

RESEARCH ARTICLE

Ig-VAE: Generative modeling of protein structure by direct 3D coordinate generation

Raphael R. Eguchi^{1,2}, Christian A. Choe³, Po-Ssu Huang^{3*}

1 Department of Biochemistry, Stanford University, Stanford, California, United States of America, **2** Department of Statistics, Stanford University, Stanford, California, United States of America, **3** Department of Bioengineering, Stanford University, Stanford, California, United States of America

* possu@stanford.edu

Abstract

While deep learning models have seen increasing applications in protein science, few have been implemented for protein backbone generation—an important task in structure-based problems such as active site and interface design. We present a new approach to building class-specific backbones, using a variational auto-encoder to directly generate the 3D coordinates of immunoglobulins. Our model is torsion- and distance-aware, learns a high-resolution embedding of the dataset, and generates novel, high-quality structures compatible with existing design tools. We show that the Ig-VAE can be used with Rosetta to create a computational model of a SARS-CoV2-RBD binder via latent space sampling. We further demonstrate that the model's generative prior is a powerful tool for guiding computational protein design, motivating a new paradigm under which backbone design is solved as constrained optimization problem in the latent space of a generative model.

OPEN ACCESS

Citation: Eguchi RR, Choe CA, Huang P-S (2022) Ig-VAE: Generative modeling of protein structure by direct 3D coordinate generation. PLoS Comput Biol 18(6): e1010271. <https://doi.org/10.1371/journal.pcbi.1010271>

Editor: Joanna Slusky, University of Kansas, UNITED STATES

Received: February 15, 2022

Accepted: June 1, 2022

Published: June 27, 2022

Copyright: © 2022 Eguchi et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: 1. Code for the working model is provided at the GitHub repository: <https://github.com/ProteinDesignLab/IgVAE> 2. The full training dataset is available at: tinyurl.com/igvaedataset.

Funding: P.-S.H was supported by startup funds from the Stanford Schools of Engineering and Medicine and the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Scientific Discovery through Advanced Computing (SciDAC) program. R.R.E was supported by the Stanford ChEM-H Chemistry/

Author summary

Many essential biochemical processes are governed by protein-protein interactions (PPIs), and our ability to make binding proteins that modulate PPIs is crucial to the creation of therapeutics and the study of cell-signaling. One critical aspect of PPI design is to capture protein conformational flexibility. Deep generative models are a class of mathematical models that are able to synthesize novel data from a finite set of training examples. Here, we make advances in computational protein design methodology by developing a deep generative model that creates protein backbones adopting the immunoglobulin fold, which is found in natural binding proteins such as antibodies. While generative models have been powerful in tasks such as image generation, using them to create proteins has remained a challenge. We solve this problem with a new model that allows for the direct generation of novel 3D molecules and show that they are of high chemical accuracy. Generated structures work well with existing protein design methods such as Rosetta, providing access to a large collection of novel immunoglobulin structures. Finally, we present a new protein design framework, called “generative design,” that shows how deep generative models such as ours can be applied to virtually any protein design problem.

Biology Interface Predoctoral Training Program and the National Institute of General Medical Sciences of the National Institutes of Health under Award Number T32GM120007. This project is also based upon work supported by Google Cloud credits awarded to P.-S.H, R.R.E and C.A.C. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

1. Introduction

Over the past two decades, structure-based protein design has provided novel solutions to challenging problems such as enzyme catalysis, viral inhibition, de novo structure generation and more [1,2,3,4,5,6,7,8,9]. In the vast majority of successful design examples using the popular Rosetta framework, the protein design process consists of two steps: (1) generation of a protein backbone, and (2) design of a sequence that minimizes the folded-state energy of the generated backbone. During the first step, backbone conformations are massively sampled without amino acid identity to search for templates that can host features such as catalytic sites, or loop conformations suited to particular functions. In the second step, sequences are chosen to both sustain the folded-state structure and provide the chemical elements that perform the desired function.

While many methods have been developed to perform sequence design [10,11], few are capable of generating backbones. Currently the majority of design templates are generated using fragment sampling in combination with expert-specified topologies of loop and secondary structure elements. Creating novel backbones for which there exist a foldable sequence remains one of the greatest challenges in protein engineering, and most engineering endeavors rely only on native structures as templates. In this study, we focus on this task of protein backbone generation, and seek to develop a method that allows us to: (1) generate novel, designable structures from the limited set of existing structural data, and (2) generate backbones that satisfy user-specified design criteria.

With recent advances in deep learning technology, machine learning tools have seen increasing applications in protein science, with deep neural networks being applied to tasks such as sequence design [11], fold recognition [12], binding site prediction [13], and structure prediction [14,15]. Generative models, which approximate the distributions of the data they are trained on, have garnered interest as a data-driven way to create novel proteins. Unfortunately the majority of protein-generators create 1D amino acid sequences [16,17,18,19] making them unsuitable for problems that require structure-based solutions such as designing protein-protein interfaces.

A major challenge in the field of 3D deep learning arises from the fact that 3D coordinates lack rotational and translational invariance, making generalizable feature learning and generation difficult. To address this challenge, our own group was the first to report a Generative Adversarial Network (GAN) that generated 64-residue peptide backbones using a distance matrix representation [20] that preserved the desired invariances. 3D coordinates were recovered using a convex optimization algorithm [20] and later, a learned coordinate recovery module [21]. Despite its novelty, the GAN method was accompanied by several difficulties. First, the generated distance matrices were not Euclidean-valid, and thus it was not possible to recover 3D coordinates that perfectly satisfied the generated distances. Second, because of the redundancy of the distance matrix representation under reflection, the quality of the torsion distributions were often degraded, leading to loss of important biochemical features, such as hydrogen-bonding, in many outputs. Ultimately, the pure distance-based representation of our GAN yielded structures that, while novel, were often chemically unrealistic and unsuitable for design [22]. Although a few other algorithms that generate contacts have been reported [23,20,21,24], all of these methods require external tools to build or recover 3D coordinates.

Here, we present a new variational autoencoder (VAE) [25] architecture that is the first to perform direct 3D coordinate generation of full-atom protein backbones, circumventing the need to recover coordinates from pairwise distance constraints and avoiding the problem of distance matrix validity [26]. Our model is novel in that it provides a loss function that is rotationally and translationally invariant, while also having torsional awareness and the ability to

directly generate 3D coordinates. We train our VAE using a unique loss function, which formulates coordinate generation as the solution to the joint problems of distance matrix reconstruction and torsion angle inference, both of which preserve the desired invariances. With the intention of delivering high-quality, immediately designable templates, our model performs class-specific generation of immunoglobulin (Ig) proteins, which are comprised of a two-layer β -sandwich structure supporting variable loop regions. Immunoglobulins are highly versatile in their target-binding capabilities, serving as core components in naturally occurring antibodies and in biologics such as ScFv's [27], nanobodies [28,29,30], and more.

Importantly, our model motivates a conceptually new way of solving protein design problems. Because the Ig-VAE generates coordinates directly, all of its outputs are fully differentiable through the coordinate representation. This allows us to use the Ig-VAE's generative prior to constrain structure generation with any differentiable coordinate-based heuristic, such as Rosetta energy, backbone shape constraints, packing metrics, and more. By constraining and optimizing a structure in the VAE latent space, designers are able to specify any desired structural features while the model creates the rest of the molecule. As an example of this approach, we use the Ig-VAE to perform constrained loop generation, towards epitope-specific antibody design. Our technology ultimately paves the way for a novel approach to protein design in which backbone construction is solved via constrained optimization in the latent space of a generative model. In contrast to conventional methods [31,32], we term this approach "generative design."

2. Methods

2.1 Dataset

All training data were collected from the antibody structure database AbDb/abYbank [33]. The training set is comprised of 10768 individual immunoglobulin domains, 4154 of which are non-sequence-redundant, including single-domain antibodies. Redundancy was determined by specification in the AbDb/abYbank database, which is based on sequence identity in the loop regions [33]. To construct a non-redundant dataset, single chains were sampled from each AbDb/abYbank cluster and sequences were further checked for redundancy in the loop regions to ensure that the single-domain sequences adhered to the AbDb/abYbank criteria. Domains that were missing residues were excluded, and sequence-redundant structures were included in the training set to allow the network to learn small backbone fluctuations, which are highly relevant to backbone design. For benchmarking we used only sets of structures sampled from the non-redundant subset of AbDb/abYbank. The training set covers 99.935% of the AbDb/abYbank database as of July 2019, and only 7 structures were excluded due to mid-chain missing residues. Structures in the dataset vary in length from 89 to 138 residues, with most falling between 114 and 130. Since the input of our model was fixed at 128 residues (512 atoms), structures larger than 128 were center-cropped. Structures smaller than 128 were "structurally padded" by using RosettaRemodel to append dummy residues to the N and C termini. The reconstruction loss of the padded regions was down-weighted to zero over the course of training (see S1 Methods in S1 File), and treated analogously to conventional zero-padding during analysis. All structures were idealized and relaxed under the Rosetta energy function with constraints to starting coordinates [10]. This relaxation step was done to remove any potential confounding factors resulting from various crystal structure optimization procedures.

