

---

# Graph Denoising Diffusion for Inverse Protein Folding

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Inverse protein folding is challenging due to its inherent one-to-many mapping  
2 characteristic, where numerous possible amino acid sequences can fold into a single,  
3 identical protein backbone. This task involves not only identifying viable sequences  
4 but also representing the sheer diversity of potential solutions. However, existing  
5 discriminative models, such as transformer-based auto-regressive models, struggle  
6 to encapsulate the diverse range of plausible solutions. In contrast, diffusion  
7 probabilistic models, as an emerging genre of generative approaches, offer the  
8 potential to generate a diverse set of sequence candidates for determined protein  
9 backbones. We propose a novel graph denoising diffusion model for inverse  
10 protein folding, where a given protein backbone guides the diffusion process on  
11 the corresponding amino acid residue types. The model infers the joint distribution  
12 of amino acids conditioned on the nodes' physiochemical properties and local  
13 environment. Moreover, we utilize amino acid replacement matrices for the  
14 diffusion forward process, encoding the biologically-meaningful prior knowledge  
15 of amino acids from their spatial and sequential neighbors as well as themselves,  
16 which reduces the sampling space of the generative process. Our model achieves  
17 state-of-the-art performance over a set of popular baseline methods in sequence  
18 recovery and exhibits great potential in generating diverse protein sequences for a  
19 determined protein backbone structure.

## 20 1 Introduction

21 Inverse protein folding, or inverse folding, aims to predict feasible amino acid (AA) sequences that  
22 can fold into a specified 3D protein structure [21]. The results from inverse folding can facilitate  
23 the design of novel proteins with desired structural and functional characteristics. These proteins  
24 can serve numerous applications, ranging from targeted drug delivery to enzyme design for both  
25 academic and industrial purposes [24, 30, 37]. In this paper, we develop a diffusion model tailored  
26 for graph node denoising to obtain new AA sequences given a protein backbone.

27 Despite its importance, inverse folding remains challenging due to the immense sequence space to  
28 explore, coupled with the complexity of protein folding. On top of energy-based physical reasoning  
29 of a protein's folded state [1], recent advancements in deep learning yield significant progress in  
30 learning the mapping from protein structures to AA sequences directly. For example, discriminative  
31 models formulate this problem as the prediction of the most likely sequence for a given structure via  
32 Transformer-based models [6, 16, 22, 32]. However, they have struggled to accurately capture the  
33 one-to-many mapping from the protein structure to non-unique AA sequences.

34 Due to their powerful learning ability, *diffusion probabilistic models* have gained increasing attention.  
35 They are capable of generating a diverse range of molecule outputs from a fixed set of conditions  
36 given the inherent stochastic nature. For example, Torsion Diffusion [18] learns the distribution  
37 of torsion angles of heavy atoms to simulate conformations for small molecules. Concurrently,  
38 SMCDIFF [42] enhances protein folding tasks by learning the stable scaffold distribution supporting

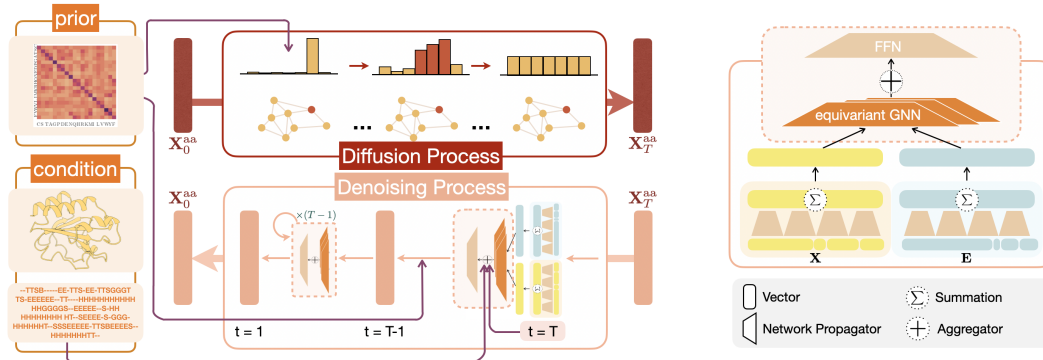


Figure 1: Overview of GRADE-IF. In the diffusion process, the original amino acid is stochastically transitioned to other amino acids, leveraging BLOSUM with varied temperatures as the transition kernel. During the denoising generation phase, initial node features are randomly sampled across the 20 amino acids with a uniform distribution. This is followed by a gradual denoising process, conditional on the graph structure and protein secondary structure at different time points. We employ a roto-translation equivariant graph neural network as the denoising network.

39 a target motif with diffusion. Similarly, DIFFDOCK [4] adopts a generative approach to protein-ligand  
 40 docking, creating a range of possible ligand binding poses for a target pocket structure.

41 Despite the widespread use of diffusion models, their comprehensive potential within the context  
 42 of protein inverse folding remains relatively unexplored. Current methods in sequence design are  
 43 primarily anchored in language models, encompassing *Masked Language Models* (MLMs) [24, 22]  
 44 and *autoregressive generative models* [16, 25, 27]. By tokenizing AAs, MLMs formulate the  
 45 sequence generation tasks as masked token enrichment. These models usually operate by drawing  
 46 an initial sequence with a certain number of tokens masked as a specific schedule and then learning  
 47 to predict the masked tokens from the given context. Intriguingly, this procedure can be viewed as  
 48 a discrete diffusion-absorbing model when trained by a parameterized objective. Autoregressive  
 49 models, conversely, can be perceived as deterministic diffusion processes. It induces conditional  
 50 distribution to each token, but the overall dependency along the entire AA sequence is recast via an  
 51 independently-executed diffusion process.

52 On the contrary, diffusion probabilistic models employ an iterative prediction methodology that  
 53 generates less noisy samples and demonstrates potential in capturing the diversity inherent in real  
 54 data distributions. This unique characteristic further underscores the promising role diffusion models  
 55 could play in advancing the field of protein sequence design. To bridge the gap, we make the first  
 56 attempt at a diffusion model for inverse folding. We model the inverse problem as a denoising  
 57 problem where the randomly assigned AA types in a protein (backbone) graph is recovered to the  
 58 wild type. The protein graph which contains the spatial and biochemical information of all AAs  
 59 is represented by equivariant graph neural networks, and diffusion process takes places on graph  
 60 nodes. In real inverse folding tasks, the proposed model achieves SOTA recovery rate, improve 4.2%  
 61 and 5.4% on recovery rate for single-chain proteins and short sequences, respectively, , especially  
 62 for conserved region which has a biological significance. Moreover, the predicted structure of  
 63 generated sequence is identical to the structure of native sequence.

64 The preservation of the desired functionalities is achieved by innovatively conditioning the model on  
 65 both secondary and third structures in the form of residue graphs and corresponding node features.  
 66 The major contributions of this paper are three-fold. Firstly, we propose GRADE-IF, a diffusion  
 67 model backed by roto-translation equivariant graph neural network for inverse folding. It stands out  
 68 from its counterparts for its ability to produce a wide array of diverse sequence candidates. Secondly,  
 69 as a departure from conventional uniform noise in discrete diffusion models, we encode the prior  
 70 knowledge of the response of AAs to evolutionary pressures by the utilization of *Blocks Substitution*  
 71 *Matrix* as the translation kernel. Moreover, to accelerate the sampling process, we adopt Denoising  
 72 Diffusion Implicit Model (DDIM) from its original continuous form to suit the discrete circumstances  
 73 and back it with thorough theoretical analysis.

74 **2 Problem Formulation**

75 **2.1 Residue Graph by Protein Backbone**

76 A residue graph, denoted as  $\mathcal{G} = (\mathbf{X}, \mathbf{A}, \mathbf{E})$ , aims to delineate the geometric configuration of a  
 77 protein. Specifically, every node stands for an AA within the protein. Correspondingly, each node  
 78 is assigned a collection of meticulously curated node attributes  $\mathbf{X}$  to reflect its physiochemical and  
 79 topological attributes. The local environment of a given node is defined by its spatial neighbors, as  
 80 determined by the  $k$ -nearest neighbor ( $k$ NN) algorithm. Consequently, each AA node is linked to  
 81 a maximum of  $k$  other nodes within the graph, specifically those with the least Euclidean distance  
 82 amongst all nodes within a  $30\text{\AA}$  contact region. The edge attributes, represented as  $\mathbf{E} \in \mathbb{R}^{93}$ , illustrate  
 83 the relationships between connected nodes. These relationships are determined through parameters  
 84 such as inter-atomic distances, local N-C positions, and a sequential position encoding scheme. We  
 85 detail the attribute construction in Appendix C.

