Global-Context Aware Generative Protein Design for Structure-to-sequence Translation

Anonymous ACL submission

Abstract

The linear sequence of amino acids determines protein structure and function. Protein design, known as the inverse of protein structure predic-004 tion, aims to obtain a novel protein sequence that will fold into the defined structure. Recent works on computational protein design have studied designing sequences for the desired backbone structure with local positional information and achieved competitive performance. However, similar local environments in different backbone structures may result in different amino acids, which indicates the global context of protein structure matters. Thus, we 014 propose the Global-Context Aware generative protein design method (GCA), consisting of lo-016 cal and global modules. While local modules focus on relationships between neighbor amino 017 acids, global modules explicitly capture nonlocal contexts. Experimental results demonstrate that GCA achieves state-of-the-art performance on structure-based protein design. Our code and pretrained model will be released.

1 Introduction

034

040

Computational protein design, which aims to invent protein molecules with desired structures and functions automatically, has a wide range of applications in therapeutics and pharmacology (Brunette et al., 2015; Huang et al., 2016; Langan et al., 2019). Recent years have witnessed remarkable advancements (Gao et al., 2020; Kryshtafovych et al., 2019; Yang et al., 2019). While classical protein design approaches depend on composite energy functions of protein physics and sampling algorithms for exploring sequence and structure spaces, data-driven approaches apply neural networks to generate protein sequences with less prior knowledge.

Designing a protein sequence for a given structure remains challenging, as mapping the structure space to the sequence space is difficult. Current data-driven methods (Ingraham et al., 2019; Jing et al., 2021; Cao et al., 2021) agree on the assumption based on biology and physics prior knowledge that, for each amino acid, its neighborhoods have the most immediate and vital effects on itself. The majority of such methods represent protein structures as graphs with hand-crafted features and aggregate local messages. The computational protein design is formulated to learn features from 3D structures with the local message passing mechanism. However, *the similar local environment in different proteins may correspond to different amino acids*. Local neighbors do matter (Ribeiro et al., 2020), but it is not enough to obtain highquality protein sequences.



Figure 1: The comparison between the local module and global module. Information flows from adjacent nodes, and distant nodes are denoted as green arrows and red arrows, respectively. The red dashed arrows indicate the implicit information flow from distant nodes.

To fully explore the non-local information, we propose the Global-Context Aware generative protein design method (GCA) with local and global modules. While local modules are built upon graph attention networks that aggregate local messages gained from neighbors, global modules extend local graph attention to global self-attention in the form of Transformer (Vaswani et al., 2017) architectures. As shown in Fig. 1, the local module focuses on adjacent structure information though distant nodes can deliver information implicitly; the global module explicitly gathers information from distant nodes. By composing multiple blocks of local and global modules, GCA can capture highorder dependencies between protein sequences and structures in both neighbor-level and overall-level. 042

043

044

045

046

047

051

052

054

2 **Proposed method**

Preliminaries 2.1

071

073

100

101

104

105

106

107

108

Protein primary structure is the linear sequence of amino acids, typically notated as a string of letters. A protein sequence $S^N = \{(a)^N | a \in$ $\{A, R, N, ...V\}$ has N amino acids while each of them is represented by a letter of twenty possible letters such as A, R, N, D, C, Q, E, G, H, I, L, K, M, F, P, S, T, W, Y, and V. A protein sequence will fold into a protein tertiary structure $\mathcal{X}^N = \{\mathbf{x}_i^{\omega} \in \mathbb{R}^3 :$ $1 \le i \le N, \omega \in \{C\alpha, C, N, O\}\}$, where N is the number of amino acids and ω indicates the chain in protein. As shown in Fig. 2, protein design predicts the protein sequence of a given protein structure, while structure prediction is the opposite.



Figure 2: Two important tasks in protein modeling: structure prediction and protein design. The visualization is created by Mol* Viewer (Sehnal et al., 2021).

