequences at both 80 and 90% sequence identity thresholds.

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

From this large number of sequences, we then sought to choose roughly 200,000 CDR3 sequences that are both deemed fit by the model and as diverse from one another as possible to cover the largest amount of sequence space. First, we featurized these sequences into fixed length, L2 normalized k-mer vectors with

k-mers of sizes 1, 2, and 3. We then used BIRCH clustering<sup>103</sup> to find diverse members of the dataset in  $O(n)$  time. We used a diameter threshold of 0.575, resulting in 382,675 clusters. K-mer size and BIRCH diameter threshold were chosen to maximize the number of clusters within a memory constraint of 70 GB. From the cluster centroids, we chose the 185,836 most probable sequences for final library construction.

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

An oligonucleotide library was synthesized (Agilent) with the following design ACTCTGT [CDR3] ATCGT where CDR3 is a sequence for one of the computationally designed clones. Two hundred picomoles of the library were PCR amplified over 15 cycles with oligonucleotides Oligo\_library\_F and Oligo\_library\_R using Q5 polymerase (New England Biolabs). Amplified DNA was PCR purified (Qiagen) and ethanol precipitated in preparation for yeast transformation. In total,  $4.8 \times 10^8$  BJ5465 (MAT $\alpha$  ura352 trp1 leu2 $\Delta$ 1 his3 $\Delta$ 200 pep4::HIS3 prb1 $\Delta$ 1.6 R can1 GAL) yeast cells, grown to OD600 1.6, were transformed, using an ECM 830 Electroporator (BTX-Harvard Apparatus), with 2.4  $\mu$ g of NotI digested pYDS\_NbCM vector and 9.9  $\mu$ g of CDR3 library PCR product yielding  $2.7 \times 10^6$  transformants. Library aliquots of  $2.4 \times 10^8$  cells per vial were frozen in tryptophan dropout media containing 10% DMSO. A list of oligonucleotides can be found in Supplementary Table 6.

**Characterization of nanobody library.** Yeast displaying the computationally designed or combinatorial synthetic nanobody library<sup>8</sup> were grown in tryptophan dropout media with glucose as the sugar source for 1 day at 30 °C and then passaged into media with galactose as the sole sugar source to induce expression of nanobodies at 25 °C. After 2 days of induction, one million cells from each library were stained with a 1:25 dilution of anti-HA AlexaFluor647 conjugated antibody (Cell Signaling Technology) in Buffer A (20 mM HEPES pH 7.5, 150 mM NaCl, 0.1% BSA, 0.2% maltose) for 30 min at 4 °C. After staining, cells were centrifuged, the supernatant was removed, and cells were resuspended in Buffer A for flow analysis with an Accuri C6 (BD Biosciences, Supplementary Fig. 12).

To find nanobody binders to HSA one round of magnetic-activated cell sorting (MACS) followed by two rounds of FACS (FACS, with SONY SH800Z Sorter) were performed on our yeast-displayed library of nanobodies. For MACS,  $4 \times 10^7$  induced cells were resuspended in binding buffer (20 mM HEPES pH 7.5, 150 mM NaCl, 0.1% ovalbumin) along with anti-fluorescein isothiocyanate (FITC) microbeads (Miltenyi) and FITC-labeled streptavidin for 35 min at 4 °C and then passed through an LD column (Miltenyi) to remove binders to microbeads and streptavidin. The remaining yeast was centrifuged and resuspended in binding buffer and incubated with 500 nM streptavidin-FITC and 2  $\mu$ M of biotinylated HSA for 1 h at 4 °C. Yeast was then centrifuged and resuspended in binding buffer containing anti-FITC microbeads for 15 min at 4 °C before passing them into an LS column and eluting and collecting the bound yeast. For the first round of FACS, induced yeast was first stained with 1  $\mu$ M of biotinylated HSA for 45 min at 4 °C and then briefly stained with 500 nM of streptavidin tetramer along with antiHA-488 to assess expression levels. Both yeast stainings were performed in FACS buffer (20 mM HEPES pH 7.5, 150 mM NaCl, 0.1% ovalbumin, 0.2% maltose). In total,  $5 \times 10^6$  yeast were sorted and 28,000 were collected and expanded for the second round of FACS. The second round of FACS was performed under the same conditions as the first and from  $3.8 \times 10^6$  sorted yeast 21,455 were collected. Nanobody Nb174684 was isolated from a screen of 36 clones for binding to HSA using a flow cytometer and then sequenced. In order to characterize the binding of Nb174684, yeast displaying Nb174684 were stained with varying amounts of AlexaFluor 488 labeled HSA and fluorescence was analyzed with a flow cytometer. FACS measurements were analyzed using FlowJo and the python package FlowCytometryTools.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