86 **2.2 Inverse Folding as a Denoising Problem**

87 The objective of inverse folding is to engineer sequences that can fold to a pre-specified desired  
 88 structure. we utilize the coordinates of  $C\alpha$  atoms to represent the 3D positions of AAs in  
 89 Euclidean space, thereby embodying the protein backbone. Based on the naturally existing  
 90 protein structures, our model is constructed to generate a protein’s native sequence based on the  
 91 coordinates of its backbone atoms. Formally we represent this problem as learning the conditional  
 92 distribution  $p(\mathbf{X}^{\text{aa}}|\mathbf{X}^{\text{pos}})$ . Given a protein of length  $n$  and a sequence of spatial coordinates  
 93  $\mathbf{X} = \{x_1^{\text{pos}}, \dots, x_i^{\text{pos}}, \dots, x_n^{\text{pos}}\}$  representing each of the backbone  $C\alpha$  atoms in the structure, the  
 94 target is to predict  $\mathbf{X}^{\text{aa}} = \{x_1^{\text{aa}}, \dots, x_i^{\text{aa}}, \dots, x_n^{\text{aa}}\}$ , the native sequence of AAs. This density is  
 95 modeled in conjunction with the other AAs along the entire chain. Our model is trained by minimizing  
 96 the negative log-likelihood of the generated AA sequence relative to the native wild-type sequence.  
 97 Sequences can then be designed either by sampling or by identifying sequences that maximize the  
 98 conditional probability given the desired secondary and tertiary structure.

99 **2.3 Discrete Denoising Diffusion Probabilistic Models**

100 Diffusion models belong to the class of generative models, where the training stage encompasses  
 101 diffusion and denoising processes. The diffusion process  $q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$   
 102 corrupts the original data  $\mathbf{x}_0 \sim q(\mathbf{x})$  into a series of latent variables  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ , with each  
 103 carrying progressively higher levels of noise. Inversely, the denoising process  $p_\theta(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_T) =$   
 104  $p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  gradually reduces the noise within these latent variables, steering them  
 105 back towards the original data distribution. The iterative denoising procedure is driven by a  
 106 differentiable operator, such as a trainable neural network.

107 While in theory there is no strict form for  $q(\mathbf{x}_t | \mathbf{x}_{t-1})$  to take, several conditions are required  
 108 to be fulfilled by  $p_\theta$  for efficient sampling: (i) The diffusion kernel  $q(\mathbf{x}_t | \mathbf{x}_0)$  requires a closed  
 109 form to sample noisy data at different time steps for parallel training. (ii) The kernel should  
 110 possess a tractable formulation for the posterior  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$ . Consequently, the posterior  
 111  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \int q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) dp_\theta(\mathbf{x}_0 | \mathbf{x}_t)$ , and  $\mathbf{x}_0$  can be used as the target of the trainable  
 112 neural network. (iii) The marginal distribution  $q(\mathbf{x}_T)$  should be independent of  $\mathbf{x}_0$ . This independence  
 113 allows us to employ  $q(\mathbf{x}_T)$  as a prior distribution for inference.

114 The aforementioned criteria are crucial for the development of suitable noise-adding modules and  
 115 training pipelines. To satisfy these prerequisites, we follow the setting in previous work [2]. For  
 116 categorical data  $\mathbf{x}_t \in \{1, \dots, K\}$ , the transition probabilities are calculated by the matrix  $[Q_t]_{ij} =$   
 117  $q(\mathbf{x}_t = j | \mathbf{x}_{t-1} = i)$ . Employing the transition matrix and on one-hot encoded categorical feature  
 118  $\mathbf{x}_t$ , we can define the transitional kernel in the diffusion process by:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathbf{x}_{t-1} Q_t \quad \text{and} \quad q(\mathbf{x}_t | \mathbf{x}) = \mathbf{x} \bar{Q}_t, \quad (1)$$

119 where  $\bar{Q}_t = Q_1 \dots Q_t$ . The Bayes rule yields that the posterior distribution can be calculated  
 120 in closed form as  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}) \propto \mathbf{x}_t Q_t^\top \odot \mathbf{x} Q_{t-1}$ . The generative probability can thus be  
 121 determined using the transition kernel, the model output at time  $t$ , and the state of the process  $\mathbf{x}_t$ .  
 122 Through iterative sampling, we eventually produce the generated output  $\mathbf{x}_0$ .

123 The prior distribution  $p(\mathbf{x}_T)$  should be independent of the observation  $\mathbf{x}_0$ . Consequently, the  
 124 construction of the transition matrix necessitates the use of a noise schedule. The most straightforward  
 125 and commonly utilized method is the uniform transition, which can be parameterized as  $\mathbf{Q}_t =$   
 126  $\alpha_t \mathbf{I} + (1 - \alpha_t) \mathbf{I} \mathbf{I}^\top / d$  with  $\mathbf{I}^\top$  be the transpose of the identity matrix  $\mathbf{I}$ . As  $t$  approaches infinity,  $\alpha$   
 127 undergoes a progressive decay until it reaches 0. Consequently, the distribution  $q(\mathbf{x}_T)$  asymptotically  
 128 approaches a uniform distribution, which is essentially independent of  $\mathbf{x}$ .

### 129 3 Graph Denoising Diffusion for Inverse Protein Folding

130 In this section, we introduce a discrete graph denoising diffusion model for protein inverse folding,  
 131 which utilizes a given graph  $\mathcal{G} = \{\mathbf{X}, \mathbf{E}\}$  with node feature  $\mathbf{X}$  and edge feature  $\mathbf{E}$  as the condition.  
 132 Specifically, the node feature depicts the AA position, AA type, and the spatial and biochemical  
 133 properties  $\mathbf{X} = [\mathbf{X}^{\text{pos}}, \mathbf{X}^{\text{aa}}, \mathbf{X}^{\text{prop}}]$ . We define a diffusion process on the AA feature  $\mathbf{X}^{\text{aa}}$ , and  
 134 denoise it conditioned on the graph structure  $\mathbf{E}$  which is encoded by *equivariant neural networks* [35].  
 135 Moreover, we incorporate protein-specific prior knowledge, including an *AA substitution scoring*  
 136 *matrix* and protein *secondary structure* during modeling. We also introduce a new acceleration  
 137 algorithm for the discrete diffusion generative process based on a transition matrix.

#### 138 3.1 Diffusion Process and Generative Denoising Process

**Diffusion Process** To capture the distribution of AA types, we independently add noise to each AA  
 node of the protein. For any given node, the transition probabilities are defined by the matrix  $\mathbf{Q}_t$ .  
 With the predefined transition matrix, we can define the forward diffusion kernel by

$$q(\mathbf{X}_t^{\text{aa}} | \mathbf{X}_{t-1}^{\text{aa}}) = \mathbf{X}_{t-1}^{\text{aa}} \mathbf{Q}_t \quad \text{and} \quad q(\mathbf{X}_t^{\text{aa}} | \mathbf{X}^{\text{aa}}) = \mathbf{X}^{\text{aa}} \bar{\mathbf{Q}}_t,$$

139 where  $\bar{\mathbf{Q}}_t = \mathbf{Q}_1 \dots \mathbf{Q}_t$  is the transition probability matrix up to step  $t$ .

**Training Denoising Networks** The second component of the diffusion model is the denoising  
 140 neural network  $f_\theta$ , parameterized by  $\theta$ . This network accepts a noisy input  $\mathcal{G}_t = (\mathbf{X}_t, \mathbf{E})$ , where  $\mathbf{X}_t$   
 141 is the concatenation of the noisy AA types and other AA properties including 20 one-hot encoded AA  
 142 type and 15 geometry properties, such as SASA, ormlized surface-aware node features, dihedral  
 143 angles of backbone atoms, and 3D positions. It aims to predict the clean type of AA  $\mathbf{X}^{\text{aa}}$ , which  
 144 allows us to model the underlying sequence diversity in the protein structure while maintaining  
 145 their inherent structural constraints. To train  $f_\theta$ , we optimize the cross-entropy loss  $L$  between the  
 146 predicted probabilities  $\hat{p}(\mathbf{X}^{\text{aa}})$  for each node’s AA type.