2.2 Represent protein as a graph

The structure of a protein is represented as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where node feature $\mathbf{v}_i \in \mathcal{V}$ corresponds to an amino acid while edge feature $\mathcal{E} = {\mathbf{e}_{ij}}_{j \in \mathcal{N}_i}$ suggests the rotation-invariant and translation-invariant relationships between each pair of nodes \mathbf{v}_i and \mathbf{v}_j . In particular, \mathcal{N}_i denotes the K-nearest neighbors of node i calculated by Euclidean distances of the backbone.

For node features, we construct three dihedral angles $\{\phi_i, \psi_i, \omega_i\}$ of the protein backbone from $C_{i-1}, N_i, C\alpha_i, C_i$, and N_{i+1} . Then these dihedral angels are embedded on the 3-torus as $\mathbf{v}_i =$ $\{\sin, \cos\} \times \{\phi_i, \psi_i, \omega_i\}.$

For edge features, we focus on describing relative spatial relationships between amino acids that satisfy rotation-invariant and translation-invariant properties. To simplify the computation, we only consider the position $\mathbf{x}_i^{C\alpha}$ of the alpha carbon $C\alpha$ as it's the central carbon atom in each amino acid. The distance $\|\mathbf{x}_{i}^{C\alpha} - \mathbf{x}_{i}^{C\alpha}\|_{2}, \forall i \neq j$ is encoded by Gaussian radial basis functions $r(\cdot)$. Then, as shown in Fig. 3, the direction is encoded by



Figure 3: A general view of how the local coordinate system is built. The final coordinate satisfies $\mathbf{b}_i \perp$ $\mathbf{n}_i, \mathbf{b}_i \times \mathbf{n}_i \perp \mathbf{n}_i, \mathbf{b}_i \times \mathbf{n}_i \perp \mathbf{b}_i.$

$$\mathbf{O}_{i}^{T} \frac{\mathbf{x}_{j}^{C\alpha} - \mathbf{x}_{i}^{C\alpha}}{\|\mathbf{x}_{j}^{C\alpha} - \mathbf{x}_{i}^{C\alpha}\|} \text{ while } \mathbf{O}_{i} = \begin{bmatrix} \mathbf{b}_{i} & \mathbf{n}_{i} & \mathbf{b}_{i} \times \mathbf{n}_{i} \end{bmatrix} \text{ de-}$$
fines a local coordinate system for each amino acid 110 by: 111

11

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

132

133

134

135

136

137

138

139

$$\mathbf{u}_{i} = \frac{\mathbf{x}_{i}^{C\alpha} - \mathbf{x}_{i-1}^{C\alpha}}{\|\mathbf{x}_{i}^{C\alpha} - \mathbf{x}_{i-1}^{C\alpha}\|}, \mathbf{b}_{i} = \frac{\mathbf{u}_{i} - \mathbf{u}_{i+1}}{\|\mathbf{u}_{i} - \mathbf{u}_{i+1}\|}, \mathbf{n}_{i} = \frac{\mathbf{u}_{i} \times \mathbf{u}_{i+1}}{\|\mathbf{u}_{i} \times \mathbf{u}_{i+1}\|}.$$
(1)

The orientation is encoded by the commonused quaternion representation of rotation matrix $q(O_i^T O_i)$. Thus, the edge feature e_{ij} is the concatenation of the distance, direction and orientation encodings as:

$$\mathbf{e}_{ij} = \left(\mathbf{r}(\|\mathbf{x}_{j}^{C\alpha} - \mathbf{x}_{i}^{C\alpha}\|_{2}), \mathbf{O}_{i}^{T} \frac{\mathbf{x}_{j}^{C\alpha} - \mathbf{x}_{i}^{C\alpha}}{\|\mathbf{x}_{j}^{C\alpha} - \mathbf{x}_{i}^{C\alpha}\|}, \mathbf{q}(\mathbf{O}_{i}^{T}\mathbf{O}_{j})\right).$$
(2)