2.2 Model architecture and training

While designing our model we sought to implement three features we feel are essential to any deep-learning-based backbone generator. First, we wanted to preserve rotational and

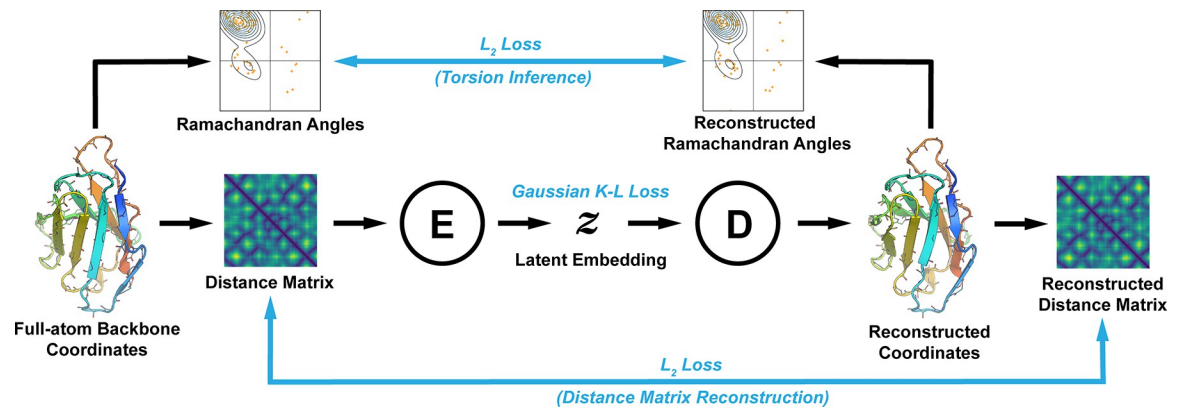


Fig 1. VAE Training Scheme. The flow of data is shown with black arrows, and losses are shown in blue. First, Ramachandran angles and distance matrices are computed from the full-atom backbone coordinates of a training example. The distance matrix is passed to the encoder network (E), which generates a latent embedding that is passed to the decoder network (D). The decoder directly generates coordinates in 3D space, from which the reconstructed Ramachandran angles and distance matrix are computed. Errors from both the angles and distance matrix are back-propagated through the 3D coordinate representation to the encoder and decoder. Note that both the torsion and distance matrix losses are rotationally and translationally invariant, and that the coordinates of the training example are never seen by the model. The shown data are real inputs and outputs of the VAE for the immunoglobulin chain in PDB:4YXH(L).

<https://doi.org/10.1371/journal.pcbi.1010271.g001>

translational invariance in our loss functions to allow for generalizable feature learning across limited datasets. Second, we wanted our model to be torsion-aware, as torsion angles play a crucial role in determining a protein's fold as well as structural quality. Last, we wanted our model to directly generate 3D coordinates to avoid any auxiliary coordinate recovery processes, and to make structure generation end-to-end differentiable. Full differentiability allows us to gradient-optimize a latent vector while subjecting generated structures to arbitrary design constraints. Many iterations of trial-and-error yielded the model training scheme shown in Fig 1, which satisfies all three properties.

Like classical VAEs [25] our model minimizes a reconstruction loss and a KL-divergence loss that constrains the latent embeddings to be isotropic gaussian. During training, a distance matrix is first obtained from the coordinates of the training data. This distance matrix is ordered by residue, and comprised of backbone atoms following the order: N, C α , C, O. Like a conventional VAE, the encoder module compresses the distance matrix into the latent space, and the decoder is tasked with reconstructing the input data. Crucially however, the output of the decoder is not a distance matrix, but rather the 3D coordinates of the reconstructed structure. A reconstructed distance matrix is computed from the coordinates using a differentiable function, and loss is backpropagated from distance matrix to distance matrix, across the coordinate representation. Note that this reconstruction loss does not specify the absolute position of the coordinates and thus preserves the required invariances. Factoring through an explicit coordinate representation allows us to integrate a torsion loss, which we formulate as a supervised-learning objective; the network infers the correct torsion distribution from the distance matrix. This loss is computed via the coordinate representation and backpropagated through to the decoder network, and is also rotationally and translationally invariant.

2.3 Loss-Function

To describe the loss function we adopt the following notation:

x : Coordinate Data.

$z \sim N(0, I)$: Latent Vector.

q_θ : Encoder.

p_ϕ : Decoder.

$x' = p_\phi(q_\theta(x))$: Reconstructed Coordinates.

$z_x = q_\theta(x)$: Latent Embedding of x

The full loss function can be written in classical VAE form as:

$$L_{\theta,\phi} = \frac{1}{n_{\text{data}}} \sum_{x \in \text{data}} \text{ReconLoss}_{\theta,\phi}(x) + \lambda \text{KL}(q_\theta(z|x) || p(z))$$

The reconstruction loss is comprised of distance matrix (*Dist*) and torsion (*Tors*) components:

$$\text{ReconLoss}_{\theta,\phi} = w_{\text{dist}} \text{Dist}(x, x') + w_{\text{tors}} \text{Tors}(x, x')$$

The *Tors* term is computed as the L_2 distance between the unit sphere projections of the real and reconstructed backbone torsion angles. The *Dist* term is comprised of three terms, each of which are L_2 losses between different components of the real and reconstructed distance matrices. Additional details are provided in the S1 Methods in [S1 File](#).

We found that early in training, the torsion loss must be up-weighted relative to the distance loss in order to achieve correct stereochemistry, as molecular handedness cannot be uniquely determined from pairwise distances alone. Decreasing the torsion weight later in training led to improvements in local structure quality. A detailed description of the loss weighting schedule is included in the S1 Methods in [S1 File](#). We note that both the distance and torsion losses are rotationally and translationally invariant, so the absolute position of the output coordinates is determined by the model itself. While coordinates of the training examples are never seen by the model directly, the model learns a natural alignment of structures along the core β -strands.

3. Results

3.1. Overview

To assess the utility our model, we studied the Ig-VAE's performance on several tasks. The first of these is data reconstruction, which reflects the ability of the model to compress structural features into a low-dimensional latent space (Section 3.2). This functionality is an underlying assumption of generative sampling, which requires that the latent space capture the scope of structure variation with sufficient resolution. Next, we assessed the quality and novelty of the generated structures, characterizing the chemical validity of the samples (Section 3.3), while also evaluating the quality of interpolations between embedded structures (Section 3.4). We visualize the distribution of embeddings within the latent space to better understand its structure, and to determine if the sampling distribution is well-supported (Section 3.4). We ultimately challenge the Ig-VAE with a real design task; specifically, generation of a novel backbone with high shape-complementarity to the ACE2 epitope of the SARS-CoV2-RBD [34] (Section 3.5.1). To evaluate the general utility of our approach, we investigated whether

we could leverage the model's generative prior to perform backbone design subject to a set of local, human-specified constraints (Section 3.5.2).

3.2 Structure embedding and reconstruction

A core feature of an effective VAE is the model's ability to embed and reconstruct data. High quality reconstructions indicate that a model is able to capture and compress structural features into a low-dimensional representation, which is a prerequisite for generation by latent space sampling. To evaluate this functionality, we reconstructed 500 randomly selected structures and compared the real and reconstructed distributions of backbone torsion angles, pairwise distances, bond lengths, and bond angles (Fig 2). The structurally padded "dummy" regions were excluded from this analysis.

The distance and torsion distributions are shown in Fig 2A, where we observe that the real and reconstructed data agree well. On average, pairwise distances smaller than 10Å tended to be reconstructed slightly smaller than the actual distances, while larger distances tended to be reconstructed slightly larger (Fig 2B, reconstructed). ϕ and ψ torsions tended within $\sim 10^\circ$ of the real angles, while ω angles tended to fall within $\sim 3^\circ$. Examples of reconstructed backbones are shown in the top row of Fig 2C. These data demonstrate that the Ig-VAE accurately performs full-atom reconstructions over a range of loop conformations. Larger versions of the reconstruction images are provided in S3 Fig.

In order to use generated backbones in conjunction with existing protein design tools, it is crucial that our model produce structures with near-chemically-valid bond lengths and bond angles. Otherwise, large movements in the backbone can occur as a result of energy-based corrections during the design process, leading to the loss of model-generated features. The bond length and bond angle distributions are depicted in Fig 2D. The majority of bond length reconstructions were within ~ 0.1 Å of ideal lengths, while bond angles tended to be within $\sim 10^\circ$ of their ideal angles. We found that a constrained optimization step using the Rosetta centroid-energy function (see S1 Methods in S1 File) could be used to effectively refine the outputs. This refinement process kept structures close to their output conformations (Fig 2C, bottom) while correcting for non-idealities in the bond lengths and bond angles (Fig 2D, green). Refinement did not improve backbone reconstruction accuracy (Fig 2B, refined), but did improve chemical validity, implying that our model outputs could be refined with Rosetta without washing-out generated structural features. An additional analysis comparing errors between the β - and loop-regions is provided in S1 Fig.

