

---

## 000 A IMPLEMENTATION DETAILS

### 001 A.1 BASELINE EVALUATION

002 In our experiments, we compare Hier-IF with several approaches, including GraphTransIngraham  
003 et al. (2019), GVPJing et al. (2020), PiFoldGao et al. (2022a), AlphaDesignGao et al. (2022b),  
004 ProteinMPNNDauparas et al. (2022), ESM-IFHsu et al. (2022), LM-DesignZheng et al. (2023) and  
005 KW-DesignGao et al. (2023). We will describe the implementation as follow:  
006

- 007 • **GraphTrans:** They introduce conditional language models for protein sequence that directly  
008 condition on a graph specification of the target structure. They can efficiently captures the  
009 complex dependencies in proteins by focusing on those that are long-range in sequence but  
010 local in 3D space.
- 011 • **GVP:** They extend standard dense layers to operate on collections of Euclidean vectors.  
012 Utilizing such layers can help model to leverage the geometric and relational aspects.
- 013 • **PiFold:** They propose a novel residue featurizer to learn multi-scale residue interactions and  
014 PiFold can generate the sequence in one shot as a result it can achieve a great effectiveness.
- 015 • **AlphaDesign:**They use AlphaFoldDB to establish a new graph-based benchmark. Fur-  
016 thermore, they use a simplified graph transformer encoder to introduce protein angles and  
017 propose a confidence-aware protein decoder. Such modules can efficiently improve the  
018 performance.
- 019 • **ProteinMPNN:** They begins from a message passing neural network which can effectively  
020 explore the information of backbones. Then, they replaced the fixed decoding order with  
021 an order agnostic autoregressive model, which means the decoding order is random. Such  
022 strategy can efficiently improve the performance for inverse folding.
- 023 • **ESM-IF:**They augments the inverse folding dataset with folding data created by AlphaFold2.  
024 Training with additional data, the performance achieves a great progress.
- 025 • **LM-Design:** They conduct a lightweight structural adapter which is implanted into protein  
026 language models and endows it with structural awareness.
- 027 • **KW-Design:** They propose a knowledge-aware module that refines low-quality residues.  
028 Furthermore, they introduce a memory-retrieval mechanism to save more training time.

### 029 A.2 DATASETS

030 In our experiments, we evaluate the proposed Hier-IF model on two widely used protein domain  
031 datasets: CATH4.2Orengo et al. (1997) and CATH4.3Hsu et al. (2022). These datasets provide a  
032 reliable benchmark for assessing the capacity of models to generalize across diverse protein topologies.  
033 For CATH4.2, we follow the experimental setup described in GraphTransIngraham et al. (2019).  
034 We utilize all protein domains from the CATH4.2 40% non-redundant set. From these domains, we  
035 extract the corresponding full protein chains, filtering out those longer than 500 residues to control  
036 computational complexity and ensure consistency with prior work. The resulting dataset is then split  
037 based on CATH topology classification codes into training, validation, and test sets, maintaining  
038 an 80/10/10 ratio. This ensures that each topology appears in only one of the three sets, promoting  
039 rigorous generalization evaluation. For CATH4.3, we adopt the protocol introduced in ESM-IFHsu  
040 et al. (2022). Specifically, we split the topology classification codes from CATH4.3 into training,  
041 validation, and test sets using the same 80/10/10 ratio. We extract full chains of no more than 500  
042 residues, ensuring consistency across datasets and alignment with computational constraints.  
043

### 044 A.3 TRAINING

045 The models are trained up to 100 epochs by default using the Adam optimizer on NVIDIA A6000s.  
046 We use the training setting as ProteinMPNN, where the batch size is set to 6000 residues. Furthermore,  
047 we employ the Noam learning rate scheduler?, which is widely used in transformer-based models.  
048 This scheduler starts with a warm-up phase where the learning rate increases linearly, followed by a  
049 decay phase where it decreases proportionally to the inverse square root of the training step.  
050

---

## B IMPLEMENTATION DETAILS OF HIER-IF

### B.1 MASKED DIFFUSION MODEL

We formulate the mask-based diffusion process over discrete or structured protein representations (e.g., residue identities or torsion angles) using a forward process that progressively masks input tokens, and a reverse process that learns to recover them.