Network architecture 2.3

2.3.1 Local module

ι

The local module is a graph neural network (GNN) that aggregates both node embeddings and local edge embeddings and updates the node embedding for further sequence generations. Considering a L-layer GNN, the key operations aggregating and updating can be formulated as follows:

$$\mathbf{h}_{\mathcal{N}_{i}}^{(l)} = aggregating^{(l)}(\{\left(\mathbf{h}_{i}^{(l-1)}, \mathbf{h}_{j}^{(l-1)}, \mathbf{h}_{\mathbf{e}_{ij}}\right) : j \in \mathcal{N}_{i}\}),$$
(3)

$$\mathbf{h}_{i}^{(l)} = updating(\mathbf{h}_{i}^{(l-1)}, \mathbf{h}_{\mathcal{N}_{i}}^{(l)}), \tag{4}$$

where $\mathbf{h}_{i}^{(l)} \in \mathbb{R}^{D}$ denotes the embedding of node i on the l-th layer, $\mathbf{h}_{\mathcal{N}_{i}}^{(l)} \in \mathbb{R}^{K \times D}$ denotes the local edge embedding of node i's neighbors on the l-th layer, K is the number of local neighbors, and Dis the dimensions of the embedding. In particular, $\mathbf{h}_{i}^{(0)} \in \mathbb{R}^{D}$ is the embedding of \mathbf{v}_{i} , and $\mathbf{h}_{\mathbf{e}_{ij}} \in \mathbb{R}^{D}$ is the embedding of the edge feature e_{ij} . The local edge information flows into node embeddings at each layer, while distant edge information flows through high-level layers.

In order to capture the relationships in local neighborhoods, we generalize graph attention scheme that take advantage of attention coefficients $\alpha \in \mathbb{R}^{N \times K}$ as strong relational inductive bias. Specifically, the attention coefficients are calculated as follows:

$$\alpha_{ij} = \frac{\exp(c_{ij})}{\sum_{k \in \mathcal{N}(i)} \exp(c_{ik})}, \forall j \in \mathcal{N}_i,$$
(5)

where c_{ij} is expressed as:

140

141

142

143

144

145

146

147

148

150

151

152

154

155

159

161

162

163

164

165

166

168

169 170

171

172

173

174

175

176

177

$$c_{ij} = \sigma \left(a^T \left[\mathbf{W} \mathbf{h}_i^{(l-1)} \mid \mid \mathbf{W} \mathbf{h}_j^{(l-1)} \mid \mid \mathbf{W} \mathbf{h}_{\mathbf{e}_{ij}} \right] \right), \forall j \in \mathcal{N}_i,$$
(6)

and $\mathbf{W} \in \mathbb{R}^{D \times D}$, $a \in \mathbb{R}^{3D}$ are learnable parameters, σ is the activation function, || is the concatenation operation.

Thus, the *aggregating* operation is adopted as:

$$\mathbf{h}_{\mathcal{N}_{i}}^{(l)} = \sum_{j \in \mathcal{N}_{i}} \alpha_{ij} \mathbf{W}_{r} \big[\mathbf{h}_{i}^{(l-1)} \mid\mid \mathbf{h}_{j}^{(l-1)} \mid\mid \mathbf{h}_{e_{ij}} \big], \quad (7)$$

where $\mathbf{W}_r \in \mathbb{R}^{D \times 3D}$ encodes the relation between *i* and *j*. The *updating* operation is simply renovating hidden layers by their local neighbors: $\mathbf{h}_i^{(l)} = \mathbf{h}_{M}^{(l)}$.

2.3.2 Global module

The global module is the fully self-attention network that generalizes Transformer (Vaswani et al., 2017) to protein graph. Specifically, the attention coefficients are calculated as follows:

$$\alpha_{ij} = \frac{\exp(c_{ij})}{\sum_{k \in \mathcal{V}} \exp(c_{ik})},\tag{8}$$

where c_{ij} is expressed as:

$$c_{ij} = \frac{1}{\sqrt{d}} \left(\mathbf{W}_q \mathbf{h}_i^{(l-1)} \right)^T \left(\mathbf{W}_k \left[\mathbf{h}_i^{(l-1)} \parallel \mathbf{h}_j^{(l-1)} \parallel \mathbf{h}_{\mathbf{e}_{ij}} \right] \right),$$
(0)