**Parameterized Generative Process** A new AA sequence is generated through the reverse diffusion  
 148 iterations on each node  $\mathbf{x}$ . The generative probability distribution  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  is estimated from the  
 149 predicted probability  $\hat{p}(\mathbf{x}^{\text{aa}} | \mathbf{x}_t)$  by the neural networks. We marginalize over the network predictions  
 150 to compute for generative distribution at each iteration:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) \propto \sum_{\hat{\mathbf{x}}^{\text{aa}}} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) \hat{p}_\theta(\mathbf{x}^{\text{aa}} | \mathbf{x}_t), \quad (2)$$

152 where the posterior

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \text{Cat} \left( \mathbf{x}_{t-1} \left| \frac{\mathbf{x}_t \mathbf{Q}_t^\top \odot \mathbf{x}^{\text{aa}} \bar{\mathbf{Q}}_{t-1}}{\mathbf{x}^{\text{aa}} \bar{\mathbf{Q}}_t \mathbf{x}_t^\top} \right. \right) \quad (3)$$

153 can be calculated from the transition matrix, state of node feature at step  $t$  and AA type  $\mathbf{x}^{\text{aa}}$ . The  $\mathbf{x}^{\text{aa}}$   
 154 is the sample of the denoising network prediction  $\hat{p}(\mathbf{x}^{\text{aa}})$ .

#### 155 3.2 Prior Distribution from Protein Observations

##### 156 3.2.1 Markov Transition Matrices

157 The transition matrix serves as a guide for a discrete diffusion model, facilitating transitions between  
 158 the states by providing the probability of moving from the current time step to the next. As it reflects  
 159 the possibility from one AA type to another, this matrix plays a critical role in both the diffusion

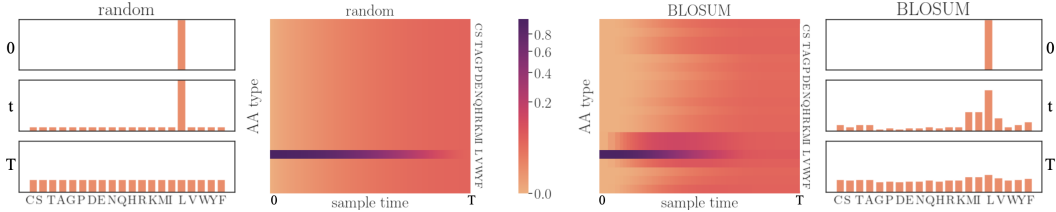


Figure 2: The middle two panels depict the transition probability of Leucine (L) from  $t = 0$  to  $T$ . Both the uniform and BLOSUM start as Dirichlet distributions and become uniform at time  $T$ . As shown in the two side figures, while the uniform matrix evenly disperses L’s probability to other AAs over time, BLOSUM favors AAs similar to L.

160 and generative processes. During the diffusion stage, the transition matrix is iteratively applied  
 161 to the observed data, which evolves over time due to inherent noise. As diffusion time increases,  
 162 the probability of the original AA type gradually decays, eventually converging towards a uniform  
 163 distribution across all AA types. In the generative stage, the conditional probability  $p(\mathbf{x}_{t-1}|\mathbf{x}_t)$   
 164 is determined by both the model’s prediction and the characteristics of the transition matrix  $\mathbf{Q}$ , as  
 165 described in Equation 2.

166 Given the biological specificity of AA substitutions, the transition probabilities between AAs are  
 167 not uniformly distributed, making it illogical to define random directions for the generative or  
 168 sampling process. As an alternative, the diffusion process could reflect evolutionary pressures by  
 169 utilizing substitution scoring matrices that conserve protein functionality, structure, or stability in  
 170 wild-type protein families. Formally, an *AA substitution scoring matrix* quantifies the rates at which  
 171 various AAs in proteins are substituted by other AAs over time [43]. In this study, we employ the  
 172 Blocks Substitution Matrix (BLOSUM) [11], which identifies conserved regions within proteins  
 173 that are presumed to have greater functional relevance. Grounded in empirical observations of  
 174 protein evolution, BLOSUM provides an estimate of the likelihood of substitutions between different  
 175 AAs. We thus incorporate BLOSUM into both the diffusion and generative processes. Initially, the  
 176 matrix is normalized into probabilities using the softmax function. Then, we use the normalized  
 177 matrix  $\mathbf{B}$  with different probability temperatures to control the noise scale of the diffusion process.  
 178 Consequently, the transition matrix at time  $t$  is given by  $\mathbf{Q}_t = \mathbf{B}^T$ . By using this matrix to refine the  
 179 transition probabilities, the generative space to be sampled is reduced effectively, thereby the model’s  
 180 predictions converge toward a meaningful subspace. See Figure 2 for a comparison of the transition  
 181 matrix over time in random and BLOSUM cases.

### 182 3.2.2 Secondary Structure

183 Protein secondary structure refers to the local spatial arrangement of AA residues in a protein chain.  
 184 The two most common types of protein secondary structure are alpha helices and beta sheets, which  
 185 are stabilized by hydrogen bonds between backbone atoms. The secondary structure of a protein  
 186 serves as a critical intermediary, bridging the gap between the AA sequence and the overall 3D  
 187 conformation of the protein. In our study, we incorporate eight distinct types of secondary structures  
 188 into AA nodes as conditions during the sampling process. This strategic approach effectively narrows  
 189 down the exploration space of potential AA sequences. Specifically, we employ DSSP (Define  
 190 Secondary Structure of Proteins) to predict the secondary structures of each AA and represent these  
 191 structures using one-hot encoding. Our neural network takes the one-hot encoding as input and  
 192 utilizes it to denoise the AA conditioned on it.

193 The imposition of motif conditions such as alpha helices and beta sheets on the search for AA  
 194 sequences not only leads to a significant reduction in the sampling space of potential sequences, but  
 195 also imparts biological implications for the generated protein sequence. By conditioning the sampling  
 196 process of AA types on their corresponding secondary structure types, we guide the resulting protein  
 197 sequence towards acquiring not only the appropriate 3D structure with feasible thermal stability but  
 198 also the capability to perform its intended function.

### 199 3.3 Equivariant Graph Denoising Network

200 Bio-molecules such as proteins and chemical compounds are structured in the 3-dimensional space,  
 201 and it is vital for the model to predict the same binding complex no matter how the input proteins are



202 positioned and oriented to encode a robust and expressive hidden representation. This property can be  
 203 guaranteed by rotation equivariance of the neural networks. A typical such a network is equivariant  
 204 graph neural network [35]. We modify its SE(3)-equivariant neural layers to update representations  
 205 for both nodes and edges, which reserves SO(3) rotation equivariance and E(3) translation invariance.  
 206 At the  $l$ th layer, an Equivariant Graph Convolution (EGC) inputs a set of  $n$  hidden node embeddings  
 207  $\mathbf{H}^{(l)} = \{\mathbf{h}_1^{(l)}, \dots, \mathbf{h}_n^{(l)}\}$  describing AA type and geometry properties, edge embedding  $\mathbf{m}_{ij}^{(l)}$  with  
 208 respect to connected nodes  $i$  and  $j$ , and  $\mathbf{X}^{\text{pos}} = \{\mathbf{x}_1^{\text{pos}}, \dots, \mathbf{x}_n^{\text{pos}}\}$  for node coordinates. The target  
 209 of a modified EGC layer is to update hidden representations  $\mathbf{H}^{(l+1)}$  for nodes and  $\mathbf{M}^{(l+1)}$  for edges.  
 210 Concisely,  $\mathbf{H}^{(l+1)}, \mathbf{M}^{(l+1)} = \text{EGC}[\mathbf{H}^{(l)}, \mathbf{X}^{\text{pos}}, \mathbf{M}^{(l)}]$ . To achieve this, an EGC layer defines

$$\begin{aligned}
 \mathbf{m}_{ij}^{(l+1)} &= \phi_e \left( \mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)}, \|\mathbf{x}_i^{(l)} - \mathbf{x}_j^{(l)}\|^2, \mathbf{m}_{ij}^{(l)} \right) \\
 \mathbf{x}_i^{(l+1)} &= \mathbf{x}_i^{(l)} + \frac{1}{n} \sum_{j \neq i} (\mathbf{x}_i^{(l)} - \mathbf{x}_j^{(l)}) \phi_x \left( \mathbf{m}_{ij}^{(l+1)} \right) \\
 \mathbf{h}_i^{(l+1)} &= \phi_h \left( \mathbf{h}_i^{(l)}, \sum_{j \neq i} \mathbf{m}_{ij}^{(l+1)} \right),
 \end{aligned} \tag{4}$$

