
Generative models for graph-based protein design

John Ingraham, Vikas K. Garg, Regina Barzilay, Tommi Jaakkola
Computer Science and Artificial Intelligence Lab, MIT
{ingraham, vgarg, regina, tommi}@csail.mit.edu

Presenter: Yukang Yang

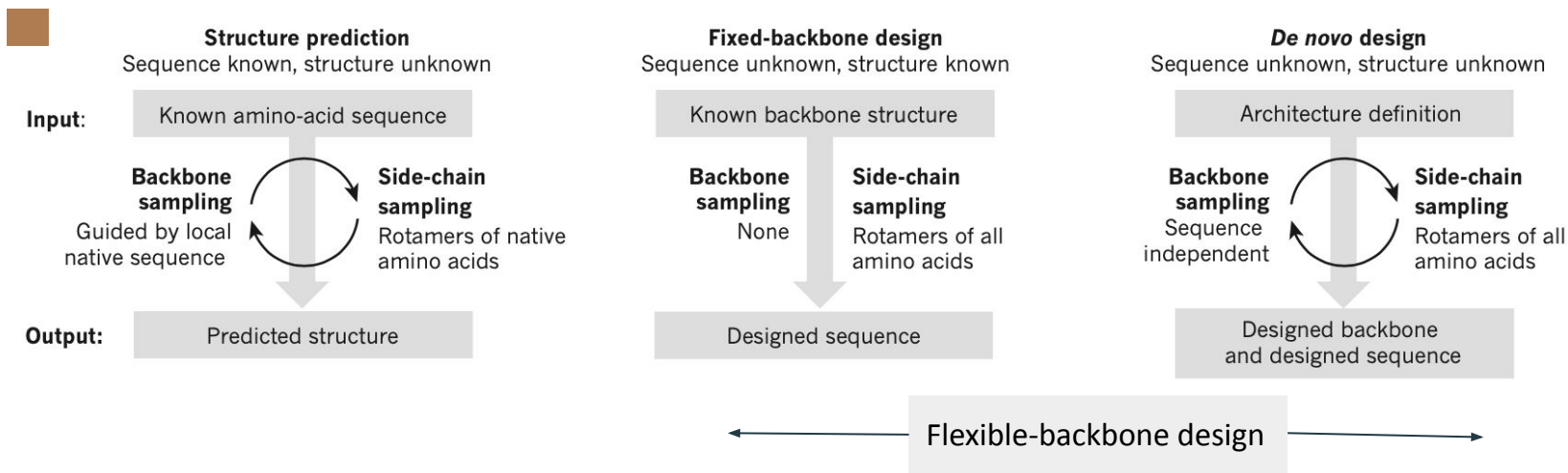
Paper: <https://openreview.net/pdf?id=ByMEAHrgLB>
Codes: <https://github.com/jingraham/neurips19-graph-protein-design>

Introduction

- Computational Protein Design (i.e., *Inverse protein folding problem*)

$$\mathcal{F}_\theta : \mathcal{X} \mapsto \mathcal{S} \quad \text{the amino acids sequence} \quad \mathcal{S} = \{s_i : 1 \leq i \leq n\}$$

- Desired **structure** $\mathcal{X} = \{x_i \in \mathbb{R}^3 : 1 \leq i \leq n\}$
- Desired functional properties



● Related Work

Bottom-up: optimizing Energy Function

Rosetta



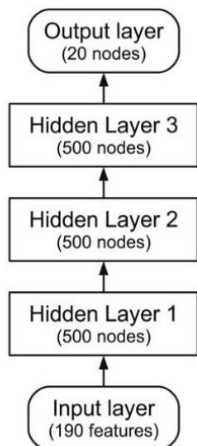
Top-down: *Conditional* Generative Model

MLP-based

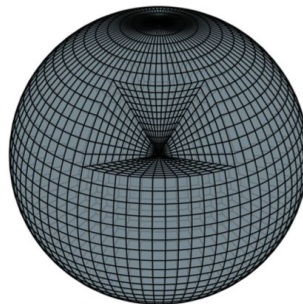
CNN-based

Graph-based

SPIN2



- Cons: Feature Design



Cons:

- Slow inference
- Complex preprocessing

Representing protein structure as a graph

Pros:

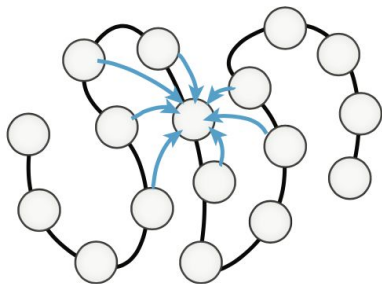
- Computational efficiency
- Inductive bias
- Representational flexibility