#### B.1.1 FORWARD PROCESS (MASKING)

Given a protein representation  $x_0 \in \mathcal{X}^L$ , where  $L$  is the sequence length, the forward process defines a sequence of increasingly masked inputs  $\{x_t\}_{t=0}^T$ . At each timestep  $t$ , a binary mask  $m_t \in \{0, 1\}^L$  is sampled from a masking schedule  $q(m_t | t)$ , and we define:

$$x_t = m_t \odot [\text{MASK}] + (1 - m_t) \odot x_{t-1}$$

where  $\odot$  denotes elementwise multiplication, and  $[\text{MASK}]$  is a special token representing an unknown residue or feature.

The marginal distribution of  $x_t$  given  $x_0$  can be written as:

$$q(x_t | x_0) = \mathbb{E}_{m_t}[x_t | x_0] = m_t \odot [\text{MASK}] + (1 - m_t) \odot x_0$$

#### B.1.2 REVERSE PROCESS (UNMASKING)

The denoising model  $p_\theta$  learns to reconstruct the original input by predicting the masked elements:

$$p_\theta(x_{t-1} | x_t) = \prod_{i \in \mathcal{M}_t} p_\theta(x_{t-1}^{(i)} | x_t)$$

where  $\mathcal{M}_t = \{i | m_t^{(i)} = 1\}$  is the set of positions masked at timestep  $t$ .

For protein sequence modeling, each term is modeled by a categorical distribution over the 20 standard amino acids:

$$p_\theta(x_{t-1}^{(i)} = a | x_t) = \text{softmax}(f_\theta(x_t))_a \quad (1)$$

where  $f_\theta$  is a neural network (e.g., a transformer) that outputs logits for each residue position.

### B.2 MASK GENERATOR

To enable a structure-aware and learnable masking strategy, we introduce a *mask generator* that produces a soft masking score for each residue at every diffusion timestep. The mask generator is implemented as a 3-layer multi-layer perceptron (MLP), which takes as input both the DSSP-derived secondary structure and the current diffusion timestep.

Each residue’s secondary structure  $s_i$  is encoded as a one-hot vector:

$$\text{onehot}(s_i) = \begin{cases} [1, 0, 0], & \text{if } s_i = \text{H (helix)} \\ [0, 1, 0], & \text{if } s_i = \text{E (strand)} \\ [0, 0, 1], & \text{if } s_i = \text{C (coil)} \end{cases}$$

The mask generator outputs a soft mask  $\tilde{m} \in (0, 1)^L$ , where each element  $\tilde{m}_i$  represents the probability of masking the  $i$ -th residue. To encourage discrete masking behavior during training, a hard mask  $m \in \{0, 1\}^L$  is sampled from  $\tilde{m}$  according to a Bernoulli distribution:

$$m_i \sim \text{Bernoulli}(\tilde{m}_i).$$

Since this sampling operation is inherently non-differentiable, it blocks gradient backpropagation through  $m$ . To mitigate this, we employ the Straight-Through Estimator (STE) technique, which treats the sampling step as an identity function during the backward pass. Concretely, the gradient is directly propagated from the loss with respect to the hard mask  $m$  back to the soft mask  $\tilde{m}$ :

$$\frac{\partial \mathcal{L}}{\partial \tilde{m}} \approx \frac{\partial \mathcal{L}}{\partial m}.$$

---

108 Alternatively, training can be conducted solely with the soft mask  $\tilde{m}$  to maintain full differentiability,  
109 at the cost of losing discrete mask interpretability.

110 This approach enables the mask generator to learn meaningful, time-dependent masking strategies  
111 that balance discrete perturbations with smooth optimization.  
112

113 —  
114 The diffusion timestep  $t$  is represented as a scalar, either as an integer in the diffusion schedule  
115  $t \in \{0, \dots, T\}$  or normalized to a continuous value in the range  $[0, 1]$ . This scalar is then mapped to  
116 a higher-dimensional embedding via a small learnable MLP:

$$117 \tau(t) = W_2 \cdot \text{ReLU}(W_1 t + b_1) + b_2, \quad \tau(t) \in \mathbb{R}^{d_t},$$