211 where  $\phi_e, \phi_h$  are the edge and node propagation operations, respectively. The  $\phi_x$  is an additional  
 212 operation that projects the vector edge embedding  $\mathbf{m}_{ij}$  to a scalar. The modified EGC layer preserves  
 213 equivariance to rotations and translations on the set of 3D node coordinates  $\mathbf{X}^{\text{pos}}$  and performs  
 214 invariance to permutations on the nodes set identical to any other GNNs.

### 215 3.4 DDIM Sampling Process

216 A significant drawback of diffusion models lies in the speed of generation process, which is typically  
 217 characterized by numerous incremental steps and can be quite slow. Deterministic Denoising Implicit  
 218 Models (DDIM) [39] are frequently utilized to counter this issue in continuous variable diffusion  
 219 generative models. DDIM operates on a non-Markovian forward diffusion process, consistently  
 220 conditioning on the input rather than the previous step. By setting the noise variance on each step to  
 221 0, the reverse generative process becomes entirely deterministic, given an initial prior sample.

222 Similarly, since we possess the closed form of generative probability  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  in terms of a  
 223 predicted  $\mathbf{x}^{\text{aa}}$  and the posterior distribution  $p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}^{\text{aa}})$ , we can also render the generative model  
 224 deterministic by controlling the sampling temperature of  $p(\mathbf{x}^{\text{aa}}|\mathbf{x}_t)$ . Consequently, we can define the  
 225 multi-step generative process by

$$p_\theta(\mathbf{x}_{t-k} | \mathbf{x}_t) \propto \sum_{\hat{\mathbf{x}}^{\text{aa}}} q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) \hat{p}^T(\mathbf{x}^{\text{aa}} | \mathbf{x}_t) \tag{5}$$

226 where the temperature  $T$  controls whether it is deterministic or stochastic, and the multi-step posterior  
 227 distribution is

$$q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \text{Cat} \left( \mathbf{x}_{t-k} \mid \frac{\mathbf{x}_t Q_t^\top \cdots Q_{t-k}^\top \odot \mathbf{x}^{\text{aa}} \bar{Q}_{t-k}}{\mathbf{x}^{\text{aa}} \bar{Q}_t \mathbf{x}_t^\top} \right). \tag{6}$$

## 228 4 Experiments

229 We validate our GRADE-IF on recovering native protein sequences in CATH [29]. The performance is  
 230 mainly compared with structure-aware SOTA models. The implementations at [https://anonymous.4open.science/r/GraDe\\_IF-9574/](https://anonymous.4open.science/r/GraDe_IF-9574/)  
 231 are programmed with PyTorch-Geometric (ver 2.2.0) and PyTorch (ver 1.12.1) and executed on an NVIDIA<sup>®</sup> Tesla V100 GPU with 5, 120 CUDA cores and  
 232 32GB HBM2 installed on an HPC cluster.  
 233

### 234 4.1 Experimental Protocol

235 **Training Setup** We employ CATH v4.3.0 based partitioning as conducted by GRAPHTRANS [17]  
 236 and GVP [19]. Proteins are categorized based on CATH topology classification, leading to a division

Table 1: Recovery rate performance of **CATH** on zero-shot models.

Model	Perplexity ↓			Recovery Rate % ↑			CATH version	
	Short	Single-chain	All	Short	Single-chain	All	4.2	4.3
STRUCTGNN [17]	8.29	8.74	6.40	29.44	28.26	35.91	✓	
GRAPHTRANS [17]	8.39	8.83	6.63	28.14	28.46	35.82	✓	
GCA [41]	7.09	7.49	6.05	32.62	31.10	37.64	✓	
GVP [19]	7.23	7.84	5.36	30.60	28.95	39.47	✓	
GVP-large [16]	7.68	6.12	6.17	32.6	39.4	39.2		✓
ALPHADESIGN [8]	7.32	7.63	6.30	34.16	32.66	41.31	✓	
ESM-IF1 [16]	8.18	6.33	6.44	31.3	38.5	38.3		✓
ProteinMPNN [5]	6.21	6.68	4.61	36.35	34.43	45.96	✓	
PIFOLD [9]	6.04	6.31	4.55	39.84	38.53	51.66	✓	
<b>GRADE-IF</b>	<b>5.49</b>	<b>6.21</b>	<b>4.35</b>	<b>45.27</b>	<b>42.77</b>	<b>52.21</b>	✓	



Figure 3: Recovery rate on core and surface residues and different secondary structure

237 of 18, 024 proteins for training, 608 for validation, and 1, 120 for testing. To evaluate the generative  
 238 quality of different proteins, we test our model across three distinct categories: *short*, *single-chain*,  
 239 and *all* proteins. The short category includes proteins with sequence lengths shorter than 100. The  
 240 single-chain category encompasses proteins composed of a single chain. In addition, the total time  
 241 step of the diffusion model is configured as 500, adhering to a cosine schedule for noise [26]. For the  
 242 denoising network, we implement six stacked EGNN blocks, each possessing a hidden dimension of  
 243 128. Our model undergoes training for default of 200 epochs, making use of the Adam optimizer.  
 244 A batch size of 64 and a learning rate of 0.0005 are applied during training. Moreover, to prevent  
 245 overfitting, we incorporate a dropout rate of 0.1 into our model’s architecture.

246 **Evaluation Metric** The quality of recovered protein sequences is quantified by *perplexity* and  
 247 *recovery rate*. The former measures how well the model’s predicted AA probabilities match the  
 248 actual AA at each position in the sequence. A lower perplexity indicates a better fit of the model to  
 249 the data. The recovery rate assesses the model’s ability to recover the correct AA sequence given the  
 250 protein’s 3D structure. It is typically computed as the proportion of AAs in the predicted sequence  
 251 that matches the original sequence. A higher recovery rate indicates a better capability of the model  
 252 to predict the original sequence from the structure.

## 253 4.2 Inverse Folding

254 Table 1 compares GRADE-IF’s performance on recovering proteins in **CATH**, with the last column  
 255 indicating the training dataset of each baseline method. To generate high-confidence sequences,  
 256 GRADE-IF integrates out uncertainties in the prior by approximating the probability  $p(\mathbf{x}^{\text{aa}}) \approx$   
 257  $\sum_{i=1}^N p(\mathbf{x}^{\text{aa}} | \mathbf{x}_T^i) p(\mathbf{x}_T^i)$ . Notably, we observed an improvement of 4.2% and 5.4% in the recovery  
 258 rate for single-chain proteins and short sequences, respectively. We also conducted evaluations on  
 259 **TS50** and **TS500** datasets, with results included in Appendix E for further reference.

260 Upon subdividing the recovery performance based on buried and surface AAs, we find that the more  
 261 conserved core residues exhibit a higher native sequence recovery rate. In contrast, the active surface  
 262 AAs demonstrate a lower sequence recovery rate. Figure 7 examines AA conservation by Solvent

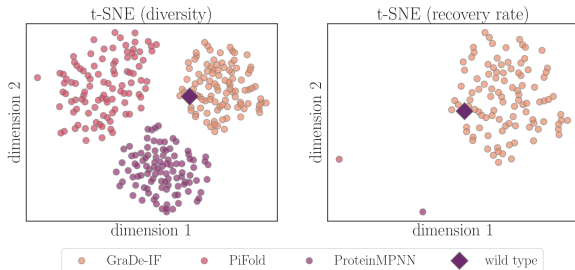


Figure 4: t-SNE of the generated sequences of GRADE-IF compared to PiFOLD and PROTEINMPNN.