Method

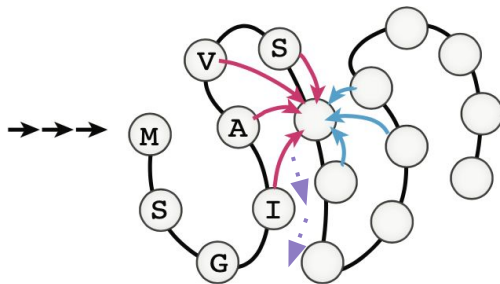
Structured Transformer

A

Structure Encoder



Sequence Decoder (autoregressive)



Information flow

→ Structure

→ Structure and sequence

○ Node (amino acid)

∩∩ Backbone

$$\mathcal{G} = (\mathcal{V}, \mathcal{E})$$

$$\mathcal{V} = \{v_1, \dots, v_N\}$$

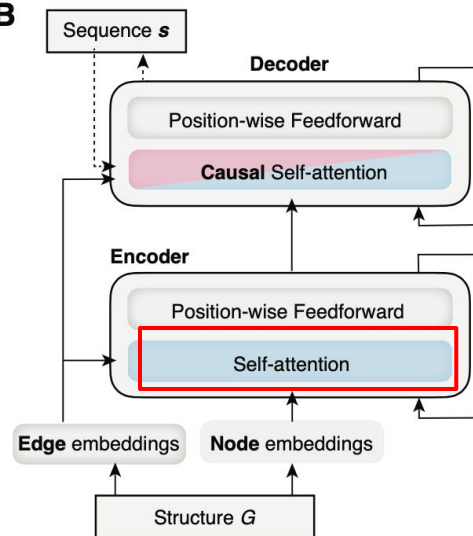
Node Features

$$\mathcal{E} = \{e_{ij}\}_{i \neq j}$$

Edge Features

$$p(\mathbf{s}|\mathbf{x}) = \prod_i p(s_i|\mathbf{x}, \mathbf{s}_{<i})$$

B



Inductive Bias:

- Invariance
- Locally Informative

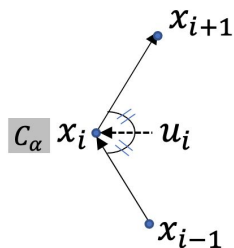


Relative Spatial Encodings

Local Frame

$$O_i = [b_i \quad n_i \quad b_i \times n_i]$$

$$u_i = \frac{x_i - x_{i-1}}{\|x_i - x_{i-1}\|}, \quad b_i = \frac{u_i - u_{i+1}}{\|u_i - u_{i+1}\|}, \quad n_i = \frac{u_i \times u_{i+1}}{\|u_i \times u_{i+1}\|}$$

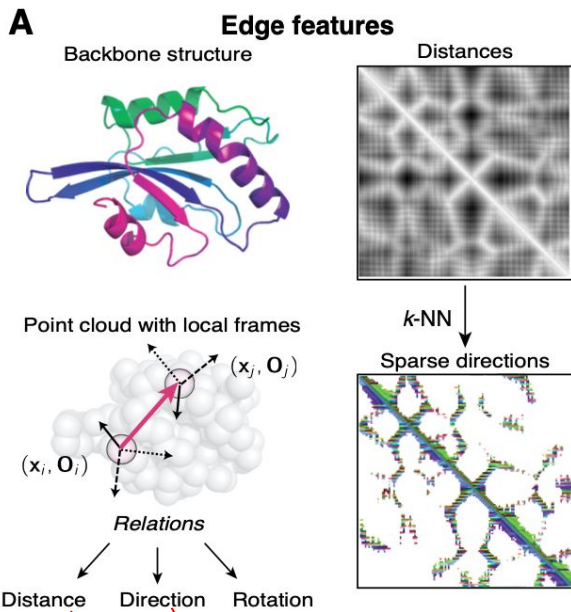


Structural Encoding

$$e_{ij}^{(s)} = \left(\mathbf{r}(\|x_j - x_i\|), \quad O_i^T \frac{x_j - x_i}{\|x_j - x_i\|}, \quad \mathbf{q}(O_i^T O_j) \right)$$

Edge Features

Positional Encoding



k-NN (sparsity)

$$\{e_{ij}\}_{j \in N(i,k)} \quad (k=30)$$

long-range dependencies in sequence

short-range in 3D space



Attention:

- Encoder

$$a_{ij}^{(\ell)} = \frac{\exp(m_{ij}^{(\ell)})}{\sum_{j' \in N(i,k)} \exp(m_{ij'}^{(\ell)})} \quad m_{ij}^{(\ell)} = \frac{\mathbf{q}_i^{(\ell)\top} \mathbf{z}_{ij}^{(\ell)}}{\sqrt{d}}$$

$$\mathbf{h}_i^{(\ell)} = \sum_{j \in N(i,k)} a_{ij}^{(\ell)} \mathbf{v}_{ij}^{(\ell)}$$