All data generated and analyzed during the study are available in this published article, its supplementary information files, and on the github repository (<https://github.com/debbiemarkslab/SeqDesign>; <https://doi.org/10.5281/zenodo.4606785><sup>104</sup>).

## Code availability

All code used for model training, mutation effect prediction, sequence generation, and library generation is also available on the github repository (<https://github.com/debbiemarkslab/SeqDesign>; <https://doi.org/10.5281/zenodo.4606785><sup>104</sup>).

Received: 28 February 2021; Accepted: 26 March 2021;

Published online: 23 April 2021

**References**

- Romero, P. A. & Arnold, F. H. Exploring protein fitness landscapes by directed evolution. *Nat. Rev. Mol. Cell Biol.* **10**, 866–876 (2009).
- Dougherty, M. J. & Arnold, F. H. Directed evolution: new parts and optimized function. *Curr. Opin. Biotechnol.* **20**, 486–491 (2009).
- Baker, D. An exciting but challenging road ahead for computational enzyme design. *Protein Sci.* **19**, 1817–1819 (2010).
- Huang, P.-S., Boyken, S. E. & Baker, D. The coming of age of de novo protein design. *Nature* **537**, 320–327 (2016).
- Fowler, D. M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nat. Methods* **11**, 801–807 (2014).
- Muyldermans, S. Nanobodies: natural single-domain antibodies. *Annu. Rev. Biochem.* **82**, 775–797 (2013).
- Sall, A. et al. Generation and analyses of human synthetic antibody libraries and their application for protein microarrays. *Protein Eng. Des. Sel.* **29**, 427–437 (2016).
- McMahon, C. et al. Yeast surface display platform for rapid discovery of conformationally selective nanobodies. *Nat. Struct. Mol. Biol.* **25**, 289–296 (2018).
- Bradbury, A. R., Sidhu, S., Dubel, S. & McCafferty, J. Beyond natural antibodies: the power of in vitro display technologies. *Nat. Biotechnol.* **29**, 245–254 (2011).
- Schoof, M. et al. An ultra-potent synthetic nanobody neutralizes SARS-CoV-2 by locking Spike into an inactive conformation. *bioRxiv*, 2020.2008.2008.238469, <https://doi.org/10.1101/2020.08.08.238469> (2020).
- Miho, E., Roskar, R., Greiff, V. & Reddy, S. T. Large-scale network analysis reveals the sequence space architecture of antibody repertoires. *Nat. Commun.* **10**, 1321 (2019).
- Jain, T. et al. Biophysical properties of the clinical-stage antibody landscape. *Proc. Natl Acad. Sci. USA* **114**, 944–949 (2017).
- Marks, C. & Deane, C. M. How repertoire data are changing antibody science. *J. Biol. Chem.* **295**, 9823–9837 (2020).
- Asti, L., Uguzzoni, G., Marcatili, P. & Pagnani, A. Maximum-entropy models of sequenced immune repertoires predict antigen-antibody affinity. *PLoS Comput. Biol.* **12**, e1004870 (2016).
- Mora, T., Walczak, A. M., Bialek, W. & Callan, C. G. Jr Maximum entropy models for antibody diversity. *Proc. Natl Acad. Sci. USA* **107**, 5405–5410 (2010).
- Marcou, Q., Mora, T. & Walczak, A. M. High-throughput immune repertoire analysis with IGoR. *Nat. Commun.* **9**, 561 (2018).
- Liu, G. et al. Antibody complementarity determining region design using high-capacity machine learning. *Bioinformatics* **36**, 2126–2133 (2020).
