
Joint Design of Protein Sequence and Structure based on Motifs

Zhenqiao Song¹, Yunlong Zhao^{2,3}, Yufei Song¹, Wenxian Shi⁴, Yang Yang^{2*}, Lei Li^{5*}

¹Department of Computer Science, University of California Santa Barbara

²Department of Chemistry and Biochemistry, University of California Santa Barbara

³Department of Chemistry, Massachusetts Institute of Technology

⁴Department of EECS, Massachusetts Institute of Technology

⁵Language Technology Institute, Carnegie Mellon University

{zhenqiao,yufei_song,yang89}@ucsb.edu, leili@cs.cmu.edu

{yunlongz,wxsh}@mit.edu

Abstract

Designing novel proteins with desired functions is crucial in biology and chemistry. However, most existing work focus on protein sequence design, leaving protein sequence and structure co-design underexplored. In this paper, we propose GeoPro, a method to design protein backbone structure and sequence jointly. Our motivation is that protein sequence and its backbone structure constrain each other, and thus joint design of both can not only avoid nonfolding and misfolding but also produce more diverse candidates with desired functions. To this end, GeoPro is powered by an equivariant encoder for three-dimensional (3D) backbone structure and a protein sequence decoder guided by 3D geometry. Experimental results on two biologically significant metalloprotein datasets, including β -lactamases and myoglobins, show that our proposed GeoPro outperforms several strong baselines on most metrics. Remarkably, our method discovers novel β -lactamases and myoglobins which are not present in protein data bank (PDB) and UniProt. These proteins exhibit stable folding and active site environments reminiscent of those of natural proteins, demonstrating their excellent potential to be biologically functional.

1 Introduction

A fundamental problem in protein engineering is designing novel proteins with desired biochemical functions such as catalytic activity [17], therapeutic efficacy [32], and fluorescence [9]. Proteins embody their function through spontaneous folding of amino acid sequences into three dimensional (3D) structures [18, 12, 48]. In particular, protein’s biochemical function is controlled by a subset of residues known as functional sites, or motifs [51]. Therefore, designing stably-folded proteins given a set of motifs is a promising direction to functional protein design.

In early representative work [26, 45], manually prepared rules are applied to discover motifs. Those rules are written by people with domain knowledge gained from the careful investigation of a specific protein family, which makes it hard to scale up to a wide number of protein families. Additionally, these efforts necessitate laborious trials and errors, making the overall process resource- and time-consuming.

To address these problems, machine learning methods have been employed to automate and improve the efficiency of protein design. Recent efforts in this field primarily focus on either protein sequence design guided by fitness landscapes [10, 19, 25, 40] or sequence generation from given 3D structures,

*Corresponding author.

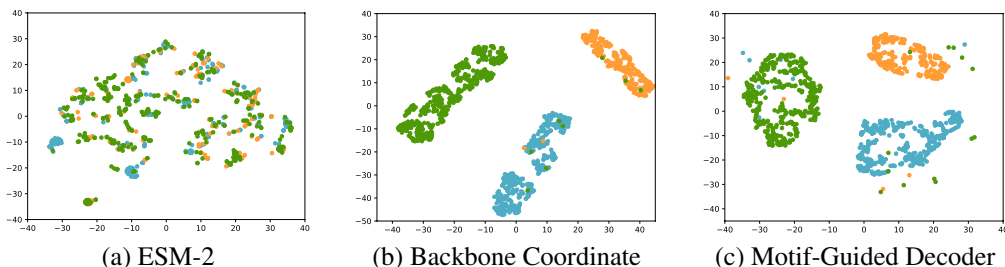


Figure 1: Visualization of three proteins of myoglobin, each containing many instances obtained by different crystal methods. Figure 1(a) is the sequence representation from ESM-2 [33]. Figure 1(b) shows the same proteins obtained by different crystal methods are closer to each other in 3D space. Figure 1(c) demonstrates protein sequence representations with closer 3D structure would also be clustered together after being revised by functional motifs.

also known as the inverse-folding problem [24, 23, 15, 27]. However, the former approaches do not consider any structural information, oftentimes leading to unfolded or misfolded proteins, while the latter one depend on an entire protein structure, leading to limited candidate diversity and novelty. A more natural way is to co-design protein sequence and structure based on critical motifs. However, current work either require manually selected motifs by biochemical experts [5, 51], which is difficult to generalize to arbitrary proteins, or to extend design based on secondary structure topology [1, 44], which can not guarantee the designed proteins exhibit desired functions.

In this paper, we propose GeoPro, an approach to design protein sequence and backbone structure jointly. Our motivation is that protein sequence and its backbone structure constrain each other, and thus joint design of both can not only avoid nonfolding and misfolding but also produce more diverse candidates with desired functions. GeoPro first uses a fine-tuned ESM-2 [33] to encode an initial protein sequence with only motifs into contextual representation. We design an equivariant graph neural network (EGNN) to encode and predict 3D backbone structure. This backbone structure encoder will refine protein residue representations through their neighboring interactions in 3D space. We further design a sequence decoder to generate the full protein sequence given the EGNN representations. GeoPro predicts the backbone structure while incorporating contextual residue representations, and inpaints the protein sequence based on functional motifs. Our approach is advantageous because the mutual constraints between backbone structure and sequence facilitates the design of stably-folded functional proteins. We provide a theoretical proof that structures within a bounded variance lead to sequences that belong to the same protein, as illustrated in Figure 1. This finding verifies that our discovered novel proteins are highly likely to be biologically functional.

We carry out extensive experiments on two metalloproteins, including β -lactamase and myoglobin, and compare the proposed model with several strong baselines. The contribution of this paper are listed as follows:

- We propose GeoPro to co-design protein sequence and backbone structure. It is able to design diverse, novel, and functional proteins that are not recorded in PDB and UniProt ².
- Experiments show that GeoPro achieves highest performance on most metrics. The designed β -lactamases and myoglobins exhibit stable folding and can respectively bind their metallofactors including zinc and heme, validating their excellent potential to be functional proteins.

2 Related Work

Generative Protein Design Protein sequence design has been studied with a wide variety of methods, including traditional directed evolution [6, 13, 39, 7] and machine learning methods [8, 4, 38, 49]. Following the success of deep generative models, there are some work focusing on protein sequence design with specific functions, aka. fitness. They either search satisfactory sequences using deep generative models [11, 10, 35, 31, 14, 22, 37, 5, 40], or directly generate protein sequences applying deep generative models [25, 46]. Another class of methods focus on inverse-folding problem [16, 24,

²UniProt is a huge protein sequence database.

52, 36, 23], which targets at producing a protein sequence that can fold into a given structure. Both approaches lack consideration for 3D structure design, resulting in constrained accuracy and novelty in the design outcomes.