where $\mathbf{W}_q \in \mathbb{R}^{D \times D}$, $\mathbf{W}_k \in \mathbb{R}^{D \times 3D}$ are parameter matrices for the query and key, and d is a scale factor. Then, the *aggregating* operation is formulated as:

$$\mathbf{h}_{\mathcal{N}_i}^{(l)} = \sum_{j \in \mathcal{V}} \alpha_{ij} \mathbf{W}_r \mathbf{h}_j^{(l-1)}, \qquad (10)$$

the *updating* operation is defined by employing layer normalization (LayerNorm), dropout (DropOut) and fully connected networks (FFN):

$$\mathbf{h}_{i}^{(l)} = \text{LayerNorm}(\mathbf{h}_{\mathcal{N}_{i}}^{(l)} + \text{DropOut}(\text{FFN}(\mathbf{h}_{\mathcal{N}_{i}}^{(l)}))). (11)$$

The overall architecture with stacked local modules and global modules is shown in Fig. 4.

3 Experiments

3.1 Experimental settings

3.1.1 Dataset

We use the CATH 4.2 dataset collected by (Ingraham et al., 2019) for evaluation. This dataset obtains full chains up to length 500, and structures have been partitioned with 40% non-redundancy by their CATH (Class, Architecture, Topology, Homologous). With no CAT overlap between sets, there are 18024 chains in the training set, 608 chains in the validation set, and 1120 chains in the test set, respectively. Two subsets of the entire test set are evaluated simultaneously: a 'Short' subset containing chains up to length 100 and a 'Single chain' subset for comparing with baselines that only use the single chain. We also consider a smaller dataset TS50, which is the standard benchmark introduced by (Li et al., 2014). Though the model is still trained on the CATH 4.2 dataset, we filter the training and validation sets to ensure there is no overlap with TS50.

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

199

200

201

202

204

205

207

208

209

210

211

212

213

214

215

216

217

218

219

3.1.2 Measurement

Perplexity Following (Ingraham et al., 2019; Madani et al., 2020), we define the perplexity that evaluates the predicted protein sequences from natural language perspective:

$$\operatorname{PERP}(\mathcal{S}^{N}, \mathcal{X}^{N}) = \exp\left(-\frac{1}{N}\sum_{i=1}^{N}\mathcal{S}_{i}^{N}\log p\left(\mathcal{S}_{i}^{N} \mid \mathcal{X}_{i}^{N}\right)\right),$$
(12)

where $(\mathcal{S}^N, \mathcal{X}^N)$ is the sequence-structure pair of a protein with N amino acids. $\mathcal{S}_i^N, \mathcal{X}_i^N$ denote the *i*-th amino acid in sequence and structure respectively. $p(\mathcal{S}_i^N | \mathcal{X}_i^N)$ is the output probability from the model.

Recovery To evaluate the predicting accuracy of the protein sequence at per-residue level, we consider the recovery:

$$\operatorname{REC}(\mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathcal{X}^N, \mathcal{S}^N) \in \mathcal{D}} \frac{1}{N} \sum_{i=1}^{N}$$
(13)
$$\mathbb{1}[S_i^N = \arg\max p(S_i^N \mid \mathcal{X}_i^N)],$$

where \mathcal{D} denotes the whole dataset.

3.1.3 Model architecture and optimization

In all experiments, GCA model is built by three blocks of local and global modules for the encoder and decoder with the hidden dimension of 128. The Adam optimizer with learning rate of $1e^{-3}$ is employed. Models are trained for 100 epochs while each batch contains up to 2,500 characters.



Figure 4: The architecture of our proposed method.

3.2 Experimental results

We present the median of PERP in Table 1. While the structure-free language model LSTMs produce confusing protein sequences, structure-based models obtain less-perplex protein sequences, indicating the importance of structural features. GCA outperforms other models as global contexts of protein structures are taken into account.