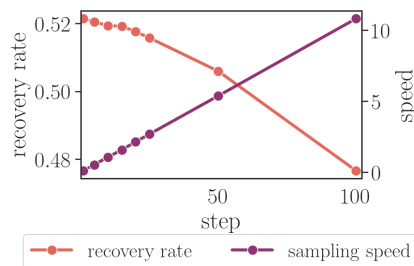


Figure 5: Trade-off of sampling speed and recovery rate.

263 Accessible Surface Area (SASA) (with  $SASA < 0.25$  indicating internal AAs) and contact number  
 264 (with the number of neighboring AAs within 8-Å in 3D space) [10]. The recovery rate of internal  
 265 residues significantly exceeds that of external residues across all three protein sequence classes, with  
 266 the recovery rate increasing in conjunction with the contact number. We also present the recovery  
 267 rate for different secondary structures, where we achieve high recovery for the majority of secondary  
 268 structures, with the exception of a minor 5-turn helix structure that occurs infrequently.

269 We further compare the diversity of GRADE-IF with PiFOLD and PROTEINMPNN in Figure 4.  
 270 For a given backbone, we generate 100 sequences with a self-similarity less than 50% and employ  
 271 t-SNE [44] for projection into a 2-dimensional space. At the same level of diversity, GRADE-  
 272 IF encompasses the wild-type sequence, whereas the other two methods fail to include the wild-type  
 273 within their sample region. Furthermore, inspiring at a recovery rate threshold of 45% for this protein,  
 274 GRADE-IF manages to generate a substantial number of samples, whereas the other two methods  
 275 revert to deterministic results. This further substantiates the superiority of our model in terms of  
 276 achieving sequence diversity and a high recovery rate concurrently.

277 We also evaluated the speed-up sampling algorithm within this dataset, as depicted in Figure 5. As  
 278 outlined in Equation equation 5, we can bypass  $k$  steps during the sampling phase. We selected a  
 279 range of step sizes and assessed their performance in terms of the recovery rate and the time required  
 280 to sample 1200 sequences. The recovery rate mildly declines with the increment in step size, reaching  
 281 48.13% at a step size of 100. However, the sampling speed at a step size of 100 is effectively 100  
 282 times faster than at a step size of 1, demonstrating a considerable speed-up.

### 283 4.3 Folding Prediction on Generated Sequences

284 We extend our investigation to the foldability of sequences generated at various sequence recovery  
 285 rates. Figure 6 contrasts the crystal structure of a native protein (PDB ID: 3FKF) with three structures  
 286 folded by ALPHAFOLD2 [20], each derived from a different GRADE-IF-generated sequence. The  
 287 resolution of the crystal structure stands at 2.2Å, suggesting that the folded structures of all generated  
 288 sequences are nearly identical to the native one, boasting an RMSD of approximately 1-Å over 139  
 289 residues. The average pLDDT score is 0.835, which, when compared to the native protein’s pLDDT of  
 290 0.91, underscores the reliability of their folded structures. In conjunction with the evidence presented  
 291 in Figure 7, indicating our method’s superior performance in generating more identical results within  
 292 conserved regions, we confidently posit that GRADE-IF can generate biologically plausible novel  
 293 sequences for given protein structures. We supplement more folding results in Appendix G.

## 294 5 Related Work

295 **Deep Learning models for protein sequence design** Self-supervised models have emerged as a  
 296 pivotal tool in the field of computational biology, providing a robust method for training extensive  
 297 protein sequences for representation learning. These models are typically divided into two categories:  
 298 structure-based generative models and sequence-based generative models. The former approaches  
 299 protein design by formulating the problem of fixed-backbone protein design as a conditional sequence  
 300 generation problem. They predict node labels, which represent AA types, with invariant or equivariant  
 301 graph neural networks [16, 17, 19, 40]. Alternatively, the latter sequence-based generative models



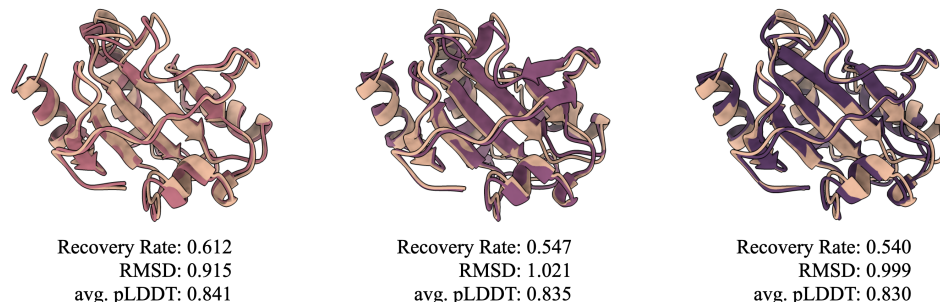


Figure 6: Folding prediction of generated protein sequence by GRADE-IF with respect to the native protein (PDB ID: 3FKF, colored in nude).

draw parallels between protein sequences and natural language processing. They employ attention-based methods to infer residue-wise relationships within the protein structure. These methods typically recover protein sequences autoregressively conditioned on the last inferred AA [24, 28, 36], or employing a BERT-style generative framework with masked language modeling objectives and enable the model to predict missing or masked parts of the protein sequence [22, 25, 33, 45].

**Denosing Diffusion models** The Diffusion Generative Model, initially introduced by Sohl-Dickstein *et al.* [38] and further developed by Ho *et al.* [12], has emerged as a potent instrument for a myriad of generative tasks in continuous time spaces. Its applications span diverse domains, from image synthesis [34] to audio generation [48], and it has also found utility in the creation of high-quality animations [13], the generation of realistic 3D objects [23], and drug design [4, 42]. Discrete adaptations of the diffusion model, on the other hand, have demonstrated efficacy in a variety of contexts, including but not limited to, text generation [2], image segmentation [14], and graph generation [15, 47]. Two distinct strategies have been proposed to establish a discrete variable diffusion process. The first approach involves the transformation of categorical data into a continuous space and then applying Gaussian diffusion [3, 15]. The alternative strategy is to define the diffusion process directly on the categorical data, an approach notably utilized in developing the D3PM model for text generation [2]. D3PM has been further extended to graph generation, facilitating the joint generation of node features and graph structure [46].

## 6 Conclusion

Deep learning approaches have striven to address a multitude of critical issues in bioengineering, such as protein folding, rigid-body docking, and property prediction. However, only a few methods have successfully generated diverse sequences for fixed backbones. In this study, we offered a viable solution by developing a denosing diffusion model to generate plausible protein sequences for a predetermined backbone structure. Our method, referred to as GRADE-IF, leverages substitution matrices for both diffusion and sampling processes, thereby exploring a practical search space for defining proteins. The iterative denosing process is predicated on the protein backbone revealing both the secondary and tertiary structure. The 3D geometry is analyzed by a modified equivariant graph neural network, which applies roto-translation equivariance to protein graphs without the necessity for intensive data augmentation. Given a protein backbone, our method successfully generated a diverse set of protein sequences, demonstrating a significant recovery rate. Importantly, these newly generated sequences are generally biologically meaningful, preserving more natural designs in the protein’s conserved regions and demonstrating a high likelihood of folding back into a structure highly similar to the native protein. The design of novel proteins with desired structural and functional characteristics is of paramount importance in the biotechnology and pharmaceutical industries, where such proteins can serve diverse purposes, ranging from targeted drug delivery to enzyme design for industrial applications. Additionally, understanding how varied sequences can yield identical structures propels the exploration of protein folding principles, thereby helping to decipher the rules that govern protein folding and misfolding. Furthermore, resolving the inverse folding problem allows the identification of different sequences that fold into the same structure, shedding light on the evolutionary history of proteins by enhancing our understanding of how proteins have evolved and diversified over time while preserving their functions.