Our analysis of the reconstructions reveal that the Ig-VAE can be used to obtain high-resolution structure embeddings that are likely useful in various learning tasks on 3D protein data. These results support our later conclusion that the Ig-VAE embedding space can be leveraged to generate structures with high atomic precision, while also showing that the KL-regularization imposed on the latent embeddings does not overpower the autoencoding functionality of the VAE.

3.3 Structure generation

To determine whether the Ig-VAE could generate novel, realistic Ig backbones, we sampled 500 structures from the latent space of the model and compared their feature distributions to 500 non-redundant structures from the dataset. Each generated structure was cropped based on its nearest neighbor in the training set (See S1 Methods in S1 File). In Fig 3A we show overlays of the distance and torsion distributions for the real and generated structures. The generated torsions were more variable than the real torsions, with more residues falling outside the range of the training data. The real and generated distance distributions agreed well. We note

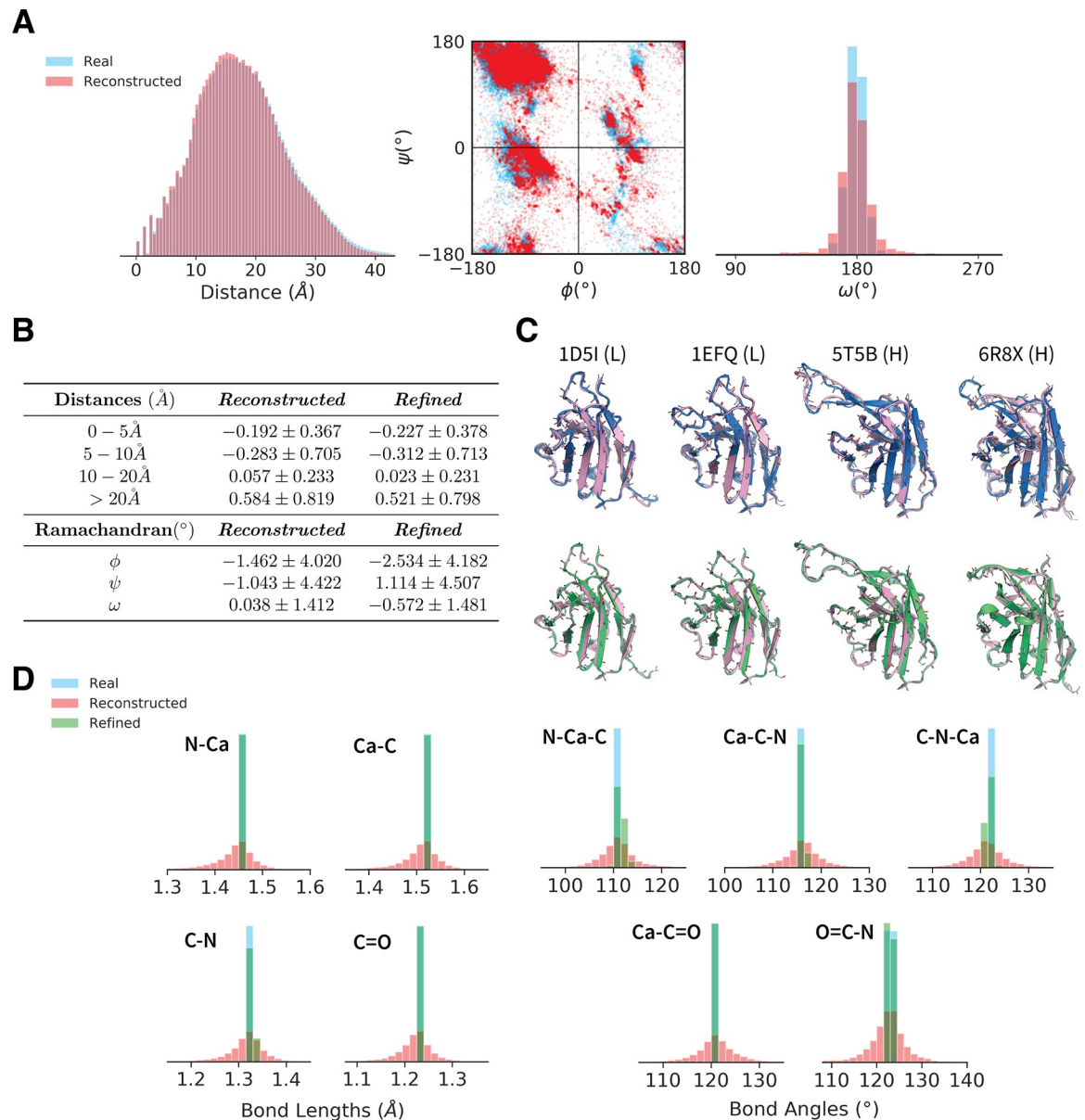


Fig 2. Analysis of Full-Atom Reconstructions. Reconstruction data for 500 randomly chosen, non-redundant structures in the training set. (A) Overlays of the pairwise distance and Ramachandran distributions of the real and reconstructed data. (B) A table of the reconstruction errors in pairwise distance and Ramachandran angle before and after refinement. The distance errors are reported as per-pairwise-distance error averaged over all structures in the dataset, and analogously for the angle errors. (C) Overlays of the real (blue), reconstructed (pink) and refined (green) structures. (D) Overlays of the bond length and bond angle distributions of the real, reconstructed and refined data. Overall, structures are accurately reconstructed, and errors in atom placement are small enough that they can be corrected with minimal changes to the model outputs.

<https://doi.org/10.1371/journal.pcbi.1010271.g002>

that in the $<5\text{\AA}$ regime there are several spikes in the distance distribution found in the real, reconstructed, and generated distributions (Figs 2A, 3A, S1 and S2). These spikes are not artifacts, but rather characteristic regularities found in real proteins that serve as good indicators of structure resolution. For example, in real structures $C\alpha$ atoms are found at regular intervals of $\sim 3.8\text{\AA}$, and $C\alpha$ -N (non-bonded) atoms are always $\sim 2.4\text{\AA}$ apart. Bond lengths, β -strands, helices are also examples of recurring-distance structures in this region.