Key & Value

$$\mathbf{r}_{ij} = (\mathbf{h}_j, \mathbf{e}_{ij})$$

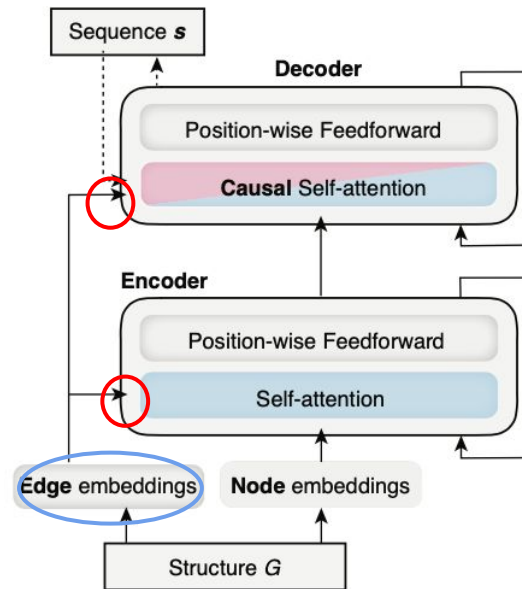
$$\mathbf{v}_{ij}^{(\ell)} = W_v^{(\ell)} \mathbf{r}_{ij}$$

$$\mathbf{z}_{ij}^{(\ell)} = W_z^{(\ell)} \mathbf{r}_{ij}$$

- Decoder

Causally consistent

$$\mathbf{r}_{ij}^{(\text{dec})} = \begin{cases} (\mathbf{h}_j^{(\text{dec})}, \mathbf{e}_{ij}, \mathbf{g}(s_j)) & i > j \\ (\mathbf{h}_j^{(\text{enc})}, \mathbf{e}_{ij}, \mathbf{0}) & i \leq j \end{cases}$$



Results

- Dataset: **CATH 4.2**

a **hierarchical** domain classification of the three-dimensional (3D) structures of proteins

- ❖ **C**lass: secondary structure content
- ❖ **A**rchitecture: shape revealed by the orientations of the secondary structure units
- ❖ **T**opology: sequential connectivity of secondary structure elements
- ❖ **H**omologous superfamily: whether the domains are evolutionarily related

The screenshot shows the homepage of the CATH / Gene3D v4.3 database. At the top, there is a navigation bar with links for Home, Search, Browse, Download, About, and Support, along with a search bar labeled 'Search CATH by keywords or ID'. The main header features the text 'CATH / Gene3D v4.3' and '151 million protein domains classified into 5,841 superfamilies'. Below this is another search bar with the text 'Search by keywords, PDB code, GO term, etc' and a 'Search' button. A banner below the search bar states: 'Core classification files for the latest version of CATH-Plus (v4.3) are now available to download. Daily updates of our very latest classifications are also available.' The page is divided into three main sections: '3D Structure' with a 'Find out more' button and a 'Go' button, 'Protein Evolution' with a 'Find out more' button and a 'Go' button, and 'Protein Function' with a 'Find out more' button and a 'Go' button. Each section includes a small icon representing its respective topic.

Structure-split setting

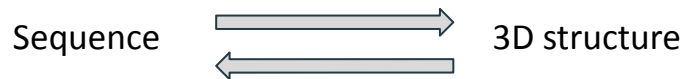
- ❖ Training set: 18024 chains
- ❖ Validation set: 608 chains
- ❖ Test set: 1120 chains

No CAT overlap

<http://www.cathdb.info>

- Evaluation

Single mutations may cause a protein to break or misfold



Many protein sequences may design the **same** 3D structure

Sequence Similarity **x**

- ❖ Likelihood-based: Perplexity
- ❖ Native sequence recovery
- ❖ Experimental comparison

❖ Likelihood-based: Perplexity

Table 1: **Null perplexities** for common statistical models of proteins.

Null model	Perplexity	Conditioned on
Uniform	20.00	-
Natural frequencies	17.83	Random position in a natural protein
Pfam HMM profiles	11.64	Specific position in a specific protein family

Perplexity $\propto 1/\text{probability}$

$$\log(\text{perplexity}(S)) = -\frac{1}{m} \sum_{i=1}^m \log p(w_i | w_1, \dots, w_{i-1})$$

Table 2: **Per-residue perplexities for protein language modeling** (lower is better). The protein chains have been cluster-split by CATH topology, such that test includes only unseen 3D folds. While a structure-conditioned language model can generalize in this structure-split setting, unconditional language models struggle.