3D Protein Design Wang et al. [51] propose inpainting method to recover both missing protein sequence and structure based on given motif segments. However, they need to pre-specify a possible inpainting length range, which is hard even for a biochemical expert. Trippe et al. [50] frame motif-scaffolding problem as a conditional sampling process in diffusion models. However, they only consider designing novel structures, which may not fully utilize the inherent correlation between protein sequence and structure. Anand and Achim [1] first propose to co-design protein sequence and structure conditioning on given secondary structures (SS). Following their work, Shi et al. [44] propose to realize design conditioning on SS and contact map. However, knowing the topology of a protein before design process is difficult and also cannot guarantee the designed proteins have the desired functions which are mostly determined by side-chain of residues at functional sites.

In this paper, we focus on joint design of protein sequence and structure conditioning on critical motifs, which leverages the relationship between protein sequence and backbone structure to help design not only stably-folded but also diverse and novel proteins with desired functions.

3 Background

3.1 Equivariance and Invariance

Equivariant Function A function f is said to be equivariant to the action of a group \mathcal{G} if $T_g(f(x)) = f(S_g(x))$ for all $g \in \mathcal{G}$, where S_g, T_g are linear representations related to the group element [43]. In this work, we consider the Euclidean group $E(3)$ generated by translations, rotations and reflections, for which S_g and T_g can be represented by a translation t and an orthogonal matrix R that rotates or reflects coordinates. f is then equivariant to a translation t , rotation or reflection R if transforming its input results in an equivalent transformation of its output, i.e. $f(Rx + t) = Rf(x) + t$.

Invariant Distribution In our setting, a conditional distribution $p(y|x)$ is invariant to the action of rotation or reflection R when:

$$p(y|x) = p(Ry|Rx) \text{ or } p(y|x) = p(y|Rx) \quad \text{for all orthogonal matrix } R \quad (1)$$

3.2 E(n) Equivariant Graph Neural Networks (EGNNs)

EGNNs [41] are a type of Graph Neural Network that satisfies the equivariance constraint. In this work, we consider interactions between all C_α in the backbone structure, and therefore assume a fully connected graph \mathcal{G} with nodes $v_i \in V$. Each node v_i is endowed with coordinates $x_i \in R^3$ as well as d -dimensional features $h_i \in R^d$. In this setting, EGNN consists of the composition of equivariant convolutional layers (EGCL): $x^{(l+1)}, h^{(l+1)} = \text{EGCL}[x^l, h^l]$, which are defined as:

$$\begin{aligned} m_{ij} &= \phi_e(h_i^l, h_j^l, d_{ij}^2, a_{ij}), \quad h_i^{l+1} = \phi_h(h_i^l, \sum_{j \neq i} \tilde{e}_{ij} m_{ij}) \\ x_i^{l+1} &= x_i^l + \sum_{j \neq i} \frac{x_j^l - x_i^l}{d_{ij} + 1} \phi_x(h_i^l, h_j^l, d_{ij}^2, a_{ij}) \end{aligned} \quad (2)$$

where l indexes the layer, and $d_{ij} = \|x_i^l - x_j^l\|_2$ is the Euclidean distance between nodes (v_i, v_j) , and a_{ij} are optional edge attributes. $\tilde{e}_{ij} = \phi_{inf}(m_{ij})$ is the attention score. All learnable components $(\phi_e, \phi_h, \phi_x, \phi_{inf})$ are parametrized by fully connected neural networks. An entire EGNN architecture is then composed of L EGCL layers, which applies the following non-linear transformation: $\hat{x}, \hat{h} = \text{EGNN}[x^0, h^0]$. This transformation satisfies the required equivariant property:

$$R\hat{x} + t, \hat{h} = \text{EGNN}[Rx^0 + t, h^0] \quad \text{for all orthogonal matrix } R \text{ and } t \in R^3 \quad (3)$$

4 Methods

In this section, we describe our method in detail. We first give the problem formulation in 4.1. Then we detail the interactive process in which a backbone structure encoder predicts the backbone

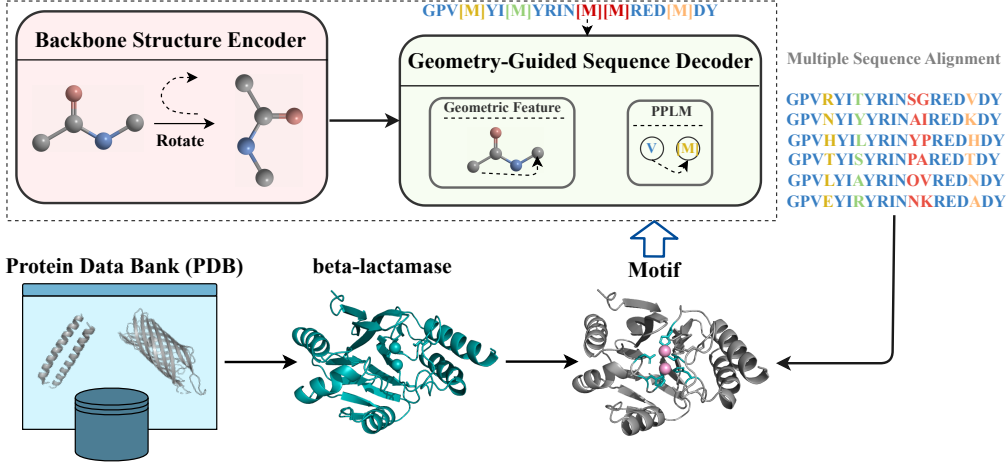


Figure 2: The overall architecture of the proposed GeoPro.

structure by incorporating the contextual residue representations in 4.2, and a geometry-guided decoder inpaints the protein sequence based on the functional geometry-guided structure information in 4.3. Finally, we provide a theoretical proof in 4.4 that protein structures within a bounded variance lead to sequences which belong to the same protein. The overall architecture is illustrated in Figure 2.

4.1 Problem Formulation

Our goal is to design a protein with desired function by co-designing its 3D backbone structure $x = \{x_1, x_2, \dots, x_L\}$ and amino acid sequence $y = \{y_1, y_2, \dots, y_L\}$ conditioning on a given motif $z = \{z_b, z_s\}$. L is the protein sequence length and $x_i \in \mathbb{R}^3$ is the 3D coordinate of the C_α of i -th amino acid $y_i \in \mathcal{V} - \mathcal{V}$ consists of 20 amino acids. A given motif z consists of backbone structure segments z_b and sequence segments z_s . To achieve this goal, we aim to maximize the joint probability of x and y conditioning on z :

$$p(x, y|z) = p(x|z; \theta)p(y|z; \phi) \quad (4)$$

where $p(x|z; \theta)$ is modeled by the backbone structure encoder with parameter θ , and $p(y|z; \phi)$ is modeled by the geometry-guided sequence decoder with parameter ϕ . We will prove the joint probability $p(x, y|z)$ is invariant to the translation and rotation actions on z in the following subsections.