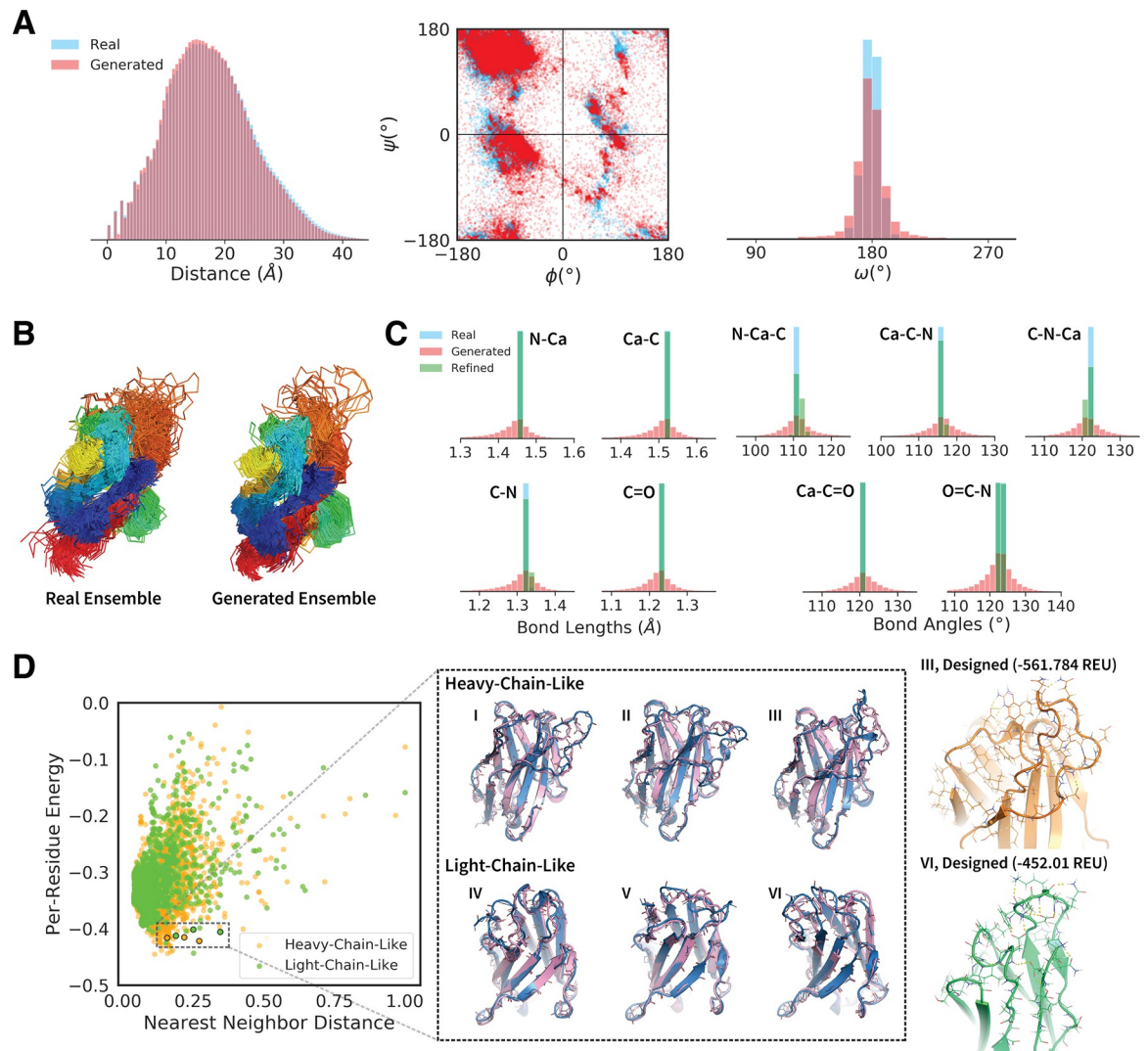


Fig 3. Analysis of Generative Sampling. Data for 500 of randomly selected non-redundant training samples and generated structures. (A) Overlays of the pairwise distance and Ramachandran distributions of the real and generated data. (B) A comparison of the real and generated structural ensembles. (C) Overlays of the bond length and bond angle distributions of the real, generated and refined data. (D) The left panel shows a plot of the post-refinement per-residue centroid energy against normalized nearest neighbor distance for the generated structures. The nearest neighbor distance is computed as the minimum Frobenius distance between the generated distance matrix and all distance matrices in the training set. Each point is colored based on whether the nearest neighbor is a heavy or light chain Ig. The center panel shows an overlay of the generated structures (pink) and their nearest neighbors (blue) in the training set. These six structures were selected using a combination of centroid energy, nearest neighbor distance, heavy/light classification, and manual inspection. The right panel shows sequence design results for structures III and VI. The energies in the left panel are centroid energies, while the energies in the right panel are full-atom Rosetta energies using the ref2015 score function.

<https://doi.org/10.1371/journal.pcbi.1010271.g003>

Visual assessment of the backbone ensembles (Fig 3B) revealed that the two were similarly variable, suggesting that our model captures much of the structural variation found in the training set. The generated structures exhibited good chemical-bond geometries (Fig 3C, red) that were slightly noisier but comparable to those of the reconstructed backbones (Fig 2D). Once again, we found that constrained refinement using Rosetta could improve chemical bond geometries (Fig 3C, green), with minimal changes to the generated

structures. An additional analysis comparing errors between the β - and loop-regions is provided in [S2 Fig](#).

To assess both the novelty and viability of the generated examples, we evaluated structures based on two criteria: (1) post-refinement energy and (2) nearest-neighbor distance. Energies were normalized by residue-count to account for variable structure sizes. Nearest-neighbor distance was computed as a length-normalized Frobenius distance between $C\alpha$ -distance matrices. For notational convenience, nearest-neighbor distances were normalized between 0 and 1. We avoid the use of the classical $C\alpha$ -RMSD, because it is neither an alignment-free nor a length-invariant metric, and because it lacks sufficient precision to make meaningful comparisons between Ig structures. Nearest neighbors are identified using $C\alpha$ distances instead of full-atom (N, $C\alpha$, C, O) distances, because $C\alpha$ distance tends to underestimate uniqueness, making our uniqueness check more stringent. In the full-atom case, small errors over more atoms tend to add up and inflate overall distance.

We found that there was a positive correlation between energy and nearest-neighbor distance ([Fig 3D](#)), implying that while our model is able to generate structures that differ from any known examples, there is a concurrent degradation in quality when structures drift too far from the training data. Despite this, we found that a significant number of generated structures had novel loop shapes, achieved favorable energies, and retained Ig-specific structural features. Six of these examples are shown as raw model outputs, overlaid with their nearest neighbors in the center panel of [Fig 3D](#). Both heavy and light-chain-like structures exhibited dynamic loop structures, and the model appears to perform well in generating both long and short loops. Larger images of the generated examples are provided in [S4 Fig](#).

To assess whether the generated loop conformations could be sustained by an amino acid sequence, we used Rosetta FastDesign [[35,36](#)] to create full-length sequences for the selected backbones. The outputs of two representative design trajectories are shown in the right panel of [Fig 3D](#). The design process yielded energetically favorable sequences with loops supported by features such as hydrophobic packing, π - π stacking and hydrogen bonding. Overall these results suggest that the Ig-VAE is capable of generating novel, high-quality backbones that are chemically accurate, and that can be used in conjunction with existing design protocols to obtain biochemically realistic sequences.

3.4 Latent space analysis and interpolation

While the results of the preceding section suggest that the Ig-VAE is able to produce novel structures, an important feature of any generative model is the ability to interpolate smoothly between examples in the latent space. In design applications this functionality allows for dense structural sampling, and modeling of transitions between distinct structural features.

A linear interpolation between two randomly selected embeddings is shown in [Fig 4A](#). The majority of interpolated structures adopt realistic conformations, retaining characteristic backbone hydrogen bonds while transitioning smoothly between different loop conformations ([Fig 4C](#)). Structures along the interpolation trajectory were able to achieve negative post-refinement energies ([Fig 4B](#)), with the highest energy structure corresponding to the most unrealistic portion of the trajectory ([Fig 4A and 4B](#), index 20).

To better understand the structure of the embedding space, we visualized the training data embeddings ([Fig 4D](#)) using two dimensionality-reduction methods: t-distributed stochastic neighbor embedding (tSNE) [[37](#)] and principal components analysis (PCA) [[38](#)]. The top panel depicts a tSNE decomposition of the embedding means (without variance) for the 4154 non-redundant structures in the dataset. K-means clustering ($k = 40$) revealed distinct clusters that roughly correlated with loop structure, suggesting a correspondence between latent space

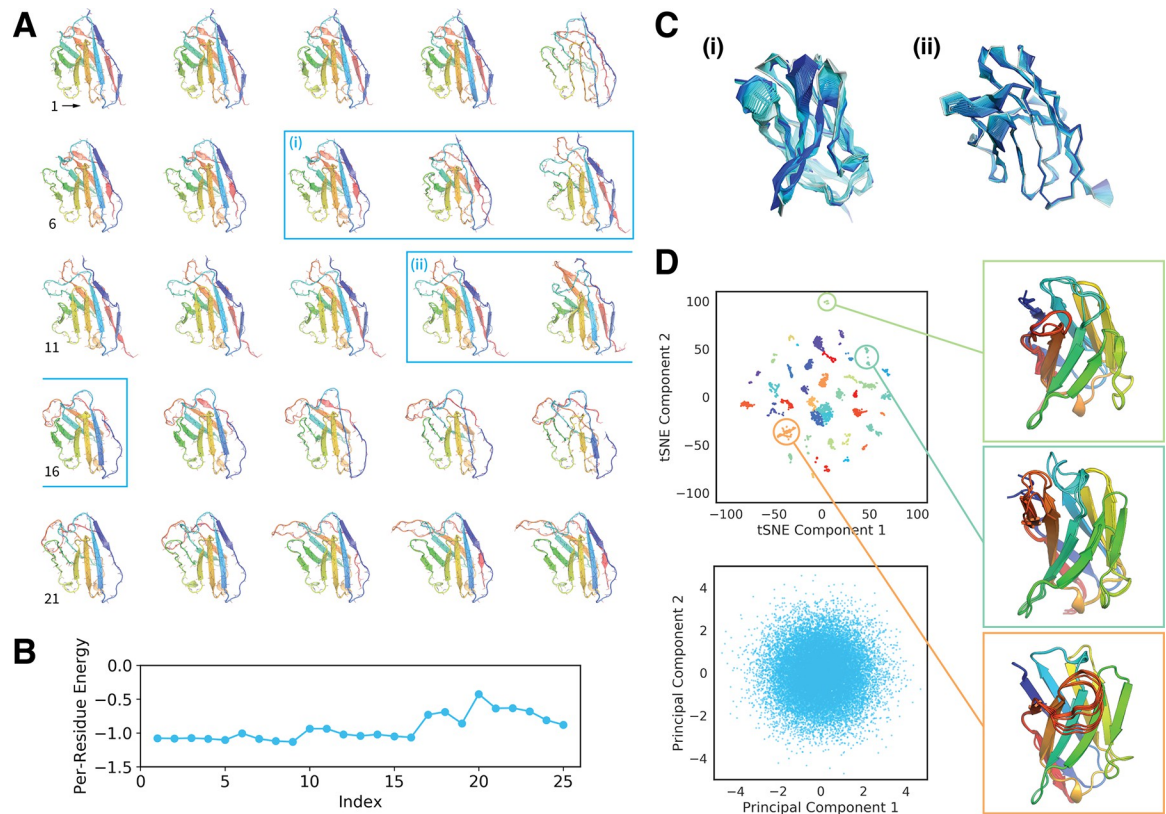


Fig 4. Latent Space Analysis and Interpolation. (A) Linear interpolation between two randomly selected embeddings. The starting and ending structures are 1TQB(H) and 5JW4(H) respectively. The backbones are unaltered, full-atom model outputs. Structures are colored by residue index in reverse rainbow order. (B) Centroid energy profiles of the structures from panel A after constrained refinement. (C) Higher frequency overlays of 80 sequential structures in the unrefined interpolation trajectory. The roman numerals correspond to the blue labels in panel A. A lighter shade of blue indicates an earlier structure while the darker shade indicates a later structure. Structure transitions are smooth and follow a near-continuous trajectory. (D) The top panel shows a tSNE dimensionality reduction of the embedding means of the 4154 non-redundant structures in the training set. Colorings correspond to k-means clusters ($k = 40$) of the post-tSNE data, and ten structures from three clusters are visualized to the right. The bottom panel depicts a principal components reduction of the latent space, showing five sampled data points per non-redundant structure.