118 where  $W_1 \in \mathbb{R}^{h \times 1}$ ,  $W_2 \in \mathbb{R}^{d_t \times h}$ , and  $h$  denotes the hidden dimension size.  
119

120 For each residue  $i$ , the input vector to the mask generator is formed by concatenating the one-hot  
121 encoded secondary structure and the time embedding:

$$122 z_i = \text{concat}(\text{onehot}(s_i), \tau(t)) \in \mathbb{R}^{3+d_t}.$$

123  
124 The mask generator  $g_\phi$  is a 3-layer MLP that maps  $z_i$  to the soft mask score:

$$125 \tilde{m}_i = \sigma(W_3 \cdot \text{ReLU}(W_2 \cdot \text{ReLU}(W_1 z_i + b_1) + b_2) + b_3),$$

126 where  $\sigma(\cdot)$  denotes the sigmoid activation, ensuring  $\tilde{m}_i \in (0, 1)$  represents the masking probability  
127 for residue  $i$ .  
128

129 —  
130 The resulting soft mask  $\tilde{m} \in (0, 1)^L$  is sampled during training to produce a hard mask  $m$ , while at  
131 inference time, a deterministic mask can be obtained by thresholding or selecting the top- $k$  residues  
132 with highest mask scores.  
133

134 This design allows the model to learn residue-level, structure- and time-aware masking strategies that  
135 adapt dynamically according to both the diffusion stage and the protein secondary structure context.  
136

### 137 B.3 SOFT VS. HARD MASKING

138 We distinguish between two types of masking strategies: *soft* masking and *hard* masking, based on  
139 how the mask is applied during training and inference.  
140

141 **Soft masking.** In soft masking, the output  $\tilde{m} \in (0, 1)^L$  from the mask generator is directly  
142 interpreted as a probabilistic or fractional mask. This enables smooth, differentiable updates, and is  
143 especially useful when the mask is applied via interpolation:  
144

$$145 \hat{x} = \tilde{m} \odot \text{noise} + (1 - \tilde{m}) \odot x$$

146 Here, each position is softly corrupted in proportion to its predicted mask score. This allows gradients  
147 to flow through the masking mechanism during backpropagation, and facilitates end-to-end training.  
148

149 **Hard masking.** In hard masking, the soft mask  $\tilde{m}$  is used to sample a binary mask  $m \in \{0, 1\}^L$ ,  
150 typically as:

$$151 m_i \sim \text{Bernoulli}(\tilde{m}_i)$$

152 This sampled mask is then used to fully corrupt certain positions and leave others untouched.  
153 While this approach aligns more closely with classical denoising objectives and provides stronger  
154 perturbations, the sampling process is non-differentiable and introduces variance during training.  
155 Common workarounds include using the Straight-Through Estimator (STE) or Gumbel-softmax  
156 relaxation.

157 **Comparison.** Soft masking provides smoother training dynamics and supports gradient-based opti-  
158 mization of the masking policy. Hard masking, in contrast, introduces more discrete and interpretable  
159 corruption patterns, often preferred at inference time or when seeking sharper denoising effects.  
160

161 In our implementation, we employ soft masking during training to ensure gradient flow, and optionally  
apply hard masking at inference via either thresholding or top- $k$  selection based on the predicted  $\tilde{m}$ .

---

## 162 B.4 DSSP SECONDARY STRUCTURE ANNOTATION

163  
164 DSSP (Define Secondary Structure of Proteins) ? is a widely used algorithm for assigning secondary  
165 structure elements to each residue in a protein based on its 3D atomic coordinates. It categorizes  
166 residues into eight classes:

$$167 \{H, B, E, G, I, T, S, C\}$$

168 where

- 169 • H: Alpha-helix
- 170 • B: Isolated beta-bridge residue
- 171 • E: Extended strand, participates in beta ladder
- 172 • G:  $3_{10}$  helix
- 173 • I: Pi-helix
- 174 • T: Turn
- 175 • S: Bend
- 176 • C: Coil (none of the above)

179 For many applications, including ours, a simplified secondary structure representation is preferred.  
180 We map the original 8 classes into 3 broader categories, commonly referred to as *Helix (H)*, *Strand*  
181 *(E)*, and *Coil (C)*, as follows:

$$182 \begin{cases} H = \{H, G, I\} & \text{(Helices)} \\ E = \{E, B\} & \text{(Strands)} \\ C = \{T, S, C\} & \text{(Coils/Loops)} \end{cases}$$

183  
184  
185  
186  
187 This reduction simplifies the secondary structure annotation while preserving the key structural motifs  
188 relevant for downstream modeling and masking tasks.