4.2 Backbone Structure Encoder

The backbone structure encoder utilizes contextual residue representations to predict the protein backbone structure, while also refining the node features through their interactions with neighboring residues in 3D space.

Specifically, we leverage EGNN (Equation 2 and 3) as the backbone structure encoder:

$$\hat{x}, \hat{h} = \text{EGNN}[x^0, h^0; \theta] \quad (5)$$

where $x^0 = z_b \cup \hat{x}_{-z_b}^0$ is the initialized 3D backbone coordinates, in which the motif z_b keeps fixed while other nodes $\hat{x}_{-z_b}^0$ are sequentially initialized on a spherical surface to their neighboring amino acid. In particular, we find the Euclidean distance between every neighboring C_α is almost the same (around $r = 3.75\text{\AA}$). Therefore, for a node $z_{k+1} \in \hat{x}_{-z_b}^0$, we initialize it as a spherical surface with the center equals to its nearest neighbor $z_k \in z_b$:

$$\hat{x}_{k+1}^0 = [z_k[0] + r \cdot \sin \omega_1 \cos \omega_2, z_k[1] + r \cdot \sin \omega_1 \sin \omega_2, z_k[2] + r \cdot \cos \omega_1] \quad (6)$$

where $\omega_1 \sim \text{Uniform}(0, \pi)$ and $\omega_2 \sim \text{Uniform}(0, 2\pi)$. For h^0 , we directly apply a pretrained protein language model (PPLM) [33] to encode the corrupted protein sequence $\tilde{y} = z_s \cup \tilde{y}_{-z_s}$ with \tilde{y}_{-z_s} equals to [mask] symbol: $h^0 = \text{PPLM}(\tilde{y})$. We initialize node features in this way because we believe a PPLM should provide some contextual information based on the large-scale protein sequences it has seen in the pretraining stage.

After getting the revised backbone structure, we minimize the Euclidean distance to learn the model:

$$\mathcal{L}_b = \sum_{x_j \in x, x_j \notin z_b} \|x_j - \hat{x}_j\|_2^2 \quad (7)$$

which equals to maximize a three-dimensional Gaussian distribution whose mean vector equals to x_j and covariance matrix equals to identity matrix for each node x_j :

$$\begin{aligned} \log p(x|z; \theta) &= \log \prod_{x_j \in x, x_j \notin z_b} p(x_j|z; \theta) = \sum_{x_j \in x, x_j \notin z_b} \log p(x_j|z; \theta) \\ \log p(x_j|z; \theta) &= \log \left\{ \frac{1}{\sqrt{(2\pi)^3}} \exp \left(-\frac{1}{2} (x_j - \hat{x}_j)^T (x_j - \hat{x}_j) \right) \right\} = -\|x_j - \hat{x}_j\|_2^2 + \text{const} \end{aligned} \quad (8)$$

As shown in Equation 3, EGNN satisfies the equivariance constraint, so we have:

$$\begin{aligned} p(Rx_j + t|Rz + t; \theta) &= \frac{1}{\sqrt{(2\pi)^3}} \exp \left(-\frac{1}{2} (Rx_j + t - R\hat{x}_j - t)^T (Rx_j + t - R\hat{x}_j - t) \right) \\ &= \frac{1}{\sqrt{(2\pi)^3}} \exp \left(-\frac{1}{2} (x_j - \hat{x}_j)^T R^T R (x_j - \hat{x}_j) \right) = p(x_j|z; \theta) \end{aligned} \quad (9)$$

where R is an orthogonal matrix. Therefore, we can see the distribution of predicting backbone structure is invariant to the rotation or translation actions on motif.

The backbone structure encoder enables us to not only predict the backbone structure but also enhance the residue representations by considering their interactions with neighboring residues in 3D space. This refinement process proves advantageous for the subsequent sequence inpainting procedure.

4.3 Geometry-Guided Sequence Decoder

The geometry-guided sequence decoder (GSD) reconstructs the protein sequence by utilizing functional-geometry structure information, thereby facilitating the discovery of novel proteins that exhibit similar biological functions to natural ones.

In particular, we leverage the residue features \hat{h} revised by motif geometry to reconstruct the original protein sequence as follows:

$$p(y|z; \phi) = \prod_{y_j \in y, y_j \notin z_s} p(y_j|z; \phi) = \text{GSD}(f(\hat{h}); \phi), \quad f(\hat{h}_j) = \begin{cases} \hat{h}_j, & y_j \in z_s, \\ \text{Emb}[\text{mask}], & \text{otherwise} \end{cases} \quad (10)$$

where \hat{h}_j is the output of backbone structure encoder of the j -th token in Equation 5. Here we use a finetuned ESM-2 [33] to initialize GSD. Then to optimize GSD, we minimize the negative log likelihood of the none-motif parts:

$$\mathcal{L}_s = \sum_{y_j \in y, y_j \notin z_s} -\log p(y_j|z; \phi) \quad (11)$$

As shown in Equation 3, \hat{h} is invariant to the rotation or translation actions on x^0 , and thus we can verify that $p(y|z; \phi)$ is invariant:

$$p(y|Rz + t; \phi) = \text{GSD}(f(\hat{h}); \phi) = p(y|z; \phi) \quad (12)$$

Combining the conclusions in Equation 9 and 12, we can prove the joint probability $p(x, y|z)$ is invariant to the translation and rotation actions on z :

$$p(Rx + t, y|Rz + t) = p(Rx + t|Rz + t; \theta) p(y|Rz + t; \phi) = p(x|z; \theta) p(y|z; \phi) = p(x, y|z) \quad (13)$$

Accordingly, our overall training objective is defined as:

$$\mathcal{L} = \alpha * \mathcal{L}_b + \beta * \mathcal{L}_s \quad (14)$$

where α and β respectively control the importance of these two terms. Jointly minimizing objective 14 equals to maximize the joint probability $p(x, y|z)$.

4.4 Theoretical Analysis

For each protein in PDB, there might be several instances due to different crystal methods, which have the same or different sequences as well as slightly different coordinates. Sequence representations incorporating this kind of geometric information would get much closer to each other than other different proteins as shown in Figure 1(c). We can prove that our functional geometry-guided sequence decoding loss has an upper bound. It means proteins having the similar structures are mapped to the geometry-guided representations in a manner that are well-separated from others. This theoretically verifies that our discovered novel proteins are highly possible to be realistic ones due to their similar active site environments to natural proteins.

In particular, we can regard our sequence reconstruction process as an auto-encoding process [30] with some perturbation \mathcal{C} : $\tilde{y} = \mathcal{C}(y)$. Here the encoding process is $g = E(y) = \text{EGNN}(\text{PPLM}(\mathcal{C}(y)))$, and decoding performs $y = \text{GSD}(g)$. Let σ denote the sigmoid function.