- DeKosky, B. J. et al. Large-scale sequence and structural comparisons of human naive and antigen-experienced antibody repertoires. *Proc. Natl Acad. Sci. USA* **113**, E2636–E2645 (2016).
- Muyldermans, S. A guide to: generation and design of nanobodies. *FEBS J* **288**, 2084–2102 (2020).
- Zimmermann, I. et al. Synthetic single domain antibodies for the conformational trapping of membrane proteins. *Elife* **7**, <https://doi.org/10.7554/eLife.34317> (2018).
- Marks, D. S. et al. Protein 3D structure computed from evolutionary sequence variation. *PLoS ONE* **6**, e28766 (2011).
- Hopf, T. A. et al. Three-dimensional structures of membrane proteins from genomic sequencing. *Cell* **149**, 1607–1621 (2012).
- Marks, D. S., Hopf, T. A. & Sander, C. Protein structure prediction from sequence variation. *Nat. Biotechnol.* **30**, 1072–1080 (2012).
- Kamisetty, H., Ovchinnikov, S. & Baker, D. Assessing the utility of coevolution-based residue-residue contact predictions in a sequence- and structure-rich era. *Proc. Natl Acad. Sci. USA* **110**, 15674–15679 (2013).
- Ovchinnikov, S., Kamisetty, H. & Baker, D. Robust and accurate prediction of residue-residue interactions across protein interfaces using evolutionary information. *Elife* **3**, e02030 (2014).
- Hopf, T. A. et al. Sequence co-evolution gives 3D contacts and structures of protein complexes. *eLife* **3**, e03430 (2014).
- Cong, Q., Anishchenko, I., Ovchinnikov, S. & Baker, D. Protein interaction networks revealed by proteome coevolution. *Science* **365**, 185 (2019).
- Green, A. G. et al. Proteome-scale discovery of protein interactions with residue-level resolution using sequence coevolution. *Nat Commun* **12**, 1396 (2019).
- Hopf, T. A. et al. Mutation effects predicted from sequence co-variation. *Nat. Biotechnol.* **35**, 128–135 (2017).
- Riesselman, A. J., Ingraham, J. B. & Marks, D. S. Deep generative models of genetic variation capture the effects of mutations. *Nat. Methods* **15**, 816–822 (2018).
- Gray, V. E., Hause, R. J., Luebeck, J., Shendure, J. & Fowler, D. M. Quantitative missense variant effect prediction using large-scale mutagenesis data. *Cell Syst.* **6**, 116–124 e113 (2018).
- Mann, J. K. et al. The fitness landscape of HIV-1 gag: advanced modeling approaches and validation of model predictions by in vitro testing. *PLoS Comput Biol.* **10**, e1003776 (2014).
- Figliuzzi, M., Jacquier, H., Schug, A., Tenaillon, O. & Weigt, M. Coevolutionary landscape inference and the context-dependence of mutations in beta-lactamase TEM-1. *Mol. Biol. Evol.* **33**, 268–280 (2015).
- Livesey, B. J. & Marsh, J. A. Using deep mutational scanning to benchmark variant effect predictors and identify disease mutations. *Mol. Syst. Biol.* **16**, e9380 (2020).
- Sim, N. L. et al. SIFT web server: predicting effects of amino acid substitutions on proteins. *Nucleic Acids Res.* **40**, W452–W457 (2012).
- Adzhubei, I., Jordan, D. M. & Sunyaev, S. R. Predicting functional effect of human missense mutations using PolyPhen-2. *Curr Protoc Hum Genet.*, Unit7 20, <https://doi.org/10.1002/0471142905.hg0720s76> (2013).
- Kircher, M. et al. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* **46**, 310–315 (2014).
- Reva, B., Antipin, Y. & Sander, C. Predicting the functional impact of protein mutations: application to cancer genomics. *Nucleic Acids Res.* **39**, e118 (2011).