<https://doi.org/10.1371/journal.pcbi.1010271.g004>

position and semantically meaningful features (Fig 4D, insert). In the bottom panel, we visualized sampled embeddings for the same set of structures, sampling 5 embeddings per example. PCA revealed a spherical, densely populated embedding space, suggesting that the isotropic gaussian sampling distribution is well supported. The PCA results also suggest that the KL-loss was sufficiently weighted during training.

Overall these results support the conclusions of the previous section, demonstrating that the Ig-VAE exhibits the features expected of a properly-functioning generative model, and that sampling from a gaussian prior is well motivated. The smooth interpolations agree with the observation that our model is capable of generating novel structures, which are expected to arise by sampling from interpolated regions between the various embeddings.

3.5 Towards epitope-specific generative design

3.5.1 Computational design of a SARS-CoV2 binder. While antibodies are usually comprised of two Ig domains, there also exist a large number of single-domain antibodies in the form of camelids [39] and Bence Jones proteins [40]. To test the utility of our model in a real-

world design problem, we challenged the Ig-VAE to generate a single-domain-binder to the ACE2 epitope of the SARS-CoV2 receptor binding domain (RBD), an epitope which is of significant interest in efforts towards resolving the 2019/20 coronavirus pandemic [34].

To do this, we first generated 5000 backbones by random sampling of the latent space distribution. To find candidates with high shape complementarity to the ACE2 epitope, we used PatchDock [41] to dock each generated structure against the RBD. To make the search sequence agnostic, both proteins were simulated as poly-valines during this step. We then selected Ig's that bound the ACE2 epitope specifically, and used FastDesign to optimize the sequences of the binding interfaces, and design full-length sequences across the entire protein. Two Ig's that exhibited good shape complementarity to the ACE2 epitope and adopted unique loop conformations are shown in Fig 5A. After sequence design these candidates achieved favorable energies and complex ddG's of -37.6 and -53.1 Rosetta Energy Units (Fig 5A, designed). Using RosettaDock [42,43], we were able to accurately recover the designed interfaces as the energy minimum of a blind global docking trajectory, suggesting that the binders are specific to their cognate epitopes (Fig 5A, recovered, docking).

These results demonstrate that the latent space of our generative model can be leveraged to create novel binding proteins that are otherwise unobtainable by discrete sampling of real structures. While we believe that designed proteins must be experimentally validated, our data suggest that generative models can provide compelling design candidates by computational design standards, making the method worthy of larger-scale and broader experimental testing. Our data also demonstrate the compatibility of the Ig-VAE with established design suites like Rosetta, which have conventionally relied on real proteins as templates.

3.5.2 Generative design. The functional elements of a protein are often localized to specific regions. Antibodies are one example of this, where binding is attributable to a set of surface-localized loops, as well as enzymes, which depend on the positioning of catalytic residues to form an active site. Despite this apparent simplicity, natural proteins carry evolution-optimized features that are required to host functional elements. The protein design process often seeks to mimic this organization, requiring the engineering of supporting elements centered around a desired feature. While logical, designing supportive features is almost always a difficult task, requiring large amounts of experience and manual tuning.

Motivated by these difficulties, we sought to investigate whether the generative prior of our model could be leveraged to create structures that conform to a human-specified feature without specification of other supporting features. To test this, we specified a 12-residue antibody loop shape as pairwise Ca distances. We then sampled 100 random latent initializations and applied the constraints to the generated structures. Next, we optimized the structures via gradient descent, backpropagating constraint errors to the latent vectors through the decoder network. From 100 initializations, we were able to recover the target loop shape in 62 trajectories, with the vast majority of structures retaining high quality, realistic features. We visualize one trajectory in the center panel of Fig 5B. While the middle Ig-loop (Fig 5B, blue) is being constrained, the other loops move to adopt sterically compatible conformations and the angles of the β -strands change to support the new loop shapes. We note that the latent-vector optimization problem is non-convex, which is why we require multiple random initializations [44,45].

Importantly, the recovered backbones in the generated ensemble differ from the originating structure (Fig 5B, orange). These data suggest that our model can be used to create backbones that satisfy specific design constraints, while also providing a distribution of compatible supporting elements. We emphasize that this procedure is not limited to distance constraints, and can be done using any differentiable coordinate-based heuristic such as shape complementarity [46], volume constraints, Rosetta energy [47], and more. With a well-formulated loss function, which warrants a study in itself, it is possible to "mold" the loops of an antibody to a