Assumption 4.1. The decoder D can approximate arbitrary $p(y|g)$ so long as it remains sufficiently Lipschitz continuous in g . Namely, there exists $L > 0$ such that decoder D obtainable via training satisfies: $\forall y \in \mathcal{Y}, \forall g^1, g^2 \in \mathcal{Z}, |\log p(y|g^1) - \log p(y|g^2)| \leq L\|g^1 - g^2\|$. (We denote this set of decoders \mathcal{D}_L .)

Theorem 4.2. Suppose $\{y^1, y^2, \dots, y^n\}$ are belonging to n/K different proteins of equal size K instances, with s_i denoting the specific protein of y^i . Suppose $p(y^i|\tilde{y}^j) = 1/K$ if $s_i = s_j$ and 0 otherwise. With a deterministic encoder mapping E from $\{y^1, y^2, \dots, y^n\}$ to $\{g^1, g^2, \dots, g^n\}$, the denoising objective $\max_{D \in \mathcal{D}_L} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n p(y^j|\tilde{y}^i) \log p(y^i|E(y^j))$ has an upper bound: $\frac{1}{n^2} \sum_{i,j:s_i \neq s_j} \log \sigma(L\|E(y^i) - E(y^j)\|) - \log K$.

Note that y^i here denotes the i -th protein sequence in the dataset. \mathcal{Y} denotes the set of all protein sequences. We provide the proof in appendix A.

5 Experiments

5.1 Datasets

We evaluate our method on two metalloproteins: β -lactamase binding zinc ion, and myoglobin binding heme. We discuss the significance of both proteins in Appendix B.1. To obtain the data, we first collect these two kinds of proteins from PDB and then extract chain A from these proteins. Only proteins capable of binding the corresponding metallofactors are retained. Next we perform length filtering for both proteins, i.e., reserving β -lactamases longer than 200 and myoglobins longer than 100. Then we run MSAs using ClustalW2 method [2]. Based on the results, we set positions whose alignment frequencies are higher than a given threshold λ as the motif parts and others as flexible parts. Then we randomly split each dataset into training/validation/test sets with the ratio 8 : 1 : 1. The detailed data statistics are shown in Appendix Table 3.

5.2 Experimental Details

Training Details We use two-layer EGNN [41] with hidden size equal to 320 as the backbone structure encoder and a finetuned ESM-2 [33] to initialize the GSD. The sequences are decoded using sampling strategy with top-3. More implementation details are provided in Appendix B.3.

Baseline Models We compare the proposed GeoPro against the following representative baselines: (1) **Hallucination** [5] uses MCMC [3] incorporating a motif constraint into the acceptance score calculation. (2) **Inpainting** [51] recovers both sequence and structure based on the given protein segments. (3) **SMCDiff** [50]+**ProteinMPNN** [15]: We first apply SMCDiff to design a protein structure based on given motifs and use ProteinMPNN to generate a sequence based on the given structure.

To better analyze the influence of different components in our model, we also conduct ablation tests: (7) **GeoPro-w/o-ctx**: The node feature of the backbone structure encoder is randomly initialized without any contextual residue information. (9) **GeoPro-w/o-geo** directly applies the residue embeddings from the finetuned ESM-2 instead of the functional geometry revised representations. (9) **GeoPro-ESM** directly applies the pretrained ESM-2 without any further training.

	Models	AAR (%, \uparrow)	RMSD (\AA , \downarrow)	pLDDT (\uparrow)	TM-score (\uparrow)
β -lactamase	Hallucination	4.79	--	30.5511	0.2918
	Inpainting	16.73	4.0599	61.7679	0.3790
	SMCDiff+PrteinMPNN	19.94	10.3960	42.0375	0.3458
	GeoPro	43.41	2.9825	62.7349	0.4256
myoglobin	Hallucination	4.81	--	38.2817	0.2754
	Inpainting	39.59	3.3751	67.0813	0.4391
	SMCDiff+ProteinMPNN	12.47	8.0067	34.5914	0.2235
	GeoPro	51.12	2.9891	77.3399	0.4656

Table 1: Model performance on two metalloprotein datasets. GeoPro achieves the best performance on both datasets.

	Models	AAR (%, \uparrow)	RMSD (\AA , \downarrow)	pLDDT (\uparrow)	TM-score (\uparrow)
β -lactamase	GeoPro	43.41	2.9825	62.7349	0.4256
	- w/o-ctx	45.18	3.8759	57.2319	0.4125
	- w/o-geo	42.07	3.5755	57.8732	0.4109
	- ESM	39.98	3.7218	55.9130	0.4019
myoglobin	GeoPro	51.12	2.9891	77.3399	0.4656
	- w/o-ctx	46.19	3.7615	68.7901	0.4502
	- w/o-geo	49.56	3.3132	63.4686	0.4419
	- ESM	42.13	4.0289	61.6971	0.4209

Table 2: Ablation study results: Either removing sequence contextual or geometric guidance would lead to performance degradation, highlighting how co-design can enhance superior protein design.

Evaluation Metrics We use the following automatic metrics to evaluate the quality of the designed proteins: (1) **AAR** assesses how similar the designed sequence is to the target sequence. (2) **RMSD** evaluates how close our designed structure is to the target structure. (3) **pLDDT** [28] provides an overall confidence score that a designed protein sequence can fold into a structure which is similar to natural proteins. We apply ESMFold [34] to calculate pLDDT due to its much higher efficiency than AlphaFold2 [28]. (4) **TM-score** [53] evaluates how similar structure predicted by the designed sequence is to the target structure. We use ESMFold to predict the structure of the designed sequence.

5.3 Main Results

Table 1 reports the performance of all models.

GeoPro can design more realistic proteins. Table 1 shows our proposed GeoPro achieves highest pLDDT and TM-score among all competitors on both datasets, demonstrating our model can generate more realistic proteins with relatively high confidence. Our interpretation is that our GeoPro leverages the inherent correlation between protein sequence and backbone structure, which can not only constrain each other to avoid misfolding but also benefit from each other to design a biologically functional protein with higher potential.

GeoPro has the smallest recovery error. As shown in Table 1, our model outperforms all competitors on AAR and RMSD, exhibiting its best performance on both sequence and structure reconstruction. It is because our GeoPro takes advantage of the relationship between protein sequence and backbone structure. Co-designing these two parts leads to more accurate sequence and structure through their interaction with each other during the design process.