Methods	Short	Single chain	All
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
SPIN2	12.11	12.61	-
Structure-based models			
StructTrans	8.56	8.97	7.14
StructGNN	8.40	8.84	6.69
GCA	7.09	7.49	6.06

Table 1: Performance of different methods on CATH 4.2 dataset assessed by PERP (lower is better).

Though PERP matters from the perspective of natural language, REC that evaluates the ability of models in inferring sequences given determined structures is also crucial. We compare GCA with other structure-based models in Table 2.

Methods	Short	Single chain	All
StructTrans	31.59	30.35	33.90
StructGNN	30.90	30.85	35.25
GCA	32.25	33.04	36.11

Table 2: Performance of different methods on CATH4.2 dataset assessed by REC (higher is better).

GCA obtains the highest REC on all three sets among these structure-based methods. Moreover, the recovery of StructGNN and StructTrans drops significantly in 'Short' and 'Single chain' sets, which suggests they are overfitting on long sequences and multiple chains, while GCA performs consistently well on them. As few structural features can be explored in short sequnce and single chain, the prediction is relatively difficult. However, the global information in GCA makes up for the deficiency of structural features of short chains, making performance significantly improved.

245

247

248

249

250

251

252

253

254

255

257

259

260

261

262

263

264

265

266

268

269

270

271

272

To compare with other methods, we conduct experiments on the standard TS50 dataset and show the results in Table 3. The methods for comparison include the CNN-based ProDCoNN (Zhang et al., 2020), the distance-map-based SPROF (Chen et al., 2019), the graph-based GVP (Jing et al., 2021) the sequential representation method SPIN (Li et al., 2014) and SPIN2 (O'Connell et al., 2018), the constraint satisfaction method ProteinSolver (Strokach et al., 2020), and the popular method Rosetta. GCA achieves remarkable performance and outperforms other methods by a large margin.

Methods	REC
Rosetta	30.0
SPIN	30.3
ProteinSolver	30.8
SPIN2	33.6
StructTrans	36.1
StructGNN	38.0
SPROF	39.2
ProDCoNN	40.7
GVP	44.1
GCA	47.0

Table 3: Performance of different methods on TS50dataset assessed by REC (higher is better).

4 Conclusion

We introduce the consideration of global information and propose the global-context aware generative protein design method, consisting of local modules and global modules. The local module propagates neighborhood messages across layers, and the global module emphasizes long-term dependencies. Experimental results show that our proposed GCA method outperforms state-of-theart methods on benchmark datasets. In 'Short' and 'Single chain' sets, the global-context aware mechanism significantly improves the performance, indicating the potentials to promote structure-based protein design.

222

273

5

Limitations

References

Though our proposed GCA method has taken the

global context into consideration, there is still room

for further improvement in efficiency. A possible

solution is learning a global context vector for each

TJ Brunette, Fabio Parmeggiani, Po-Ssu Huang, Gira

Bhabha, Damian C Ekiert, Susan E Tsutakawa,

Greg L Hura, John A Tainer, and David Baker. 2015.

Exploring the repeat protein universe through compu-

tational protein design. Nature, 528(7583):580-584.

Chen, Igor Melnyk, and Yang Shen. 2021. Fold2seq:

A joint sequence(1d)-fold(3d) embedding-based gen-

erative model for protein design. In Proceedings of

the 38th International Conference on Machine Learn-

ing, volume 139 of Proceedings of Machine Learning

Sheng Chen, Zhe Sun, Lihua Lin, Zifeng Liu, Xun Liu,

Yutian Chong, Yutong Lu, Huiying Zhao, and Yue-

dong Yang. 2019. To improve protein sequence pro-

file prediction through image captioning on pairwise

residue distance map. Journal of chemical informa-

Wenhao Gao, Sai Pooja Mahajan, Jeremias Sulam, and

tural modeling and design. Patterns, 1(9):100142.

Po-Ssu Huang, Scott E Boyken, and David Baker. 2016.