- Lefranc, M. P. et al. IMGT unique numbering for immunoglobulin and T cell receptor variable domains and Ig superfamily V-like domains. *Dev. Comp. Immunol.* **27**, 55–77 (2003).
- van der Lee, R. et al. Classification of intrinsically disordered regions and proteins. *Chem. Rev.* **114**, 6589–6631 (2014).
- Mullaney, J. M., Mills, R. E., Pittard, W. S. & Devine, S. E. Small insertions and deletions (INDELs) in human genomes. *Hum. Mol. Genet.* **19**, R131–R136 (2010).
- Lin, M. et al. Effects of short indels on protein structure and function in human genomes. *Sci. Rep.* **7**, 9313 (2017).
- Mills, R. E. et al. Natural genetic variation caused by small insertions and deletions in the human genome. *Genome Res.* **21**, 830–839 (2011).
- Pentony, M. M. & Jones, D. T. Modularity of intrinsic disorder in the human proteome. *Proteins* **78**, 212–221 (2010).
- Campbell, P. J. et al. Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
- Turajlic, S. et al. Insertion-and-deletion-derived tumour-specific neoantigens and the immunogenic phenotype: a pan-cancer analysis. *Lancet Oncol.* **18**, 1009–1021 (2017).
- Deiana, A., Forcelloni, S., Porrello, A. & Giansanti, A. Intrinsically disordered proteins and structured proteins with intrinsically disordered regions have different functional roles in the cell. *PLoS ONE* **14**, e0217889 (2019).
- Uversky, V. N. et al. Unfoldomics of human diseases: linking protein intrinsic disorder with diseases. *BMC Genom.* **10**(Suppl 1), S7 (2009).
- Graves, A., Mohamed, A. & Hinton, G. Speech recognition with deep recurrent neural networks. In *Proc. IEEE International Conference on Acoustics, Speech and Singal Processing*, 6645–6649 (2013).
- Wang, Y. et al. Tacotron: Towards end-to-end speech synthesis. *arXiv* **1703**, 10135 (2017).
- Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv* **1409**, 0473 (2014).
- Sutskever, I., Vinyals, O. & Le, Q. V. Sequence to sequence learning with neural networks. *Adv. Neural Inf. Process. Syst.* **27**, 3104–3112 (2014).
- Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods* **16**, 1315–1322 (2019).
- Linder, J., Bogard, N., Rosenberg, A. B. & Seelig, G. A generative neural network for maximizing fitness and diversity of synthetic DNA and protein sequences. *Cell Syst.* **11**, 49–62 e16 (2020).
- Strodthoff, N., Wagner, P., Wenzel, M. & Samek, W. UDSPProt: universal deep sequence models for protein classification. *Bioinformatics* **36**, 2401–2409 (2020).
- Brookes, D. H., Park, H. & Listgarten, J. Conditioning by adaptive sampling for robust design. In *Proc. 36th International Conference on Machine Learning*, 773–782 (2019).
- McCoy, L. E. et al. Molecular evolution of broadly neutralizing Llama antibodies to the CD4-binding site of HIV-1. *PLoS Pathog.* **10**, e1004552 (2014).
- Chan, S. K., Rahumatullah, A., Lai, J. Y. & Lim, T. S. Naive human antibody libraries for infectious diseases. *Adv. Exp. Med Biol.* **1053**, 35–59 (2017).
- Ravikumar, A., Arzumanyan, G. A., Obadi, M. K. A., Javanpour, A. A. & Liu, C. C. Scalable, continuous evolution of genes at mutation rates above genomic error thresholds. *Cell* **175**, 1946–1957 e1913 (2018).