5.4 Ablation Study

Table 2 shows the results of ablation study. Without finetuning the sequence decoder, the performance drops the most on all metrics (GeoPro-ESM). It validates that related information can be better explored from a pretrained protein language model by tuning the parameter to the specific protein family. Without the contextual sequence representation initialization, GeoPro-w/o-ctx achieves worse RMSD scores than GeoPro-w/o-geo. Instead, after removing the geometric information guidance,

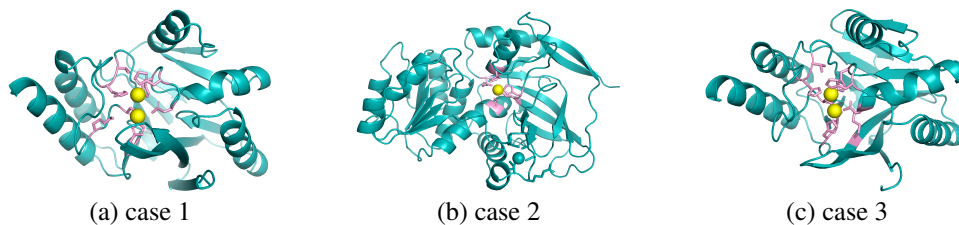


Figure 3: Designed proteins of β -lactamase which belong to different subclasses (a) B1, (b) B2, (c) B3 metal-dependent β -lactamases.

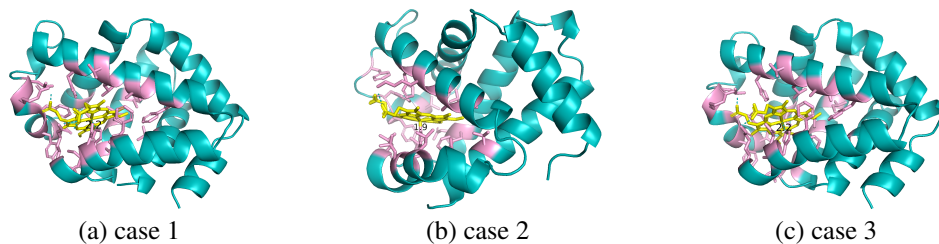


Figure 4: Designed proteins of myoglobin which bind heme through hydrogen bonds.

GeoPro-w/o-geo obtains lower amino acid recovery rate than GeoPro-w/o-ctx. Once again, these phenomena validate the mutual benefits between protein sequence and structure, highlighting how co-design can enhance the generation of superior protein candidates.

6 Analysis

To intuitively know how well our model can perform, we visualize the designed proteins. Specifically, we first randomly select 3 cases from the top-100 candidates according to pLDDT and then use AlphaFold2 to predict the protein structure, which will be subsequently used to predict the relevant ligand by AlphaFill [21]. Finally we apply PyMOL [42] to visualize the final results. From Figure 3 and 4, we can see that all the designed β -lactamases and myoglobins have active site environments highly similar to natural proteins, and can respectively bind zinc ion and heme (yellow parts), demonstrating our model can design functional proteins. Furthermore, all these cases are not included by PDB, and some of them even have lower identity overlapping with UniProt sequences (e.g., 61.0% AAR of Figure 3 (b)), verifying our GeoPro is able to design novel proteins. Additionally, the three β -lactamases belong to three different subclasses of metal-dependent β -lactamases featuring complementary Zn coordination chemistries, validating that our model can design diverse proteins. We provide more cases in Appendix C and all the cases are supplied in supplementary material.

7 Conclusion

This paper proposes GeoPro, a method to co-design protein sequence and backbone structure. The proposed model leverages the inherent correlation between protein sequence and backbone structure, and is powered by an equivariant 3D backbone structure encoder and a geometry-guided sequence decoder. Experimental results show that our proposed GeoPro outperforms several strong baselines on most metrics. Additionally, our model discovers novel β -lactamases and myoglobins which are not recorded by PDB and UniProt. One limitation of this work is, although our model has demonstrated promising results, the designed proteins have not undergone wet-lab testing. Consequently, we cannot provide complete assurance regarding the ability of the designed metalloproteins to bind the corresponding metal cofactor and perform their desired functions. Future work will involve the wet-lab testing to verify the metal binding ability and biological function of our designed metalloproteins.

References

- [1] Namrata Anand and Tudor Achim. Protein structure and sequence generation with equivariant denoising diffusion probabilistic models. *arXiv preprint arXiv:2205.15019*, 2022.
- [2] Catherine L Anderson, Cory L Strobe, and Etsuko N Moriyama. Suitensa: visual tools for multiple sequence alignment comparison and molecular sequence simulation. *BMC bioinformatics*, 12(1):1–14, 2011.
- [3] Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. *Machine learning*, 50:5–43, 2003.
- [4] Christof Angermueller, David Dohan, David Belanger, Ramya Deshpande, Kevin Murphy, and Lucy Colwell. Model-based reinforcement learning for biological sequence design. In *International conference on learning representations*, 2019.
- [5] Ivan Anishchenko, Samuel J Pellock, Tamuka M Chidyausiku, Theresa A Ramelot, Sergey Ovchinnikov, Jingzhou Hao, Khushboo Bafna, Christoffer Norn, Alex Kang, Asim K Bera, et al. De novo protein design by deep network hallucination. *Nature*, 600(7889):547–552, 2021.
- [6] Frances H Arnold. Design by directed evolution. *Accounts of chemical research*, 31(3):125–131, 1998.
- [7] Frances H Arnold. Directed evolution: bringing new chemistry to life. *Angewandte Chemie International Edition*, 57(16):4143–4148, 2018.
- [8] David Belanger, Suhani Vora, Zeldia Mariet, Ramya Deshpande, David Dohan, Christof Angermueller, Kevin Murphy, Olivier Chapelle, and Lucy Colwell. Biological sequences design using batched bayesian optimization. 2019.
- [9] Surojit Biswas, Grigory Khimulya, Ethan C Alley, Kevin M Esvelt, and George M Church. Low-n protein engineering with data-efficient deep learning. *Nature methods*, 18(4):389–396, 2021.
- [10] David Brookes, Hahnbeom Park, and Jennifer Listgarten. Conditioning by adaptive sampling for robust design. In *International conference on machine learning*, pages 773–782. PMLR, 2019.
- [11] David H Brookes and Jennifer Listgarten. Design by adaptive sampling. *arXiv preprint arXiv:1810.03714*, 2018.
- [12] Cyrus Chothia. Principles that determine the structure of proteins. *Annual review of biochemistry*, 53(1):537–572, 1984.
- [13] Paul A Dalby. Strategy and success for the directed evolution of enzymes. *Current opinion in structural biology*, 21(4):473–480, 2011.
- [14] Payel Das, Tom Sercu, Kahini Wadhawan, Inkit Padhi, Sebastian Gehrmann, Flaviu Cipcigan, Vijil Chenthamarakshan, Hendrik Strobelt, Cicero Dos Santos, Pin-Yu Chen, et al. Accelerated antimicrobial discovery via deep generative models and molecular dynamics simulations. *Nature Biomedical Engineering*, 5(6):613–623, 2021.
- [15] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- [16] Sarel J Fleishman, Andrew Leaver-Fay, Jacob E Corn, Eva-Maria Strauch, Sagar D Khare, Nobuyasu Koga, Justin Ashworth, Paul Murphy, Florian Richter, Gordon Lemmon, et al. Rosettascripts: a scripting language interface to the rosetta macromolecular modeling suite. *PLoS one*, 6(6):e20161, 2011.
- [17] Richard J Fox, S Christopher Davis, Emily C Mundorff, Lisa M Newman, Vesna Gavrilovic, Steven K Ma, Loleta M Chung, Charlene Ching, Sarena Tam, Sheela Muley, et al. Improving catalytic function by prosar-driven enzyme evolution. *Nature biotechnology*, 25(3):338–344, 2007.