## 189 C EVALUATION METRIC

190  
191  
192 **Median Recovery** Median Recovery measures the similarity between generated protein sequences  
193 and target reference sequences, commonly used to evaluate reconstruction accuracy. Given a set of  
194 sequence pairs  $\{(S_i^{\text{gen}}, S_i^{\text{ref}})\}_{i=1}^N$ , the recovery for each pair is defined as the proportion of matching  
195 amino acids at corresponding positions:

$$196 \text{Recovery}_i = \frac{1}{L_i} \sum_{j=1}^{L_i} \mathbb{I}(S_{i,j}^{\text{gen}} = S_{i,j}^{\text{ref}}),$$

197  
198 where  $L_i$  is the length of the  $i$ -th sequence and  $\mathbb{I}(\cdot)$  is the indicator function. Median Recovery is the  
199 median of all individual recoveries:

$$200 \text{Median Recovery} = \text{median}(\{\text{Recovery}_i\}_{i=1}^N).$$

201  
202 A higher median recovery indicates better agreement between generated and reference sequences.

203  
204 **Perplexity** Perplexity is a widely used metric to assess the quality of probabilistic sequence models,  
205 reflecting how well a model predicts a given protein sequence. For a sequence  $S = (s_1, s_2, \dots, s_L)$   
206 and a model parameterized by  $\theta$  that estimates conditional probabilities  $p_\theta(s_j | s_{<j})$ , the log-likelihood  
207 is:

$$208 \log p_\theta(S) = \sum_{j=1}^L \log p_\theta(s_j | s_{<j}).$$

209  
210 The perplexity of the sequence is defined as:

$$211 \text{Perplexity}(S) = \exp\left(-\frac{1}{L} \log p_\theta(S)\right).$$

212  
213 Lower perplexity indicates that the model predicts the sequence with higher confidence and better  
214 accuracy.  
215

**Designability** A protein backbone is considered *designable* if there exists at least one amino acid sequence that folds into that structure. Our evaluation of designability follows the methodology outlined by ?. For each backbone generated by a model, we produce eight candidate sequences using ProteinMPNN Dauparas et al. (2022) with a sampling temperature of 0.1. Each designed sequence is then folded using ESMFold Lin et al. (2023) to predict its 3D structure.

We calculate the root mean square deviation (RMSD) between each predicted structure and the model’s original backbone. The lowest RMSD among the eight predictions, termed the self-consistency RMSD (scRMSD), is used to evaluate designability. A backbone is classified as designable if its scRMSD is below a threshold of 2 Å.

The overall designability score for a model is computed as the fraction of generated backbones that meet this criterion:

$$\text{Designability} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}(\text{scRMSD}(S_i) < 2 \text{ \AA}),$$

where  $N$  is the total number of generated backbones and  $\mathbf{1}(\cdot)$  is the indicator function.

This approach quantitatively measures the practical feasibility of designing sequences that reliably fold back to the intended structures, thus reflecting the quality of the generative model.

**Structural Diversity via Foldseek** Structural diversity measures the extent to which generated protein sequences fold into a variety of distinct structural conformations, which is important for exploring novel folds and functions.

Foldseek ? is a state-of-the-art tool for fast and sensitive structural comparison and clustering of large protein datasets. After predicting 3D structures for designed sequences, Foldseek is used to perform pairwise structural alignments and group proteins into clusters based on their similarity.

Diversity is then quantified by metrics such as:

- The number of distinct structural clusters detected at a given similarity cutoff, reflecting the breadth of structural exploration.
- Average pairwise structural dissimilarity (e.g., average RMSD or 1 - TM-score) among the generated structures.