60. Vaswani, A. et al. Attention is all you need. *Adv. Neural Inf. Process. Syst.* **30**, 5998–6008 (2017).
61. Sutskever, I., Martens, J. & Hinton, G. Generating text with recurrent neural networks. In *Proc. 28th International Conference on Machine Learning (ICML-11)*, 1017–1024 (2011).
62. Pascanu, R., Mikolov, T. & Begio, Y. On the difficulty of training recurrent neural networks. In *Proc. International Conference on Machine Learning*, 1310–1318 (2013).
63. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997).
64. Cho, K. et al. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proc. 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1724–1734 (2014).
65. van den Oord, A. et al. Wavenet: a generative model for raw audio. *arXiv* **1609**, 03449 (2016).
66. Kalchbrenner, N. et al. Neural machine translation in linear time. *arXiv* **1610**, 100099 (2016).
67. Gupta, A. & Rush, A. Dilated convolutions for modeling long-distance genomic dependencies. *arXiv* **1710**, 01278 (2017).
68. Grimm, D. G. et al. The evaluation of tools used to predict the impact of missense variants is hindered by two types of circularity. *Hum. Mutat.* **36**, 513–523 (2015).
69. Kunz, P. et al. Exploiting sequence and stability information for directing nanobody stability engineering. *Biochim Biophys. Acta Gen. Subj.* **1861**, 2196–2205 (2017).
70. Shriver-Lake, L. C., Zabetakis, D., Goldman, E. R. & Anderson, G. P. Evaluation of anti-botulinum neurotoxin single domain antibodies with additional optimization for improved production and stability. *Toxicon* **135**, 51–58 (2017).
71. Turner, K. B. et al. Improving the biophysical properties of anti-ricin single-domain antibodies. *Biotechnol. Rep.* **6**, 27–35 (2015).
72. Zabetakis, D., Anderson, G. P., Bayya, N. & Goldman, E. R. Contributions of the complementarity determining regions to the thermal stability of a single-domain antibody. *PLoS ONE* **8**, e77678 (2013).
73. Durbin, R., Eddy, S. R., Krogh, A. & Mitchison, G. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. (Cambridge university press, 1998).
74. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
75. Mighell, T. L., Evans-Dutson, S. & O’Roak, B. J. A saturation mutagenesis approach to understanding PTEN lipid phosphatase activity and genotype-phenotype relationships. *Am. J. Hum. Genet.* **102**, 943–955 (2018).
76. Pokusaeva, V. O. et al. An experimental assay of the interactions of amino acids from orthologous sequences shaping a complex fitness landscape. *PLoS Genet.* **15**, e1008079 (2019).
77. Choi, Y., Sims, G. E., Murphy, S., Miller, J. R. & Chan, A. P. Predicting the functional effect of amino acid substitutions and indels. *PLoS ONE* **7**, e46688 (2012).
78. Landrum, M. J. et al. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
79. M. A. P. T. Alzforum. Retrieved August 12, 2020. from <https://www.alzforum.org/mutations/mapt>.
80. Harmsen, M. M. et al. Llama heavy-chain V regions consist of at least four distinct subfamilies revealing novel sequence features. *Mol. Immunol.* **37**, 579–590 (2000).
81. Beerli, R. R. & Rader, C. Mining human antibody repertoires. *MAbs* **2**, 365–378 (2010).
82. Dimitrov, J. D., Pashov, A. D. & Vassilev, T. L. Antibody polyspecificity: what does it matter? *Adv. Exp. Med. Biol.* **750**, 213–226 (2012).
83. Dimitrov, J. D. et al. Antibody polyreactivity in health and disease: statu variabilis. *J. Immunol.* **191**, 993–999 (2013).
84. Kelly, R. L., Zhao, J., Le, D. & Wittrup, K. D. Nonspecificity in a nonimmune human scFv repertoire. *MAbs* **9**, 1029–1035 (2017).
85. Lim, C. C., Choong, Y. S. & Lim, T. S. Cognizance of molecular methods for the generation of mutagenic phage display antibody libraries for affinity maturation. *Int J Mol Sci.* **20**, <https://doi.org/10.3390/ijms20081861> (2019).
86. Pashova, S., Schneider, C., von Gunten, S. & Pashov, A. Antibody repertoire profiling with mimotope arrays. *Hum. Vacc Immunother.* **13**, 314–322 (2017).
87. Wardemann, H. et al. Predominant autoantibody production by early human B cell precursors. *Science* **301**, 1374–1377 (2003).