Test set	Short	Single chain	All
Structure-conditioned models			
Structured Transformer (ours)	8.54	9.03	6.85
SPIN2 [8]	12.11	12.61	-
Language models			
LSTM ($h = 128$)	16.06	16.38	17.13
LSTM ($h = 256$)	16.08	16.37	17.12
LSTM ($h = 512$)	15.98	16.38	17.13
Test set size	94	103	1120

$$\text{Perplexity} = e^{\frac{\sum_i^m L_{nll}}{m}}$$

Ablations:

ProteinMPNN

Message Passing Neural Networks (MPNN)

$$\Delta h_i = \sum_j \text{MLP}(h_i, h_j, e_{ij})$$

Node features	Edge features	Aggregation	Short	Single chain	All
Rigid backbone					
Dihedrals	Distances, Orientations	Attention	8.54	9.03	6.85
Dihedrals	Distances, Orientations	PairMLP	8.33	8.86	6.55
C _α angles	Distances, Orientations	Attention	9.16	9.37	7.83
Dihedrals	Distances	Attention	9.11	9.63	7.87
Flexible backbone					
C _α angles	Contacts, Hydrogen bonds	Attention	11.71	11.81	11.51
	SPIN2 [8]		12.11	12.61	-

❖ Native sequence recovery

Method	Recovery (%)	Speed (AA/s) CPU	Speed (AA/s) GPU
Rosetta 3.10 fixbb	17.9	4.88×10^{-1}	N/A
Ours ($T = 0.1$)	27.6	2.22×10^2	1.04×10^4

(a) Single chain test set (103 proteins)

Method	Recovery (%)
Rosetta, fixbb 1	33.1
Rosetta, fixbb 2	38.4
Ours ($T = 0.1$)	39.2

(b) Ollikainen benchmark (40 proteins)

- More Accurate
- Faster

Table 4: **Improved reliability and speed compared to Rosetta.** (a) On the ‘single chain’ test set, our model more accurately recovers native sequences than Rosetta fixbb with greater speed (CPU: single core of Intel Xeon Gold 5115, GPU: NVIDIA RTX 2080). This set includes NMR-based structures for which Rosetta is known to not be robust [46]. (b) Our model also performs favorably on a prior benchmark of 40 proteins. All results reported as median of average over 100 designs.

$$p(\mathbf{s}|\mathbf{x}) = \prod_i p(s_i|\mathbf{x}, \mathbf{s}_{<i})$$



$$p^{(T)}(\mathbf{s}|\mathbf{x}) = \prod_i \frac{p(s_i|\mathbf{x}, \mathbf{s}_{<i})^{1/T}}{\sum_a p(a|\mathbf{x}, \mathbf{s}_{<i})^{1/T}}$$

Biased sampling

Rosetta : state-of-the-art framework for computational protein design

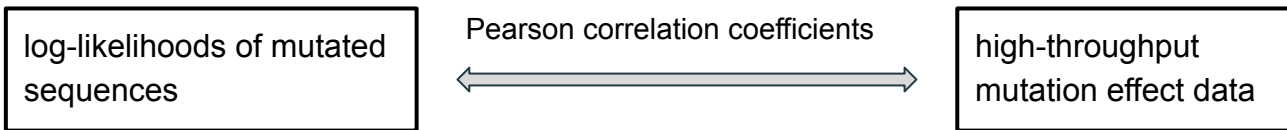
❖ Experimental comparison

Mutation effects

Table 5: **Structure-conditioned likelihoods correlate with mutation effects in *de novo*-designed miniproteins.** Shown are Pearson correlation coefficients (R , higher is better) between the log-likelihoods of mutated sequences and high-throughput mutation effect data from a systematic design of miniproteins [6]. Each design (column) includes 775 experimentally tested mutant protein sequences.

Design	$\beta\beta\alpha\beta\beta_{37}$	$\beta\beta\alpha\beta\beta_{1498}$	$\beta\beta\alpha\beta\beta_{1702}$	$\beta\beta\alpha\beta\beta_{1716}$	$\alpha\beta\beta\alpha_{779}$
Rigid backbone	0.47	0.45	0.12	0.47	0.57
Flexible backbone	0.50	0.44	0.17	0.40	0.56

Design	$\alpha\beta\beta\alpha_{223}$	$\alpha\beta\beta\alpha_{726}$	$\alpha\beta\beta\alpha_{872}$	$\alpha\alpha\alpha_{134}$	$\alpha\alpha\alpha_{138}$
Rigid backbone	0.36	0.11	0.21	0.24	0.33
Flexible backbone	0.33	0.21	0.23	0.36	0.41



Take-away

Structured Transformer

- **Graph-based** Transformer
- + Inductive Bias: 3D **structural** encodings, spatial **locality**
- Improved **perplexities**
- Compared to the SOTA protein design program, more **accurate and faster**

*Showing the potential of being able to **efficiently** design and engineer protein sequence with **structurally-guided** generative models...*

Q & A