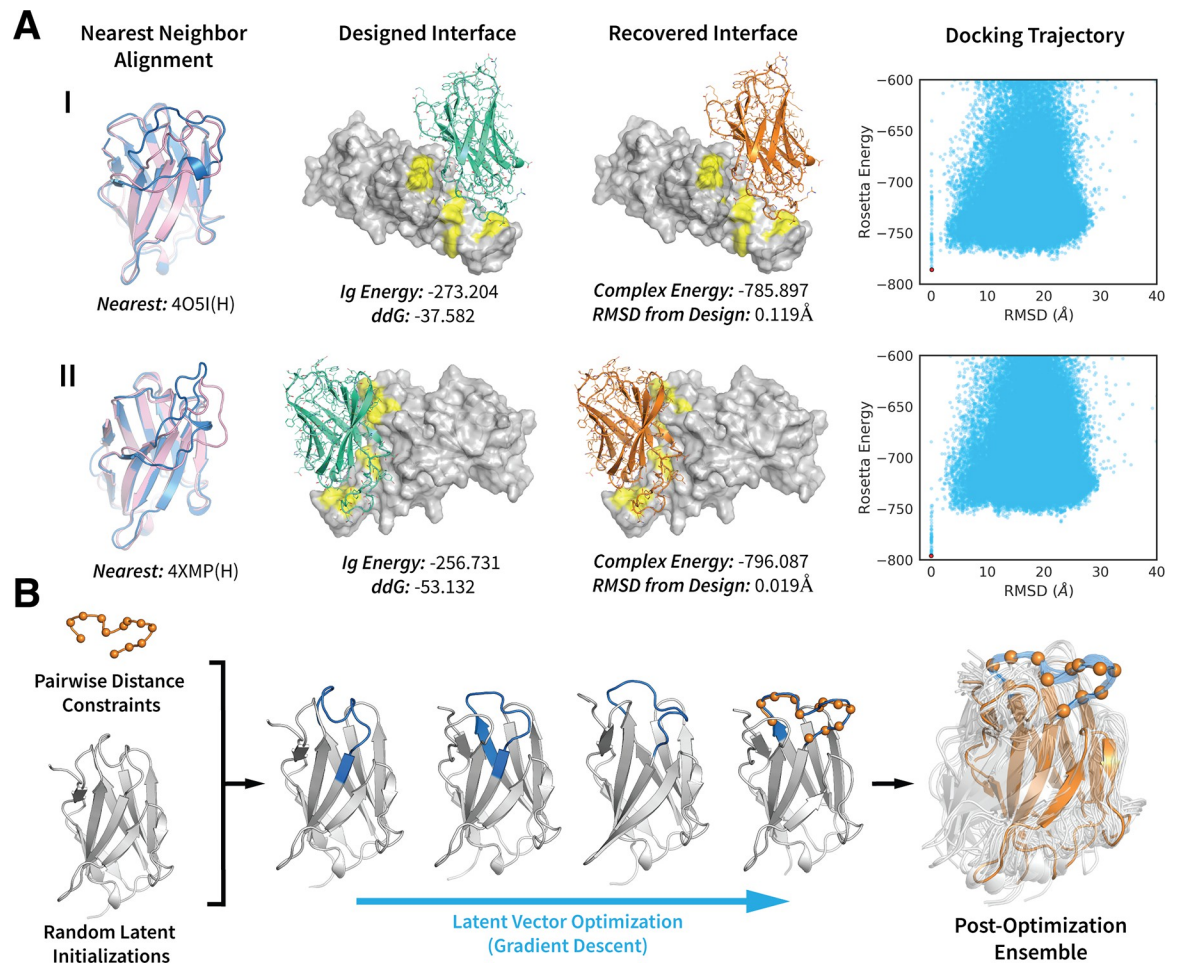


Fig 5. Towards Epitope-Specific Generative Design. (A) Design of two generated immunoglobulin (Ig) structures, I and II, targeting the ACE2 epitope of the SARS-CoV2 receptor binding domain (RBD). The left-most column shows an alignment of the unrefined generated structures (pink) to their respective nearest neighbors (blue) in the training set. The second column depicts the designed interfaces with the full-sequence Ig's shown in green. The SARS-Cov2-RBD is rendered as a white surface, with the ACE2 epitope shown in yellow. The ddG's and Rosetta energies (ref2015) of the Ig's are shown below each complex. The third column depicts the lowest energy structures from a RosettaDock global-docking trajectory. The recovered Ig structures are shown in orange. The right-most column shows the full docking trajectory (100,000 decoys), with the lowest energy structures shown as red dots. (B) Ig-backbone design by constrained latent vector optimization. The loop constraints are taken from PDB:5X2O(L, 32–44) with the target shape shown by orange spheres. The constrained regions of the optimization trajectory are shown in blue. In the post-optimization ensemble the full 5X2O(L) backbone is shown in orange. The ensemble depicts the outputs of 62 optimization trajectories (out of 100 initialized) that successfully recovered the target loop shape.

<https://doi.org/10.1371/journal.pcbi.1010271.g005>

target epitope, or even the backbone of an enzyme around a substrate of choice. Our example demonstrates a novel formulation of protein design as a constrained optimization problem in the latent space of a generative model. In contrast to methods that require manual curation of each part of a protein, we term our approach “generative design,” where the requirement of human-specified heuristics constitutes the “design” element, and using a generative model to fill in the details of the structure constitutes the “generative” element.

4. Discussion

While a protein's fold dictates the range of functions that it can host, both sequence features and structural features within a fold are crucial in dictating function. In the case of antibodies,

structural variation in the loop regions plays an especially important role, often determining target compatibility and specificity [48,49]. Although fragment sampling of existing structures has allowed us to engineer a wide range of folds [1,2,3], it remains a challenge to reliably generate novel Ig-loops outside the realm of known structures, making binder design for novel targets difficult.

In this study we have approached this problem by training a deep generative model that generates Ig structures with novel loop shapes, drawing only from existing structural data. We demonstrate that the structural distribution captured by our VAE can be used for backbone design, which we pose as a conditional generation problem solved via constrained optimization in the model's latent space. Specifically, the generative prior provided by the Ig-VAE restricts backbone conformations to the regime of the Ig-fold, and the latent space search effectively generates structures conforming to user-specified features. This scheme, which we term "generative design", is possible because our model generates 3D coordinates directly, as a result of innovations that allow for rotational and translational invariance, and torsional awareness during training.

We speculate that histogram and distribution-based representations used in structure-prediction methods such as trRosetta [15] and DeepAb [50] can be also used for training. These distributions are used to capture uncertainty that is important in structure prediction, while in VAE training, where data compression/reconstruction is the goal, the noise in this representation may become a disadvantage due to lack of precision. In addition, when building structures with distribution-representations [51], the process often requires some form of stochastic recovery, for example using Rosetta, making it unsuited towards design tasks requiring optimization by backpropagation through 3D coordinates. Graphs have also been used for CDR3 design via joint iterative prediction of sequence and structural coordinates [52]. Compared to our method which adheres to the *de novo* approach of full-backbone generation and backbone-conditioned sequence design, graphs appear promising for tasks that require local redesign of existing structures where the context of a pre-built backbone is important. Such methods may be used in conjunction with our own model to improve local sampling.

All in all, the Ig-VAE yields a powerful tool for creating single-domain antibodies, and allows for high-throughput construction of epitope-tailored, structure-guided libraries. With such a tool, it may be possible to circumvent screening of fully randomized libraries, a large proportion of which are usually insoluble or fundamentally incompatible with the target of interest [53,54,55,56]. Importantly, our approach is not specific to Ig's, and can be applied to any fold-class well represented in structural databases such as enzymes, nucleic acid binding domains and more.

Conceptually, our study offers important insight into the massive possibilities of deep generative modeling in protein design. Generative design provides a new approach that is highly flexible, allowing for the arbitrary specification of locally desired features, while leveraging the model to "fill in" the ancillary parts of the molecule. This functionality is mathematically grounded in the framework of conditional probability, making it fully data-guided while requiring minimal human intuition compared to conventional backbone design methods.

Overall, our work is of significant interest to both protein engineers and machine learning scientists as the first successful example of 3D deep generative modeling applied to protein backbone design. We speculate that our scheme will motivate further study of class-specific generative models, as well as development of differentiable loss functions that can be used, for example, to morph enzymes to host small molecule binding sites, or even mold the loops of antibodies to the surfaces of targets.

Supplemental information

Details regarding model training, dataset construction, and interface design are provided in the Supplemental Information. Rosetta commands and interpolation movies are provided in [S1 File](#).

Supporting information

S1 File. *interpolation.mov*. mov format of the interpolation movie for the trajectory analyzed in [Fig 4](#). *interpolation.mp4*. mp4 format of the interpolation movie for the trajectory analyzed in [Fig 4](#). *patchdock_output_I.pdb*. The raw patchdock output used for design I in [Fig 5](#). *patchdock_output_II.pdb*. The raw patchdock output used for design II in [Fig 5](#). **S1 Methods**.

Detailed documentation of model training, architecture, and all design protocols used in the study. *supplemental_methods_attachment.zip*. Scripts and commands used for design. Referenced by *supplemental_methods_rev.pdf*.

(ZIP)

S1 Fig. Comparison of Reconstruction Behavior in Beta and Loop Regions. Reconstruction data for 500 randomly chosen, non-redundant structures in the training set, plotted separately for the beta and loop regions. The top row in each section shows, overlays of the pairwise distance and Ramachandran distributions of the real and reconstructed data. The bottom rows show overlays of the bond length and bond angle distributions of the real, reconstructed and refined data. Real data are shown in blue, reconstructed in pink, and refined in green. In the bottom right of each section we show a table of the reconstruction errors in pairwise distance and Ramachandran angle before and after refinement. The distance errors are reported as per-pairwise-distance error averaged over all structures in the dataset, and analogously for the angle errors.

(TIF)

S2 Fig. Comparison of Generation Behavior in Beta and Loop Regions. Data for 500 of randomly selected non-redundant training samples and generated structures, plotted separately for the beta and loop regions. The top row in each section shows, overlays of the pairwise distance and Ramachandran distributions of the real and reconstructed data. The bottom rows show overlays of the bond length and bond angle distributions of the real, reconstructed and refined data. Real data are shown in blue, reconstructed in pink, and refined in green.

(TIF)

S3 Fig. Larger Rendering of Reconstructed Examples. A higher resolution rendering of the reconstruction examples shown in [Fig 2C](#) for two heavy-chain and two light-chain examples. Real structures are shown in blue, reconstructed structures are shown in pink, and refined structures are shown in green.

(TIF)

S4 Fig. Larger Rendering of Generated Examples. A higher resolution rendering of the generated examples shown in [Fig 3D](#) (Center), selected based on the analysis in [Fig 3D](#) (Left), for two heavy-chain and two light-chain examples. Real structures are shown in blue, reconstructed structures are shown in pink, and refined structures are shown in green. Each alignment shows an overlay of the generated structures (pink) and their nearest neighbors (blue) in the training set. The roman numerals correspond to the same structures as in [Fig 3D](#).

(TIF)

Acknowledgments

We thank Sergey Ovchinnikov for helpful discourse during early phases of this project, and for contributing initial code that became part of the torsion-reconstruction loss function. We thank Namrata Anand for helpful feedback and assistance.

Author Contributions

Conceptualization: Raphael R. Eguchi, Po-Ssu Huang.

Data curation: Raphael R. Eguchi.

Formal analysis: Raphael R. Eguchi, Christian A. Choe.

Investigation: Raphael R. Eguchi.

Methodology: Raphael R. Eguchi.

Project administration: Po-Ssu Huang.

Resources: Raphael R. Eguchi, Po-Ssu Huang.

Supervision: Po-Ssu Huang.

Validation: Raphael R. Eguchi, Christian A. Choe.

Visualization: Raphael R. Eguchi.

Writing – original draft: Raphael R. Eguchi, Po-Ssu Huang.

Writing – review & editing: Raphael R. Eguchi, Po-Ssu Huang.