- [18] Nobuhiro Go. Theoretical studies of protein folding. *Annual review of biophysics and bioengineering*, 12(1):183–210, 1983.
- [19] Anvita Gupta and James Zou. Feedback gain for dna optimizes protein functions. *Nature Machine Intelligence*, 1(2):105–111, 2019.
- [20] Varsha Gupta. Metallo beta lactamases in pseudomonas aeruginosa and acinetobacter species. *Expert opinion on investigational drugs*, 17(2):131–143, 2008.
- [21] Maarten L Hekkelman, Ida de Vries, Robbie P Joosten, and Anastassis Perrakis. Alphafill: enriching alphafold models with ligands and cofactors. *Nature Methods*, 20(2):205–213, 2023.
- [22] Samuel C Hoffman, Vijil Chenthamarakshan, Kahini Wadhawan, Pin-Yu Chen, and Payel Das. Optimizing molecules using efficient queries from property evaluations. *Nature Machine Intelligence*, 4(1):21–31, 2022.
- [23] Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. In *International Conference on Machine Learning*, pages 8946–8970. PMLR, 2022.
- [24] John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph-based protein design. *Advances in neural information processing systems*, 32, 2019.
- [25] Moksh Jain, Emmanuel Bengio, Alex Hernandez-Garcia, Jarrid Rector-Brooks, Bonaventure FP Dossou, Chanakya Ajit Ekbote, Jie Fu, Tianyu Zhang, Michael Kilgour, Dinghuai Zhang, et al. Biological sequence design with gflownets. In *International Conference on Machine Learning*, pages 9786–9801. PMLR, 2022.
- [26] Lin Jiang, Eric A Althoff, Fernando R Clemente, Lindsey Doyle, Daniela Rothlisberger, Alexandre Zanghellini, Jasmine L Gallaher, Jamie L Betker, Fujie Tanaka, Carlos F Barbas III, et al. De novo computational design of retro-aldol enzymes. *science*, 319(5868):1387–1391, 2008.
- [27] Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael JL Townshend, and Ron Dror. Learning from protein structure with geometric vector perceptrons. *arXiv preprint arXiv:2009.01411*, 2020.
- [28] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [29] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2014.
- [30] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [31] Aviral Kumar and Sergey Levine. Model inversion networks for model-based optimization. *Advances in Neural Information Processing Systems*, 33:5126–5137, 2020.
- [32] HA Daniel Lagassé, Aikaterini Alexaki, Vijaya L Simhadri, Nobuko H Katagiri, Wojciech Jankowski, Zuben E Sauna, and Chava Kimchi-Sarfaty. Recent advances in (therapeutic protein) drug development. *F1000Research*, 6, 2017.
- [33] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- [34] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- [35] Ali Madani, Bryan McCann, Nikhil Naik, Nitish Shirish Keskar, Namrata Anand, Raphael R Eguchi, Po-Ssu Huang, and Richard Socher. Progen: Language modeling for protein generation. *arXiv preprint arXiv:2004.03497*, 2020.

- [36] Matt McPartlon, Ben Lai, and Jinbo Xu. A deep se (3)-equivariant model for learning inverse protein folding. *bioRxiv*, pages 2022–04, 2022.
- [37] Igor Melnyk, Payel Das, Vijil Chenthamarakshan, and Aurelie Lozano. Benchmarking deep generative models for diverse antibody sequence design. *arXiv preprint arXiv:2111.06801*, 2021.
- [38] Henry Moss, David Leslie, Daniel Beck, Javier Gonzalez, and Paul Rayson. Boss: Bayesian optimization over string spaces. *Advances in neural information processing systems*, 33: 15476–15486, 2020.
- [39] Michael S Packer and David R Liu. Methods for the directed evolution of proteins. *Nature Reviews Genetics*, 16(7):379–394, 2015.
- [40] Zhizhou Ren, Jiahan Li, Fan Ding, Yuan Zhou, Jianzhu Ma, and Jian Peng. Proximal exploration for model-guided protein sequence design. *bioRxiv*, 2022.
- [41] Victor Garcia Satorras, Emiel Hooeboom, and Max Welling. E (n) equivariant graph neural networks. In *International conference on machine learning*, pages 9323–9332. PMLR, 2021.
- [42] LLC Schrödinger and Warren DeLano. Pymol. URL <http://www.pymol.org/pymol>.
- [43] Jean-Pierre Serre et al. *Linear representations of finite groups*, volume 42. Springer, 1977.
- [44] Chence Shi, Chuanrui Wang, Jiarui Lu, Bozitao Zhong, and Jian Tang. Protein sequence and structure co-design with equivariant translation. *arXiv preprint arXiv:2210.08761*, 2022.
- [45] Justin B Siegel, Alexandre Zanghellini, Helena M Lovick, Gert Kiss, Abigail R Lambert, Jennifer L St. Clair, Jasmine L Gallaher, Donald Hilvert, Michael H Gelb, Barry L Stoddard, et al. Computational design of an enzyme catalyst for a stereoselective bimolecular diels-alder reaction. *Science*, 329(5989):309–313, 2010.
- [46] Zhenqiao Song and Lei Li. Importance weighted expectation-maximization for protein sequence design. *arXiv preprint arXiv:2305.00386*, 2023.
- [47] Barry A Springer, Stephen G Sligar, John S Olson, and George N Jr Phillips. Mechanisms of ligand recognition in myoglobin. *Chemical Reviews*, 94(3):699–714, 1994.
- [48] Tyler N Starr and Joseph W Thornton. Exploring protein sequence–function landscapes. *Nature biotechnology*, 35(2):125–126, 2017.
- [49] Kei Terayama, Masato Sumita, Ryo Tamura, and Koji Tsuda. Black-box optimization for automated discovery. *Accounts of Chemical Research*, 54(6):1334–1346, 2021.
- [50] Brian L Trippe, Jason Yim, Doug Tischer, Tamara Broderick, David Baker, Regina Barzilay, and Tommi Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. *arXiv preprint arXiv:2206.04119*, 2022.
- [51] Jue Wang, Sidney Lisanza, David Juergens, Doug Tischer, Joseph L Watson, Karla M Castro, Robert Ragotte, Amijai Saragovi, Lukas F Milles, Minkyung Baek, et al. Scaffolding protein functional sites using deep learning. *Science*, 377(6604):387–394, 2022.
- [52] Peng Xiong, Xiuhong Hu, Bin Huang, Jiahai Zhang, Quan Chen, and Haiyan Liu. Increasing the efficiency and accuracy of the abacus protein sequence design method. *Bioinformatics*, 36(1):136–144, 2020.
- [53] Yang Zhang and Jeffrey Skolnick. Tm-align: a protein structure alignment algorithm based on the tm-score. *Nucleic acids research*, 33(7):2302–2309, 2005.