John Ingraham, Vikas Garg, Regina Barzilay, and

Tommi Jaakkola. 2019. Generative models for graph-

based protein design. In Advances in Neural Infor-

mation Processing Systems, volume 32, pages 15820-

Bowen Jing, Stephan Eismann, Patricia Suriana,

Raphael John Lamarre Townshend, and Ron Dror.

2021. Learning from protein structure with geomet-

ric vector perceptrons. In International Conference

Andriy Kryshtafovych, Torsten Schwede, Maya Topf,

Krzysztof Fidelis, and John Moult. 2019. Critical

assessment of methods of protein structure prediction

(casp)-round xiii. Proteins: Structure, Function,

Robert A Langan, Scott E Boyken, Andrew H Ng, Jen-

nifer A Samson, Galen Dods, Alexandra M West-

brook, Taylor H Nguyen, Marc J Lajoie, Zibo Chen, Stephanie Berger, et al. 2019. De novo design of bioactive protein switches. Nature, 572(7768):205-

and Bioinformatics, 87(12):1011-1020.

The coming of age of de novo protein design. Nature,

Jeffrey J. Gray. 2020. Deep learning in protein struc-

Research, pages 1261–1271. PMLR.

tion and modeling, 60(1):391-399.

537(7620):320-327.

15831. Curran Associates, Inc.

on Learning Representations.

210.

Yue Cao, Payel Das, Vijil Chenthamarakshan, Pin-Yu

protein, which we leave as future work.

- 274 275
- 276

- 279
- 284

- 290
- 293 294
- 297
- 298
- 301

304

306 307

310

311 312

- 313
- 314
- 315
- 316 317

319

320

Zhixiu Li, Yuedong Yang, Eshel Faraggi, Jian Zhan, and Yaoqi Zhou. 2014. Direct prediction of profiles of sequences compatible with a protein structure by neural networks with fragment-based local and energy-based nonlocal profiles. Proteins: Structure, Function, and Bioinformatics, 82(10):2565-2573.

325

326

328

331

332

334

335

336

337

338

339

340

341

342

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

361

362

363

364

365

366

367

369

- Ali Madani, Bryan McCann, Nikhil Naik, Nitish Shirish Keskar, Namrata Anand, Raphael R Eguchi, Po-Ssu Huang, and Richard Socher. 2020. Progen: Language modeling for protein generation. arXiv preprint arXiv:2004.03497.
- James O'Connell, Zhixiu Li, Jack Hanson, Rhys Heffernan, James Lyons, Kuldip Paliwal, Abdollah Dehzangi, Yuedong Yang, and Yaoqi Zhou. 2018. Spin2: Predicting sequence profiles from protein structures using deep neural networks. Proteins: Structure, Function, and Bioinformatics, 86(6):629-633.
- Leonardo FR Ribeiro, Yue Zhang, Claire Gardent, and Iryna Gurevych. 2020. Modeling global and local node contexts for text generation from knowledge graphs. Transactions of the Association for Computational Linguistics, 8:589-604.
- David Sehnal, Sebastian Bittrich, Mandar Deshpande, Radka Svobodová, Karel Berka, Václav Bazgier, Sameer Velankar, Stephen K Burley, Jaroslav Koča, and Alexander S Rose. 2021. Mol* viewer: modern web app for 3d visualization and analysis of large biomolecular structures. Nucleic Acids Research.
- Alexey Strokach, David Becerra, Carles Corbi-Verge, Albert Perez-Riba, and Philip M. Kim. 2020. Fast and flexible protein design using deep graph neural networks. Cell Systems, 11(4):402-411.e4.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In NeurIPS, pages 5998-6008.
- Kevin K Yang, Zachary Wu, and Frances H Arnold. 2019. Machine-learning-guided directed evolution for protein engineering. Nature methods, 16(8):687-694.
- Yuan Zhang, Yang Chen, Chenran Wang, Chun-Chao Lo, Xiuwen Liu, Wei Wu, and Jinfeng Zhang. 2020. Prodconn: Protein design using a convolutional neural network. Proteins: Structure, Function, and Bioinformatics, 88(7):819-829.