References

1. Dou J, Vorobieva AA, Sheffler W, Doyle LA, Park H, Bick MJ, et al. De novo design of a fluorescence-activating β -barrel. *Nature*. 2018 Sep; 561(7724):485–91. Available from: <https://doi.org/10.1038/s41586-018-0509-0> PMID: 30209393
2. Huang P, Feldmeier K, Parmeggiani F, Velasco DF, Höcker B, Baker D. De novo design of a four-fold symmetric TIM-barrel protein with atomic-level accuracy. *Nature Chemical Biology*. 2015; 12(1):29–34. Available from: https://www.bakerlab.org/wp-content/uploads/2015/12/Huang_NatChemBio_2015.pdf. <https://doi.org/10.1038/nchembio.1966> PMID: 26595462
3. Silva DA, Yu S, Ulge UY, Spangler JB, Jude KM, Labão-Almeida C, et al. De novo design of potent and selective mimics of IL-2 and IL-15. *Nature*. 2019 Jan; 565(7738):186–91. Available from: <https://doi.org/10.1038/s41586-018-0830-7> PMID: 30626941
4. Pejchal R, Doores KJ, Walker LM, Khayat R, Huang PS, Wang SK, et al. A Potent and Broad Neutralizing Antibody Recognizes and Penetrates the HIV Glycan Shield. *Science*. 2011 Nov; 334(6059):1097–103. Available from: <http://www.sciencemag.org/cgi/doi/10.1126/science.1213256>. PMID: 21998254
5. Röthlisberger D, Khersonsky O, Wollacott AM, Jiang L, DeChancie J, Betker J, et al. Kemp elimination catalysts by computational enzyme design. *Nature*. 2008 May; 453(7192):190–5. Available from: <http://www.nature.com/doi/10.1038/nature06879>. PMID: 18354394
6. Langan RA, Boyken SE, Ng AH, Samson JA, Dods G, Westbrook AM, et al. De novo design of bioactive protein switches. *Nature*. 2019 Aug; 572(7768):205–10. Available from: <http://www.nature.com/articles/s41586-019-1432-8>. <https://doi.org/10.1038/s41586-019-1432-8> PMID: 31341284
7. Wei KY, Moschidi D, Bick MJ, Nerli S, McShan AC, Carter LP, et al. Computational design of closely related proteins that adopt two well-defined but structurally divergent folds. *Proceedings of the National Academy of Sciences*. 2020 Mar; 117(13):7208–15. Available from: <http://www.pnas.org/lookup/doi/10.1073/pnas.1914808117>. PMID: 32188784
8. Lu P, Min D, DiMaio F, Wei KY, Vahey MD, Boyken SE, et al. Accurate computational design of multi-pass transmembrane proteins. *Science*. 2018 Mar; 359(6379):1042–6. Available from: <https://www.sciencemag.org/lookup/doi/10.1126/science.aag1739>. PMID: 29496880
9. Huang PS, Boyken SE, Baker D. The coming of age of de novo protein design. *Nature*. 2016 Sep; 537(7620):320–7. Available from: <http://www.nature.com/articles/nature19946>. <https://doi.org/10.1038/nature19946> PMID: 27629638

10. Leaver-Fay A, Tyka M, Lewis SM, Lange OF, Thompson J, Jacak R, et al. Rosetta3: An Object-Oriented Software Suite for the Simulation and Design of Macromolecules. In: Johnson ML, Brand L, editors. *Computer Methods, Part C*. vol. 487 of *Methods in Enzymology*. Academic Press; 2011. p. 545–574. Available from: <http://www.sciencedirect.com/science/article/pii/B9780123812704000196>.
11. Anand N, Eguchi RR, Derry A, Altman RB, Huang PS. Protein Sequence Design with a Learned Potential. *Bioinformatics*; 2020. Available from: <http://biorxiv.org/lookup/doi/10.1101/2020.01.06.895466>.
12. Eguchi RR, Huang PS. Multi-scale structural analysis of proteins by deep semantic segmentation. *Bioinformatics*. 2020 Mar; 36(6):1740–9. Available from: <https://academic.oup.com/bioinformatics/article/36/6/1740/5551337>. <https://doi.org/10.1093/bioinformatics/btz650> PMID: 31424530
13. Gainza P, Sverrisson F, Monti F, Rodolà E, Boscaini D, Bronstein MM, et al. Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning. *Nature Methods*. 2019 Dec. Available from: <http://www.nature.com/articles/s41592-019-0666-6>. <https://doi.org/10.1038/s41592-019-0666-6> PMID: 31819266
14. Senior AW, Evans R, Jumper J, Kirkpatrick J, Sifre L, Green T, et al. Improved protein structure prediction using potentials from deep learning. *Nature*. 2020 Jan; 577(7792):706–10. Available from: <http://www.nature.com/articles/s41586-019-1923-7>. <https://doi.org/10.1038/s41586-019-1923-7> PMID: 31942072
15. Yang J, Anishchenko I, Park H, Peng Z, Ovchinnikov S, Baker D. Improved protein structure prediction using predicted interresidue orientations. *Proceedings of the National Academy of Sciences*. 2020; 117(3):1496–503. Publisher: National Academy of Sciences _eprint: <https://www.pnas.org/content/117/3/1496.full.pdf>. Available from: <https://www.pnas.org/content/117/3/1496>. <https://doi.org/10.1073/pnas.1914677117> PMID: 31896580
16. Alley EC, Khimulya G, Biswas S, AlQuraishi M, Church GM. Unified rational protein engineering with sequence-based deep representation learning. *Nature Methods*. 2019 Dec; 16(12):1315–22. Available from: <http://www.nature.com/articles/s41592-019-0598-1>. <https://doi.org/10.1038/s41592-019-0598-1> PMID: 31636460
17. Madani A, McCann B, Naik N, Keskar NS, Anand N, Eguchi RR, et al. ProGen: Language Modeling for Protein Generation. *Synthetic Biology*; 2020. Available from: <http://biorxiv.org/lookup/doi/10.1101/2020.03.07.982272>.
18. Riesselman AJ, Ingraham JB, Marks DS. Deep generative models of genetic variation capture the effects of mutations. *Nature Methods*. 2018 Oct; 15(10):816–22. Available from: <https://doi.org/10.1038/s41592-018-0138-4>. PMID: 30250057
19. Riesselman A, Shin JE, Kollasch A, McMahon C, Simon E, Sander C, et al. Accelerating Protein Design Using Autoregressive Generative Models. *bioRxiv*. 2019. Available from: <https://www.biorxiv.org/content/early/2019/09/05/757252>.
20. Anand N, Huang P. Generative modeling for protein structures. In: Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R, editors. *Advances in Neural Information Processing Systems* 31. Curran Associates, Inc.; 2018. p. 7494–505. Available from: <http://papers.nips.cc/paper/7978-generative-modeling-for-protein-structures.pdf>.
21. Anand N, Eguchi RR, Huang PS. Fully differentiable full-atom protein backbone generation. In: DGS@ICLR; 2019.
22. Ovchinnikov S, Huang PS. Structure-based protein design with deep learning. *Current Opinion in Chemical Biology*. 2021; 65:136–44. *Mechanistic Biology * Machine Learning in Chemical Biology*. Available from: <https://www.sciencedirect.com/science/article/pii/S1367593121001125>. <https://doi.org/10.1016/j.cbpa.2021.08.004> PMID: 34547592
23. Huang H, Amor BB, Lin X, Zhu F, Fang Y. G-VAE, a Geometric Convolutional VAE for Protein Structure Generation. *arXiv*; 2021. Available from: <https://arxiv.org/abs/2106.11920>.
24. Guo X, Tadepalli S, Zhao L, Shehu A. Generating Tertiary Protein Structures via an Interpretative Variational Autoencoder. *arXiv:200407119 cs, q-bio, stat*. 2020 Apr. *ArXiv: 2004.07119*. Available from: <http://arxiv.org/abs/2004.07119>.
25. Kingma DP, Welling M. Auto-Encoding Variational Bayes. *arXiv:13126114 cs, stat*. 2014 May. *ArXiv: 1312.6114*. Available from: <http://arxiv.org/abs/1312.6114>.
26. Hoffmann M, Noé F. Generating valid Euclidean distance matrices. *arXiv:191003131 cs, stat*. 2019 Nov. *ArXiv: 1910.03131*. Available from: <http://arxiv.org/abs/1910.03131>.
27. Boder ET, Midelfort KS, Wittrup KD. Directed evolution of antibody fragments with monovalent femtomolar antigen-binding affinity. *Proceedings of the National Academy of Sciences*. 2000 Sep; 97(20):10701–5. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.170297297>. PMID: 10984501
28. Streltsov VA, Carmichael JA, Nuttall SD. Structure of a shark IgNAR antibody variable domain and modeling of an early-developmental isotype. *Protein Science*. 2005 Nov; 14(11):2901–9. Available from: <http://doi.wiley.com/10.1110/ps.051709505>. PMID: 16199666