343 **Checklist**

- 344 1. For all authors...
- 345 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
346 contributions and scope? Yes
- 347 (b) Did you describe the limitations of your work? Yes
- 348 (c) Did you discuss any potential negative societal impacts of your work? No, we believe  
349 our work has no negative societal impacts.
- 350 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
351 them? Yes
- 352 2. If you are including theoretical results...
- 353 (a) Did you state the full set of assumptions of all theoretical results? Yes
- 354 (b) Did you include complete proofs of all theoretical results? Yes
- 355 3. If you ran experiments...
- 356 (a) Did you include the code, data, and instructions needed to reproduce the main  
357 experimental results (either in the supplemental material or as a URL)? Yes
- 358 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
359 were chosen)? Yes
- 360 (c) Did you report error bars (e.g., with respect to the random seed after running  
361 experiments multiple times)? No
- 362 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
363 of GPUs, internal cluster, or cloud provider)? Yes
- 364 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 365 (a) If your work uses existing assets, did you cite the creators? Yes
- 366 (b) Did you mention the license of the assets? NA
- 367 (c) Did you include any new assets either in the supplemental material or as a URL? Yes
- 368 (d) Did you discuss whether and how consent was obtained from people whose data you're  
369 using/curating? NA
- 370 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
371 information or offensive content? NA
- 372 5. If you used crowdsourcing or conducted research with human subjects...
- 373 (a) Did you include the full text of instructions given to participants and screenshots, if  
374 applicable? NA
- 375 (b) Did you describe any potential participant risks, with links to Institutional Review  
376 Board (IRB) approvals, if applicable? NA
- 377 (c) Did you include the estimated hourly wage paid to participants and the total amount  
378 spent on participant compensation? NA

## 379 References

- 380 [1] Rebecca F Alford, Andrew Leaver-Fay, Jeliuzko R Jeliuzkov, Matthew J O’Meara, Frank P  
381 DiMaio, Hahnbeom Park, Maxim V Shapovalov, P Douglas Renfrew, Vikram K Mulligan, Kalli  
382 Kappel, et al. The Rosetta all-atom energy function for macromolecular modeling and design.  
383 *Journal of chemical theory and computation*, 13(6):3031–3048, 2017.
- 384 [2] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg.  
385 Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information*  
386 *Processing Systems*, 34:17981–17993, 2021.
- 387 [3] Ting Chen, Ruixiang ZHANG, and Geoffrey Hinton. Analog bits: Generating discrete data  
388 using diffusion models with self-conditioning. In *The Eleventh International Conference on*  
389 *Learning Representations*, 2023.
- 390 [4] Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, and Tommi S. Jaakkola. Diffdock:  
391 Diffusion steps, twists, and turns for molecular docking. In *The Eleventh International*  
392 *Conference on Learning Representations*, 2023.
- 393 [5] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F  
394 Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep  
395 learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- 396 [6] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Wang Yu, Llion Jones,  
397 Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and  
398 Burkhard Rost. ProtTrans: Towards cracking the language of lifes code through self-supervised  
399 deep learning and high performance computing. *IEEE Transactions on Pattern Analysis and*  
400 *Machine Intelligence*, 2021.
- 401 [7] Octavian-Eugen Ganea, Xinyuan Huang, Charlotte Bunne, Yatao Bian, Regina Barzilay,  
402 Tommi S Jaakkola, and Andreas Krause. Independent SE(3)-equivariant models for end-to-end  
403 rigid protein docking. In *International Conference on Learning Representations*, 2021.
- 404 [8] Zhangyang Gao, Cheng Tan, and Stan Z Li. Alphadesign: A graph protein design method and  
405 benchmark on alphafolddb. *arXiv preprint arXiv:2202.01079*, 2022.
- 406 [9] Zhangyang Gao, Cheng Tan, and Stan Z. Li. Pifold: Toward effective and efficient protein  
407 inverse folding. In *International Conference on Learning Representations*, 2023.
- 408 [10] Hai’e Gong, Haicang Zhang, Jianwei Zhu, Chao Wang, Shiwei Sun, Wei-Mou Zheng, and  
409 Dongbo Bu. Improving prediction of burial state of residues by exploiting correlation among  
410 residues. *BMC bioinformatics*, 18(3):165–175, 2017.
- 411 [11] Steven Henikoff and Jorja G Henikoff. Amino acid substitution matrices from protein blocks.  
412 *Proceedings of the National Academy of Sciences*, 89(22):10915–10919, 1992.
- 413 [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances*  
414 *in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 415 [13] Jonathan Ho, Tim Salimans, Alexey A. Gritsenko, William Chan, Mohammad Norouzi, and  
416 David J. Fleet. Video diffusion models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and  
417 Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 418 [14] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax  
419 flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural*  
420 *Information Processing Systems*, 34:12454–12465, 2021.
- 421 [15] Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant  
422 diffusion for molecule generation in 3d. In *International Conference on Machine Learning*,  
423 pages 8867–8887. PMLR, 2022.

- 424 [16] Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer,  
425 and Alexander Rives. Learning inverse folding from millions of predicted structures. In  
426 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato,  
427 editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
428 *Proceedings of Machine Learning Research*, pages 8946–8970. PMLR, 17–23 Jul 2022.
- 429 [17] John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for  
430 graph-based protein design. *Advances in Neural Information Processing Systems*, 32, 2019.
- 431 [18] Bowen Jing, Gabriele Corso, Jeffrey Chang, Regina Barzilay, and Tommi S. Jaakkola. Torsional  
432 diffusion for molecular conformer generation. In Alice H. Oh, Alekh Agarwal, Danielle  
433 Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*,  
434 2022.
- 435 [19] Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and Ron  
436 Dror. Learning from protein structure with geometric vector perceptrons. In *International  
437 Conference on Learning Representations*, 2021.
- 438 [20] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf  
439 Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al.  
440 Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021.
- 441 [21] George A Khoury, James Smadbeck, Chris A Kieslich, and Christodoulos A Floudas. Protein  
442 folding and de novo protein design for biotechnological applications. *Trends in biotechnology*,  
443 32(2):99–109, 2014.
- 444 [22] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,  
445 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level  
446 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- 447 [23] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In  
448 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
449 2837–2845, 2021.
- 450 [24] Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton,  
451 Jose Luis Olmos Jr, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language  
452 models generate functional protein sequences across diverse families. *Nature Biotechnology*,  
453 pages 1–8, 2023.
- 454 [25] Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language  
455 models enable zero-shot prediction of the effects of mutations on protein function. In *Advances  
456 in Neural Information Processing Systems*, volume 34, pages 29287–29303, 2021.
- 457 [26] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic  
458 models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- 459 [27] Erik Nijkamp, Jeffrey Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. ProGen2:  
460 exploring the boundaries of protein language models. *arXiv:2206.13517*, 2022.
- 461 [28] Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena Hurtado, Aidan N Gomez,  
462 Debora Marks, and Yarin Gal. Tranception: protein fitness prediction with autoregressive  
463 transformers and inference-time retrieval. In *International Conference on Machine Learning*,  
464 pages 16990–17017. PMLR, 2022.
- 465 [29] CA Orengo, AD Michie, S Jones, DT Jones, MB Swindells, and JM Thornton. CATH – a  
466 hierarchic classification of protein domain structures. *Structure*, 5(8):1093–1109, 1997.
- 467 [30] Robin Pearce and Yang Zhang. Deep learning techniques have significantly impacted protein  
468 structure prediction and protein design. *Current opinion in structural biology*, 68:194–207,  
469 2021.
- 470 [31] Yifei Qi and John ZH Zhang. Denscpd: improving the accuracy of neural-network-based  
471 computational protein sequence design with densenet. *Journal of chemical information and  
472 modeling*, 60(3):1245–1252, 2020.