## D DIFFERENT SAMPLING STRATEGY

Method	Recovery			Perplexity		
	Short	Single-Chain	All	Short	Single-Chain	All
Iterative Refinement	0.40	0.41	0.53	6.97	6.02	5.80
Time-aware Mask Generator	<b>0.46</b>	<b>0.46</b>	<b>0.59</b>	<b>5.46</b>	<b>5.17</b>	<b>4.01</b>

Masked diffusion models generate discrete sequences by iteratively predicting masked tokens. While these models are efficient and parallelizable, initial predictions may suffer from inconsistencies or errors, particularly when conditioned on complex structures (e.g., protein 3D conformation). To address this, **iterative refinement** is introduced as a mechanism to improve generation quality by progressively revising uncertain predictions across multiple steps.

Let a discrete sequence of length  $L$  be denoted by  $\mathbf{x} = [x_1, x_2, \dots, x_L]$ , and let  $\mathbf{x}^{(0)}$  be the initial input where some tokens are masked. The goal is to iteratively update  $\mathbf{x}^{(t)}$  to produce a coherent sequence  $\hat{\mathbf{x}}^{(T)}$  after  $T$  refinement steps.

At each iteration  $t = 1, \dots, T$ , the refinement process involves three stages:

1. **Prediction:** The masked model  $f_\theta$  predicts the token distribution given the current sequence:

$$\hat{\mathbf{x}}^{(t)} = f_\theta(\mathbf{x}^{(t-1)}, \text{cond})$$

where `cond` denotes optional conditioning information (e.g., structural constraints).

- 
- 270 2. **Remasking:** Based on prediction confidence (e.g., entropy or max softmax probability), a  
271 subset  $\mathcal{M}^{(t)}$  of uncertain positions is selected to be masked again:

$$272 \mathcal{M}^{(t)} = \text{TopKUncertain}(\hat{\mathbf{x}}^{(t)})$$

- 273  
274 3. **Update:** The input for the next step is updated as:

$$275 x_i^{(t)} = \begin{cases} \hat{x}_i^{(t)}, & \text{if } i \notin \mathcal{M}^{(t)} \\ [\text{MASK}], & \text{if } i \in \mathcal{M}^{(t)} \end{cases}$$

276  
277 This process is repeated for  $T$  iterations to progressively refine the prediction toward a coherent and  
278 high-quality output.

279 Following Zheng et al. (2023), we implement the iterative refinement strategy for sequences genera-  
280 tion. Compared with the time-aware Mask Generator proposed, we find our methods can achieve  
281 better performance. Because of the awareness of the hierarchical information, Hier-IF can design the  
282 protein sequence with higher recovery and perplexity.

## 283 REFERENCES

- 284 Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles,  
285 Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning-based  
286 protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- 287 Zhangyang Gao, Cheng Tan, Pablo Chacón, and Stan Z Li. Pifold: Toward effective and efficient  
288 protein inverse folding. *arXiv preprint arXiv:2209.12643*, 2022a.
- 289 Zhangyang Gao, Cheng Tan, and Stan Z Li. Alphadesign: A graph protein design method and  
290 benchmark on alphafolddb. *arXiv preprint arXiv:2202.01079*, 2022b.
- 291 Zhangyang Gao, Cheng Tan, and Stan Z Li. Knowledge-design: Pushing the limit of protein design  
292 via knowledge refinement. *arXiv preprint arXiv:2305.15151*, 2023.
- 293 Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander  
294 Rives. Learning inverse folding from millions of predicted structures. In *International conference  
295 on machine learning*, pp. 8946–8970. PMLR, 2022.
- 296 John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph-  
297 based protein design. *Advances in neural information processing systems*, 32, 2019.
- 298 Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael JL Townshend, and Ron Dror. Learning  
299 from protein structure with geometric vector perceptrons. *arXiv preprint arXiv:2009.01411*, 2020.
- 300 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,  
301 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level  
302 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- 303 Christine A Orengo, Alex D Michie, Susan Jones, David T Jones, Mark B Swindells, and Janet M  
304 Thornton. Cath—a hierarchic classification of protein domain structures. *Structure*, 5(8):1093–1109,  
305 1997.
- 306 Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei Ye, and Quanquan Gu. Structure-informed  
307 language models are protein designers. In *International conference on machine learning*, pp.  
308 42317–42338. PMLR, 2023.