29. Huo J, Le Bas A, Ruza RR, Duyvesteyn HME, Mikolajek H, Malinauskas T, et al. Neutralizing nanobodies bind SARS-CoV-2 spike RBD and block interaction with ACE2. *Nature Structural & Molecular Biology*. 2020 Jul. Available from: <http://www.nature.com/articles/s41594-020-0469-6>.
30. Flajnik MF, Deschacht N, Muyldermans S. A Case Of Convergence: Why Did a Simple Alternative to Canonical Antibodies Arise in Sharks and Camels? *PLoS Biology*. 2011 Aug; 9(8):e1001120. Available from: <https://doi.org/10.1371/journal.pbio.1001120> PMID: 21829328
31. Huang PS, Ban YEA, Richter F, Andre I, Vernon R, Schief WR, et al. RosettaRemodel: A Generalized Framework for Flexible Backbone Protein Design. *PLoS ONE*. 2011 Aug; 6(8):e24109. Available from: <https://doi.org/10.1371/journal.pone.0024109> PMID: 21909381
32. Koga N, Tatsumi-Koga R, Liu G, Xiao R, Acton TB, Montelione GT, et al. Principles for designing ideal protein structures. *Nature*. 2012 Nov; 491(7423):222–7. Available from: <https://doi.org/10.1038/nature11600> PMID: 23135467
33. Ferdous S, Martin ACR. AbDb: antibody structure database—a database of PDB-derived antibody structures. *Database*. 2018 Jan;2018. Available from: <https://academic.oup.com/database/article/doi/10.1093/database/bay040/4989324>. PMID: 29718130
34. Liu M, Wang T, Zhou Y, Zhao Y, Zhang Y, Li J. Potential role of ACE2 in coronavirus disease 2019 (COVID-19) prevention and management. *Journal of Translational Internal Medicine*. 2020 May; 8(1):9–19. Available from: <https://content.sciendo.com/view/journals/jtim/8/1/article-p9.xml>. <https://doi.org/10.2478/jtim-2020-0003> PMID: 32435607
35. Bhardwaj G, Mulligan VK, Bahl CD, Gilmore JM, Harvey PJ, Cheneval O, et al. Accurate de novo design of hyperstable constrained peptides. *Nature*. 2016 Oct; 538(7625):329–35. Available from: <http://www.nature.com/articles/nature19791>. <https://doi.org/10.1038/nature19791> PMID: 27626386
36. Fleishman SJ, Leaver-Fay A, Corn JE, Strauch EM, Khare SD, Koga N, et al. RosettaScripts: A Scripting Language Interface to the Rosetta Macromolecular Modeling Suite. *PLoS ONE*. 2011 Jun; 6(6):e20161. Available from: <https://doi.org/10.1371/journal.pone.0020161> PMID: 21731610
37. van der Maaten L, Hinton G. Visualizing Data using t-SNE. *Journal of Machine Learning Research*. 2008; 9:2579–605. Available from: <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
38. KPFR S. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*. 1901; 2(11):559–72. Available from: <https://doi.org/10.1080/14786440109462720>.
39. Arbabi-Ghahroudi M. Camelid Single-Domain Antibodies: Historical Perspective and Future Outlook. *Frontiers in Immunology*. 2017 Nov;8. Available from: <http://journal.frontiersin.org/article/10.3389/fimmu.2017.01589/full>. <https://doi.org/10.3389/fimmu.2017.01589> PMID: 29209322
40. Preud'homme JL. Bence Jones Proteins. In: Delves PJ, editor. *Encyclopedia of Immunology* (Second Edition). second edition ed. Oxford: Elsevier; 1998. p. 341 342. Available from: <http://www.sciencedirect.com/science/article/pii/B0122267656000931>.
41. Schneidman-Duhovny D, Inbar Y, Nussinov R, Wolfson HJ. PatchDock and SymmDock: servers for rigid and symmetric docking. *Nucleic Acids Research*. 2005 Jul; 33(Web Server):W363–7. Available from: <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gki481>. PMID: 15980490
42. Chaudhury S, Berrondo M, Weitzner BD, Muthu P, Bergman H, Gray JJ. Benchmarking and Analysis of Protein Docking Performance in Rosetta v3.2. *PLoS ONE*. 2011 Aug; 6(8):e22477. Available from: <https://doi.org/10.1371/journal.pone.0022477> PMID: 21829626
43. Gray JJ, Moughon S, Wang C, Schueler-Furman O, Kuhlman B, Rohl CA, et al. Protein–Protein Docking with Simultaneous Optimization of Rigid-body Displacement and Side-chain Conformations. *Journal of Molecular Biology*. 2003 Aug; 331(1):281–99. Available from: <https://linkinghub.elsevier.com/retrieve/pii/S0022283603006703>. [https://doi.org/10.1016/s0022-2836\(03\)00670-3](https://doi.org/10.1016/s0022-2836(03)00670-3) PMID: 12875852
44. Lipton ZC, Tripathi S. Precise Recovery of Latent Vectors from Generative Adversarial Networks. arXiv:170204782 cs, stat. 2017 Feb. ArXiv: 1702.04782. Available from: <http://arxiv.org/abs/1702.04782>.
45. Egan N, Zhang J, Shen K. Generalized Latent Variable Recovery for Generative Adversarial Networks. arXiv:181003764 cs, stat. 2018 Oct. ArXiv: 1810.03764. Available from: <http://arxiv.org/abs/1810.03764>.
46. Lawrence MC, Colman PM. Shape Complementarity at Protein/Protein Interfaces. *Journal of Molecular Biology*. 1993; 234(4):946 950. Available from: <http://www.sciencedirect.com/science/article/pii/S0022283683716487>. <https://doi.org/10.1006/jmbi.1993.1648> PMID: 8263940
47. Alford RF, Leaver-Fay A, Jeliaskov JR, O'Meara MJ, DiMaio FP, Park H, et al. The Rosetta All-Atom Energy Function for Macromolecular Modeling and Design. *Journal of Chemical Theory and Computation*. 2017 Jun; 13(6):3031–48. Available from: <https://pubs.acs.org/doi/10.1021/acs.jctc.7b00125>. PMID: 28430426

48. Ma Y, Ding Y, Song X, Ma X, Li X, Zhang N, et al. Structure-guided discovery of a single-domain antibody agonist against human apelin receptor. *Sci Adv.* 2020 Jan; 6(3):eaax7379. <https://doi.org/10.1126/sciadv.aax7379> PMID: 31998837
49. Regep C, Georges G, Shi J, Popovic B, Deane CM. The H3 loop of antibodies shows unique structural characteristics. *Proteins: Structure, Function, and Bioinformatics.* 2017; 85(7):1311–8. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/prot.25291>.
50. Ruffolo JA, Sulam J, Gray JJ. Antibody structure prediction using interpretable deep learning. *Patterns.* 2022; 3(2):100406. Available from: <https://www.sciencedirect.com/science/article/pii/S2666389921002804>. <https://doi.org/10.1016/j.patter.2021.100406> PMID: 35199061
51. Anishchenko I, Pellock SJ, Chidyausiku TM, Ramelot TA, Ovchinnikov S, Hao J, et al. De novo protein design by deep network hallucination. *Nature.* 2021 Dec; 600(7889):547–52. Available from: <https://doi.org/10.1038/s41586-021-04184-w> PMID: 34853475
52. Jin W, Wohlwend J, Barzilay R, Jaakkola T. Iterative Refinement Graph Neural Network for Antibody Sequence-Structure Co-design. *arXiv*; 2021. Available from: <https://arxiv.org/abs/2110.04624>.
53. Conley GP, Viswanathan M, Hou Y, Rank DL, Lindberg AP, Cramer SM, et al. Evaluation of protein engineering and process optimization approaches to enhance antibody drug manufacturability. *Biotechnology and Bioengineering.* 2011 Nov; 108(11):2634–44. Available from: <http://doi.wiley.com/10.1002/bit.23220>. PMID: 21618474
54. Voynov V, Chennamsetty N, Kayser V, Helk B, Trout BL. Predictive tools for stabilization of therapeutic proteins. *mAbs.* 2009 Nov; 1(6):580–2. Available from: <http://www.tandfonline.com/doi/abs/10.4161/mabs.1.6.9773> PMID: 20068399
55. Jenkins N, Murphy L, Tyther R. Post-translational Modifications of Recombinant Proteins: Significance for Biopharmaceuticals. *Molecular Biotechnology.* 2008 Jun; 39(2):113–8. Available from: <http://link.springer.com/10.1007/s12033-008-9049-4>. <https://doi.org/10.1007/s12033-008-9049-4> PMID: 18327554
56. Dudgeon K, Rouet R, Kokmeijer I, Schofield P, Stolp J, Langley D, et al. General strategy for the generation of human antibody variable domains with increased aggregation resistance. *Proceedings of the National Academy of Sciences.* 2012 Jul; 109(27):10879–84. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.1202866109>. PMID: 22745168