- 473 [32] Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu,  
474 and Alexander Rives. MSA transformer. In *International Conference on Machine Learning*,  
475 pages 8844–8856. PMLR, 2021.
- 476 [33] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,  
477 Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from  
478 scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National  
479 Academy of Sciences*, 118(15):e2016239118, 2021.
- 480 [34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
481 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF  
482 Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- 483 [35] Victor Garcia Satorras, Emiel Hooeboom, and Max Welling. E(n) equivariant graph neural  
484 networks. In *International conference on machine learning*, pages 9323–9332, 2021.
- 485 [36] Jung-Eun Shin, Adam J Riesselman, Aaron W Kollasch, Conor McMahon, Elana Simon, Chris  
486 Sander, Aashish Manglik, Andrew C Kruse, and Debora S Marks. Protein design and variant  
487 prediction using autoregressive generative models. *Nature Communications*, 12(1):2403, 2021.
- 488 [37] Paweł Śledź and Amedeo Caffisch. Protein structure-based drug design: from docking to  
489 molecular dynamics. *Current opinion in structural biology*, 48:93–102, 2018.
- 490 [38] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep  
491 unsupervised learning using nonequilibrium thermodynamics. In *International Conference on  
492 Machine Learning*, pages 2256–2265. PMLR, 2015.
- 493 [39] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In  
494 *International Conference on Learning Representations*.
- 495 [40] Alexey Strokach, David Becerra, Carles Corbi-Verge, Albert Perez-Riba, and Philip M Kim.  
496 Fast and flexible protein design using deep graph neural networks. *Cell Systems*, 11(4):402–411,  
497 2020.
- 498 [41] Cheng Tan, Zhangyang Gao, Jun Xia, and Stan Z Li. Generative de novo protein design with  
499 global context. *arXiv preprint arXiv:2204.10673*, 2022.
- 500 [42] Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay,  
501 and Tommi S. Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the  
502 motif-scaffolding problem. In *International Conference on Learning Representations*, 2023.
- 503 [43] Rakesh Trivedi and Hampapathalu Adimurthy Nagarajaram. Substitution scoring matrices for  
504 proteins-an overview. *Protein Science*, 29(11):2150–2163, 2020.
- 505 [44] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine  
506 learning research*, 9(11), 2008.
- 507 [45] Jesse Vig, Ali Madani, Lav R Varshney, Caiming Xiong, Nazneen Rajani, et al. BERTology  
508 meets biology: Interpreting attention in protein language models. In *International Conference  
509 on Learning Representations*, 2021.
- 510 [46] Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal  
511 Frossard. Digress: Discrete denoising diffusion for graph generation. *arXiv preprint  
512 arXiv:2209.14734*, 2022.
- 513 [47] Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal  
514 Frossard. Digress: Discrete denoising diffusion for graph generation. In *The Eleventh  
515 International Conference on Learning Representations*, 2023.
- 516 [48] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu.  
517 Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Transactions on  
518 Audio, Speech, and Language Processing*, 2023.
- 519 [49] Bingxin Zhou, Outongyi Lv, Kai Yi, Xinye Xiong, Pan Tan, Liang Hong, and Yu Guang Wang.  
520 Accurate and definite mutational effect prediction with lightweight equivariant graph neural  
521 networks. *arXiv:2304.08299*, 2023.

## 522 A Broader Impact and Limitations

523 **Broader Impact** We have developed a generative model rooted in the diffusion denoising paradigm,  
 524 specifically tailored to the context of protein inverse folding. As with any other generative models, it  
 525 is capable of generating *de novo* content (protein sequences) under specified conditions (e.g., protein  
 526 tertiary structure). While this method holds substantial potential for facilitating scientific research and  
 527 biological discoveries, its misuse could pose potential risks to human society. For instance, in theory,  
 528 it possesses the capacity to generate novel viral protein sequences with enhanced functionalities. To  
 529 mitigate this potential risk, one approach could be to confine the training dataset for the model to  
 530 proteins derived from prokaryotes and/or eukaryotes, thereby excluding viral proteins. Although this  
 531 strategy may to some extent compromise the overall performance and generalizability of the trained  
 532 model, it also curtails the risk of misuse of the model by limiting the understanding and analysis of  
 533 viral protein construction.

534 **Limitations** The conditions imposed on the sampling process gently guide the generated protein  
 535 sequences. However, in certain scenarios, stringent restrictions may be necessary to produce a  
 536 functional protein. Secondary structure, as a living example, actively contributes to the protein’s  
 537 functionality. For instance, transmembrane  $\alpha$ -helices play essential roles in protein functions, such as  
 538 passing ions or other molecules and transmitting a signal across the membrane. Moreover, the current  
 539 zero-shot model is trained on a general protein database. For specific downstream applications,  
 540 such as generating new sequences for a particular protein or protein family, it may necessitate the  
 541 incorporation of auxiliary modules or the modification of training procedures to yield more fitting  
 542 sequences.

## 543 B Non-Markovian Forward Process

544 We give the derivation of posterior distribution  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}})$  for generative process from step to  
 545 step. The proof relies on the Bayes rule, Markov property, and the pre-defined transition matrix for  
 546 AAs.

547 **Proposition 1.** For  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}})$  defined in Eq 3, we have

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \text{Cat} \left( \mathbf{x}_{t-1} \middle| \frac{\mathbf{x}_t Q_t^\top \odot \mathbf{x}^{\text{aa}} \bar{Q}_{t-1}}{\mathbf{x}^{\text{aa}} \bar{Q}_t \mathbf{x}_t^\top} \right).$$

548 *Proof.* By Bayes rules, we can expand the original equation  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}})$  to

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}^{\text{aa}}) q(\mathbf{x}_{t-1} | \mathbf{x}^{\text{aa}})}{q(\mathbf{x}_t | \mathbf{x}^{\text{aa}})} = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-1}) q(\mathbf{x}_{t-1} | \mathbf{x}^{\text{aa}})}{q(\mathbf{x}_t | \mathbf{x}^{\text{aa}})}.$$

549 As pre-defined diffusion process, we get  $q(\mathbf{x}_t | \mathbf{x}^{\text{aa}}) = \mathbf{x}^{\text{aa}} \bar{Q}_t$ , and  $q(\mathbf{x}_{t-1} | \mathbf{x}^{\text{aa}}) = \mathbf{x}^{\text{aa}} \bar{Q}_{t-1}$ .

For the term of  $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}^{\text{aa}})$  by Bayes rule and Markov property, we have

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}^{\text{aa}}) = q(\mathbf{x}_t | \mathbf{x}_{t-1}) \propto q(\mathbf{x}_{t-1} | \mathbf{x}_t) \pi(\mathbf{x}_t) \propto \mathbf{x}_t Q_t^\top \odot \pi(\mathbf{x}_t)$$

550 where the normalizing constant is  $\sum_{\mathbf{x}_{t-1}} \mathbf{x}_t Q_t^\top \odot \pi(\mathbf{x}_t) = (\mathbf{x}_t \sum_{\mathbf{x}_{t-1}} Q_t^\top) \odot \pi(\mathbf{x}_t) = \mathbf{x}_t \odot \pi(\mathbf{x}_t)$

551 Then  $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}^{\text{aa}}) = \frac{\mathbf{x}_t Q_t^\top}{\mathbf{x}_t}$ , and the posterior distribution is:

$$q(\mathbf{x}_{t-1} | \mathbf{x}^{\text{aa}}, \mathbf{x}_t) = \text{Cat} \left( \mathbf{x}_{t-1} \middle| \frac{\mathbf{x}_t Q_t^\top \odot \mathbf{x}^{\text{aa}} \bar{Q}_{t-1}}{\mathbf{x}^{\text{aa}} \bar{Q}_t \mathbf{x}_t^\top} \right).$$

552 □

553 The following gives the derivation for the discrete DDIM which accelerates the generative process.

554 **Proposition 2.** For  $q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}})$  defined in Eq 6,

$$q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \text{Cat} \left( \mathbf{x}_{t-k} \middle| \frac{\mathbf{x}_t Q_t^\top \cdots Q_{t-k}^\top \odot \mathbf{x}^{\text{aa}} \bar{Q}_{t-k}}{\mathbf{x}^{\text{aa}} \bar{Q}_t \mathbf{x}_t^\top} \right).$$



555 *Proof.* By Bayes rules, we can expand the original equation  $q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}})$  to

$$q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-k}, \mathbf{x}^{\text{aa}}) q(\mathbf{x}_{t-k} | \mathbf{x}^{\text{aa}})}{q(\mathbf{x}_t | \mathbf{x}^{\text{aa}})} = \frac{q(\mathbf{x}_t | \mathbf{x}_{t-k}) q(\mathbf{x}_{t-k} | \mathbf{x}^{\text{aa}})}{q(\mathbf{x}_t | \mathbf{x}^{\text{aa}})}.$$

556 As pre-defined diffusion process, we get  $q(\mathbf{x}_t | \mathbf{x}^{\text{aa}}) = \mathbf{x}^{\text{aa}} \bar{Q}_t$ , and  $q(\mathbf{x}_{t-1} | \mathbf{x}^{\text{aa}}) = \mathbf{x}^{\text{aa}} \bar{Q}_{t-1}$ .