Appendix

A Proof of Theorem 3.2

Theorem A.1. Suppose $\{y^1, y^2, \dots, y^n\}$ are belonging to n/K different proteins of equal size K instances, with s_i denoting the specific protein of y^i . Suppose $p(y^i|\tilde{y}^j) = 1/K$ if $s_i = s_j$ and 0 otherwise. With a deterministic encoder mapping E from $\{y^1, y^2, \dots, y^n\}$ to $\{g^1, g^2, \dots, g^n\}$, the denoising objective $\max_{D \in \mathcal{D}_L} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n p(y^j|\tilde{y}^i) \log p(y^i|E(y^j))$ has an upper bound: $\frac{1}{n^2} \sum_{i,j:s_i \neq s_j} \log \sigma(L\|E(y^i) - E(y^j)\|) - \log K$.

Proof. For convinience, we denote $g^j = E(y^j)$. We consider what is the optimal decoder probability assignment $p(y^i|g^j)$ under the Lipschitz constraint. Suppose g^i, g^j satisfy that with some $0 < \delta < \zeta$: $\|g^i - g^j\| < \delta$ if $s_i = s_j$ and $\|g^i - g^j\| > \zeta$ otherwise. To lower bound the training objective, we can choose:

$$p(y^i|g^j) = \begin{cases} \frac{1-\gamma}{K}, & s_i = s_j \\ \frac{\gamma}{n-K}, & \text{otherwise} \end{cases} \quad (15)$$

with $\gamma = \sigma(-L\zeta) \in (0, \frac{1}{2})$, where σ denotes sigmoid function. Note that this choice can ensure $\sum_{i \in [n]} p(y_i|g^j) = 1$ for each $j \in [n]$, and this also does not violate Lipschitz condition.

The objective is to find optimal decoder D which:

$$\begin{aligned} & \max_{D \in \mathcal{D}_L} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n p(y^j|\tilde{y}^i) \log p(y^i|g^j) \\ &= \max_{D \in \mathcal{D}_L} \frac{1}{nK} \sum_j \sum_{i:s_i=s_j} \log p(y_i|g^j) \end{aligned} \quad (16)$$

Now let us define $P_j = \sum_{i:s_i=s_j} p(y^i|g^j) = K \cdot p(y^j|g^j)$, the above optimizing objective becomes:

$$\begin{aligned} & \max_{D \in \mathcal{D}_L} \frac{1}{nK} \sum_j \sum_{i:s_i=s_j} \log p(y_i|g^j) \\ &= \max_{D \in \mathcal{D}_L} \frac{1}{n} \sum_j \log p(y_j|g^j) \\ &= \max_{D \in \mathcal{D}_L} \frac{1}{n} \sum_j \log P_j - \log K \\ &= \max_{D \in \mathcal{D}_L} \frac{1}{2n^2} \sum_i \sum_j (\log P_i + \log P_j) - \log K \\ &\leq \frac{1}{2n^2} \sum_i \sum_j \max_{D \in \mathcal{D}_L} (\log P_i + \log P_j) - \log K \end{aligned} \quad (17)$$

For each term $\max_{D \in \mathcal{D}_L} \log P_i + \log P_j$:

$$\begin{aligned} \log P_i + \log P_j &= 2 \log K \cdot \frac{1-\gamma}{K} \\ &= 2 \log(1-\gamma) \\ &= 2 \log \frac{\exp(L\zeta)}{1 + \exp(L\zeta)} \\ &= 2 \log \sigma(-L\zeta) \\ &\leq 2 \log \sigma(-L\|g^i - g^j\|) \end{aligned} \quad (18)$$

Overall, we have:

$$\begin{aligned} & \max_{D \in \mathcal{D}_L} \frac{1}{nK} \sum_j \sum_{i:s_i=s_j} \log p(y_i|g^j) \\ &\leq \frac{1}{n^2} \sum_i \sum_{j:s_i \neq s_j} \log \sigma(-L\|g^i - g^j\|) - \log K \end{aligned} \quad (19)$$

Protein	PDB	Metal Binding	Length Filtering
β -lactamase	171,484	7,802	5,427
myoglobin	14,573	3,381	3,381

Table 3: Detailed data statistics of the two metalloproteins.

B Additional Experimental Details

B.1 Significance of Metalloproteins Studied Herein

Metalloproteins comprise almost 50% of all the naturally occurring proteins. The Zn-dependent β -lactamases and Fe-dependent myoglobins studied herein represent biologically significant metalloprotein examples. In particular, β -lactamases are enzymes produced by microorganisms to break down β -lactam antibiotics, conferring antibiotic resistance [20]. Thus, the study and design of β -lactamases hold relevance to public health and play a critical role in the development of new antibiotics. On the other hand, myoglobin is a heme-containing protein involved in oxygen storage and transport in muscle tissue [47], highlighting their biological significance.

Overall, these two proteins both have significant research values. Our experimental results show our proposed GeoPro has the capability to design novel proteins with desired functions for both general protein and enzyme, exhibiting its superior generalization.

B.2 Protein Data Statistics

Detailed data statistics for β -lactamase and myoglobin are reported in Table 3. We will release the two created metalloprotein datasets in the near future.

B.3 More Implementation Details

The mini-batch size and learning rate are set to 4 sequences and $1e-7$ respectively. The model is trained for 10 epochs with 1 NVIDIA RTX A6000 GPU card. We apply Adam [29] as the optimizer with a linear warm-up over the first 4,000 steps and linear decay for later steps. We tune the hyperparameters α and β both from 0.1 to 1.0 with step size 0.1 and find $\alpha = 0.1$ and $\beta = 1.0$ performs best on the validation set for β -lactamase. For myoglobin, we find the gradient tends to vanish when α is relatively large, and thus we tune it from 0.01 to 0.1 with step size 0.01 and find $\alpha = 0.01$ performs best on the corresponding validation set.