88. Rives, A. et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *bioRxiv*, <https://doi.org/10.1101/622803> (2019).
89. Rao, R. et al. Evaluating protein transfer learning with TAPE. In *Proc. 33rd Conference on Neural Information Processing Systems* (2019).
90. Wellner, A. et al. Rapid generation of potent antibodies by autonomous hypermutation in yeast. *bioRxiv* **2020.11.11**, 378778 (2020).
91. DeKosky, B. J. et al. High-throughput sequencing of the paired human immunoglobulin heavy and light chain repertoire. *Nat. Biotechnol.* **31**, 166–169 (2013).
92. Ingraham, J. B., Vikas, G. K., Barzilay, R. & Jaakkola, T. Generative models for graph-based protein design. In *Proc. 33rd Conference on Neural Information Processing Systems* 15794–15805 (2019).
93. Kim, Y., Wiseman, S., Miller, A. C., Sontag, D. & Rush, A. Semi-amortized variational autoencoders. *arXiv* **1802**, 02550 (2018).
94. Yang, Z., Hu, Z., Salakhutdinov, R. & Berg-Kirkpatrick, T. Improved variational autoencoders for text modeling using dilated convolutions. *arXiv* **1702**, 08139 (2017).
95. van den Oord, A. & Vinyals, O. Neural discrete representation learning. *Adv. Neural Inf. Process. Syst.* **30**, 6306–6315 (2017).
96. Greener, J. G., Moffat, L. & Jones, D. T. Design of metalloproteins and novel protein folds using variational autoencoders. *Sci. Rep.* **8**, 16189 (2018).
97. Vijayakumar, A. K. et al. Diverse beam search: decoding diverse solutions from neural sequence models. *arXiv* **1610**, 02424 (2016).
98. Salimans, T. & Kingma, D. P. Weight normalization: a simple reparametrization to accelerate training of deep neural networks. *Adv. Neural Inf. Process. Syst.* **29**, 901–909 (2016).
99. Ba, J. L., Kiros, J. R. & Hinton, G. Layer normalization. *arXiv* **1607**, 06450 (2016).
100. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014).
101. Kingma, D. P. & Ba, J. L. Adam: a method for stochastic optimization. *arXiv* **1412**, 6980 (2014).
102. Steinegger, M. & Soding, J. Clustering huge protein sequence sets in linear time. *Nat. Commun.* **9**, 2542 (2018).
103. Zhang, T., Ramakrishnan, R. & Livny, M. BIRCH: an efficient data clustering method for very large databases. *ACM SIGMOD Rec.* **25**, 103–114 (1996).
104. Shin, J.-E., Riesselman, A. J., Kollasch, A. W. & Marks, D. S. SeqDesign. <https://doi.org/10.5281/zenodo.4606785> (2021).
105. Kyte, J. & Doolittle, R. F. A simple method for displaying the hydrophobic character of a protein. *J. Mol. Biol.* **157**, 105–132 (1982).
106. Bjellqvist, B. et al. The focusing positions of polypeptides in immobilized pH gradients can be predicted from their amino acid sequences. *Electrophoresis* **14**, 1023–1031 (1993).
107. Cock, P. J. A. et al. Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).

## Acknowledgements

We would like to thank John Ingraham, members of the Marks and Sander labs, and Harvard Research Computing for their insight and feedback on our research. J.-E.S., A.W.K., A.C.K., and D.S.M. are funded by an NIH TR01 grant (R01CA260415). C.M. and A.C.K. were funded through DP5 (DP5OD021345). A.M. is funded through DP5 (DP5OD023048). C.S. is funded by the Chan Zuckerberg Foundation (CZF2019-002433).

## Author contributions

D.S.M., A.C.K., and A.J.R. conceived the project; A.J.R. constructed the model; J.-E.S., A.J.R., A.W.K., and E.S. designed and evaluated computational experiments for validation, prediction, and generation of sequences; A.M. compiled natural nanobody sequence data; C.M. constructed the library and performed experiments; J.-E.S., A.W.K., and C.M. analyzed the library experimental data; A.J.R., J.-E.S., A.W.K., C.M., C.S., A.C.K., and D.S.M. wrote the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-021-22732-w>.

**Correspondence** and requests for materials should be addressed to A.C.K. or D.S.M.

**Peer review information** *Nature Communications* thanks the anonymous reviewer(s) for their contribution to the peer review of this work.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher’s note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021