C More Designed Cases

Figure 5 and 6 respectively illustrate more designed proteins for β -lactamase and myoglobin. It shows all proteins exhibit active site environments reminiscent of those of natural proteins and can bind corresponding metallofactors, i.e., β -lactamases bind zinc ion and myoglobins bind heme, demonstrating their excellent potential to be biologically functional. Besides, the sequences of these proteins are not present in PDB, and exhibit diverse structures, demonstrating our GeoPro has the ability to design novel and diverse proteins with desired functions.

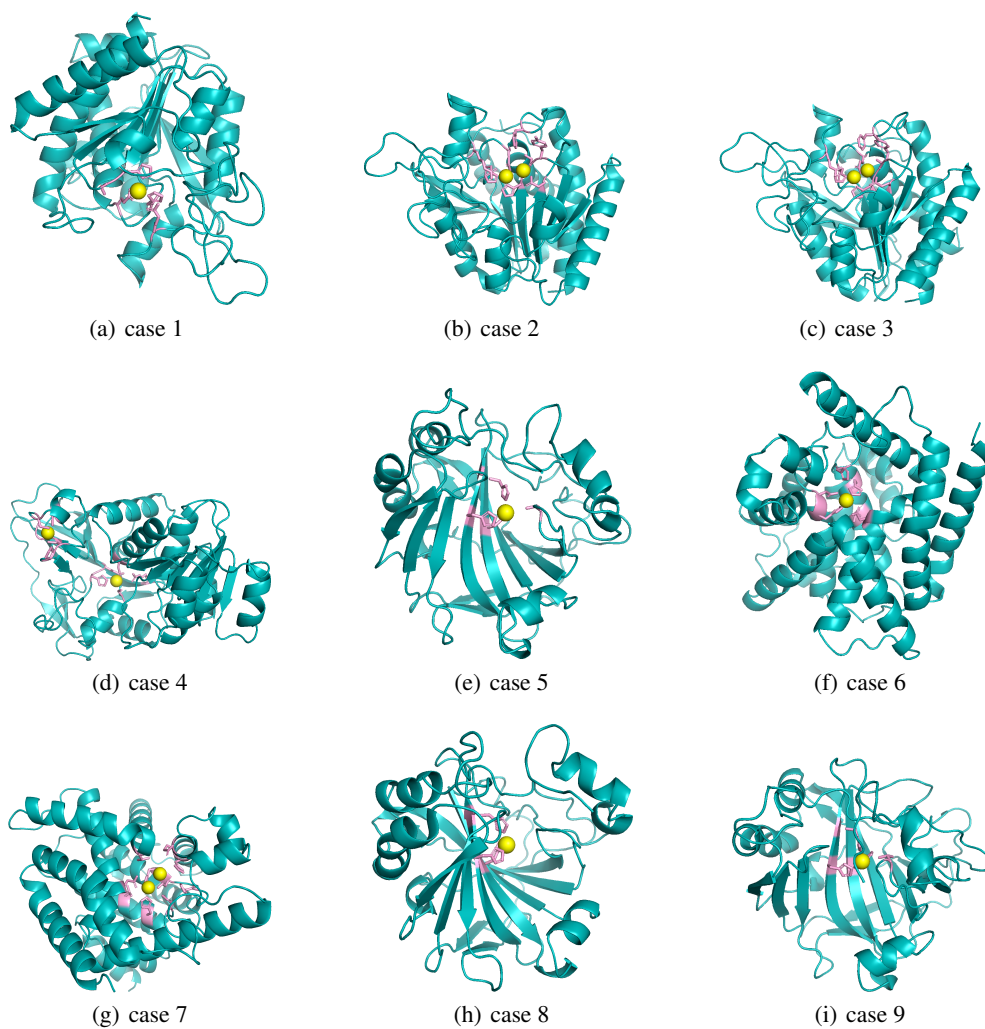


Figure 5: More designed cases for β -lactamase.

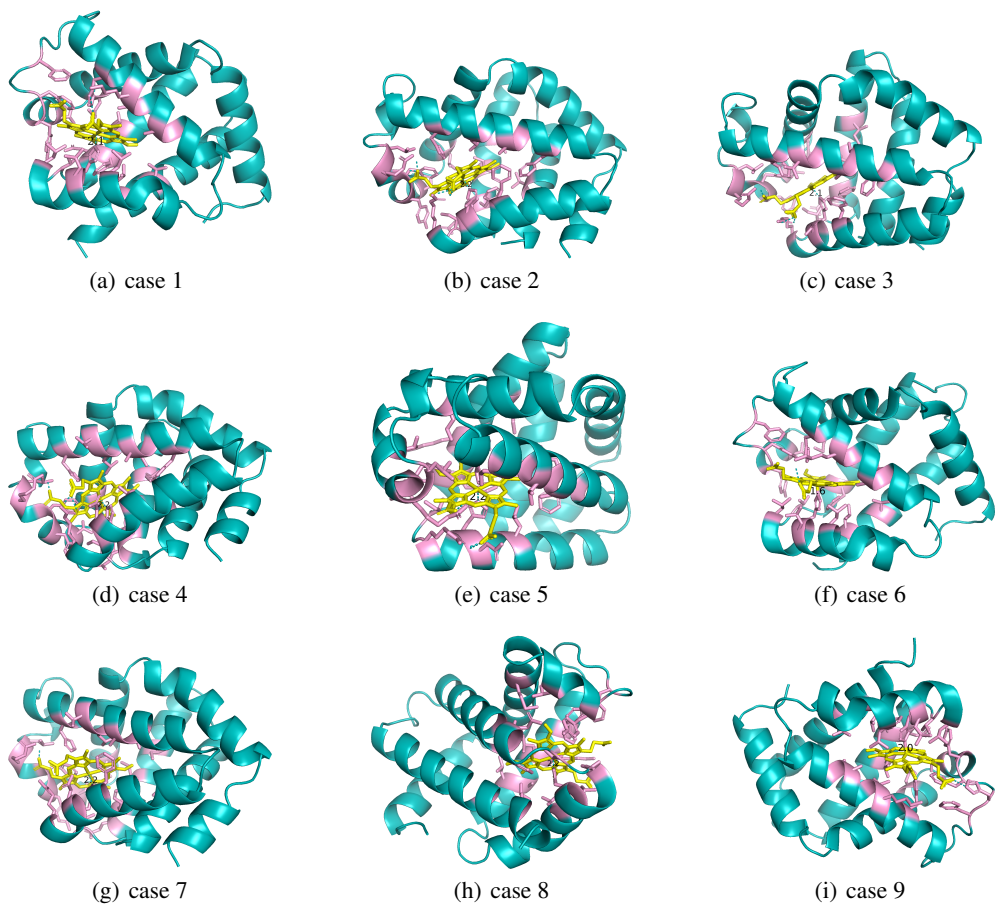


Figure 6: More designed cases for myoglobin.