557 Similarly with  $q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}^{\text{aa}})$  in Proposition 1,  $q(\mathbf{x}_t | \mathbf{x}_{t-k}, \mathbf{x}^{\text{aa}}) = \frac{\mathbf{x}_t Q_t^\top \cdots Q_{t-k}^\top}{\mathbf{x}_t}$  and the  
558 posterior is

$$q(\mathbf{x}_{t-k} | \mathbf{x}_t, \mathbf{x}^{\text{aa}}) = \text{Cat} \left( \mathbf{x}_{t-k} \mid \frac{\mathbf{x}_t Q_t^\top \cdots Q_{t-k}^\top \odot \mathbf{x}^{\text{aa}} \bar{Q}_{t-k}}{\mathbf{x}^{\text{aa}} \bar{Q}_t \mathbf{x}_t^\top} \right).$$

559 □

## 560 C Graph Representation of Folded Proteins

561 The geometry of proteins suggests higher-level structures and topological relationships, which are  
562 vital to protein functionality. For a given protein, we create a  $k$ -nearest neighbor ( $k$ NN) graph  
563  $\mathcal{G} = (\mathbf{X}, \mathbf{E})$  to describe its physiochemical and geometric properties with nodes representing AAs  
564 by  $\mathbf{X} \in \mathbb{R}^{39}$  node attributes with 20-dim AA type encoder, 16-dim AA properties, and 3-dim AA  
565 positions. The undirected edge connections are formulated via a  $k$ NN-graph with cutoff. In other  
566 words, each node is connected to up to  $k$  other nodes in the graph that has the smallest Euclidean  
567 distance over other nodes and the distance is smaller than a certain cutoff (*e.g.*, 30Å). Edge attributes  
568 are defined for connected node pairs. For instance, if node  $i$  and  $j$  are connected to each other, their  
569 relationship will be described by  $\mathbf{E}_{ij} = \mathbf{E}_{ji} \in \mathbb{R}^{93}$ .

570 The AA types are one-hot encoded to 20 binary values by  $\mathbf{X}^{\text{aa}}$ . On top of it, the properties of AAs  
571 and AAs' local environment are described by  $\mathbf{X}^{\text{prop}}$ , including the normalized crystallographic  
572 B-factor, solvent-accessible surface area (SASA), normalized surface-aware node features, dihedral  
573 angles of backbone atoms, and 3D positions. SASA measures the level of exposure of an AA to  
574 solvent in a protein by a scalar value, which provides an important indicator of active sites of proteins  
575 to locate whether a residue is on the surface of the protein. Both B-factor and SASA are standardized  
576 with AA-wise mean and standard deviation on the associate attribute. Surface-aware features [7] of  
577 an AA is non-linear projections to the weighted average distance of the central AA to its one-hop  
578 neighbors  $i' \in \mathcal{N}_i$ , *i.e.*,

$$\rho(\mathbf{x}_i; \lambda) = \frac{\left\| \sum_{i' \in \mathcal{N}_i} w_{i,i',\lambda} (\mathbf{X}^{\text{pos},i} - \mathbf{X}^{\text{pos},i'}) \right\|}{\sum_{i' \in \mathcal{N}_i} w_{i,i',\lambda} \|\mathbf{X}^{\text{pos},i} - \mathbf{X}^{\text{pos},i'}\|},$$

579 where the weights are defined by

$$w_{i,i',\lambda} = \frac{\exp\left(-\|\mathbf{X}_{\text{pos},i} - \mathbf{X}_{\text{pos},i'}\|^2 / \lambda\right)}{\sum_{i' \in \mathcal{N}_i} \exp\left(-\|\mathbf{X}_{\text{pos},i} - \mathbf{X}_{\text{pos},i'}\|^2 / \lambda\right)}$$

580 with  $\lambda \in \{1, 2, 5, 10, 30\}$ . The  $\mathbf{X}^{\text{pos},i} \in \mathbb{R}^3$  denotes the 3D coordinates of the  $i$ th residue, which  
581 is represented by the position of  $\alpha$ -carbon. We also use the backbone atom positions to define the  
582 spatial conformation of each AA in the protein chain with trigonometric values of dihedral angles  
583  $\{\sin, \cos\} \circ \{\phi_i, \psi_i, \omega_i\}$ .

584 Edge attributes  $\mathbf{E} \in \mathbb{R}^{93}$ , on the other hand, include kernel-based distances, relative spatial positions,  
585 and relative sequential distances for pairwise distance characterization. For two connected residues  $i$   
586 and  $j$ , the kernel-based distance between them is projected by Gaussian radial basis functions (RBF)  
587 of  $\exp\left\{\frac{\|\mathbf{x}_j - \mathbf{x}_i\|^2}{2\sigma_r^2}\right\}$  with  $r = 1, 2, \dots, R$ . A total number of 15 distinct distance-based features  
588 are created with  $\sigma_r = \{1.5^k \mid k = 0, 1, 2, \dots, 14\}$ . Next, local frames [7] are created from the  
589 corresponding residues' heavy atoms positions to define 12 relative positions. They represent local  
590 fine-grained relations between AAs and the rigid property of how the two residues interact with each  
591 other. Finally, the residues' sequential relationship is encoded with 66 binary features by their relative  
592 position  $d_{i,j} = |s_i - s_j|$ , where  $s_i$  and  $s_j$  are the absolute positions of the two nodes in the AA chain  
593 [49]. We further define a binary contact signal [17] to indicate whether two residues contact in the  
594 space, *i.e.*, the Euclidean distance  $\|C\alpha_i - C\alpha_j\| < 8$ .

595 **D Training and Inference**

596 In this section, we elucidate the training and inference methodologies implemented in the diffusion  
 597 generative model. As shown in Algorithm 1, training commences with a random sampling of a time  
 598 scale  $t$  from a uniform distribution between 1 and  $T$ . Subsequently, we calculate the noise posterior  
 599 and integrate noise as dictated by its respective distribution. We then utilize an equivariant graph  
 600 neural network for denoising predictions, using both the noisy amino acid and other properties as  
 601 node features, and leveraging the graph structure for geometric information. This results in the model  
 602 outputting the denoised amino acid type. Ultimately, the cross-entropy loss is computed between the  
 603 predicted and original amino acid types, providing a parameter for optimizing the neural network.

---

**Algorithm 1** Training

---

- 1: **Input:** A graph  $\mathcal{G} = \{\mathbf{X}, \mathbf{E}\}$
  - 2: Sample  $t \sim \mathcal{U}(1, T)$
  - 3: Compute  $q(\mathbf{X}_t | \mathbf{X}^{\text{aa}}) = \mathbf{X}^{\text{aa}} \bar{Q}_t$
  - 4: Sample noisy  $\mathbf{X}_t \sim q(\mathbf{X}_t | \mathbf{X}^{\text{aa}})$
  - 5: Forward pass:  $\hat{p}(\mathbf{X}^{\text{aa}}) = f_\theta(\mathbf{X}_t, t, \mathbf{E}, ss)$
  - 6: Compute cross-entropy loss:  $L = L_{\text{CE}}(\hat{p}(\mathbf{X}^{\text{aa}}), \mathbf{X})$
  - 7: Compute the gradient and optimize denoise network  $f_\theta$
- 

604 Upon completing the training, we are capable of sampling data using the neural network and the  
 605 posterior distribution  $p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}^{\text{aa}})$ . As delineated in the algorithm, we initially sample an amino  
 606 acid uniformly from 20 classes, then employ our neural network to denoise  $\mathbf{X}^{\text{aa}}$  from time  $t$ . From  
 607 here, we can calculate the forward probability utilizing the model output and the posterior distribution.  
 608 Through iterative processing, the ultimate model sample closely approximates the original data  
 609 distribution. More importantly, we illustrate how to speed up the sampling procedure using DDIM in  
 610 Algorithm 3. It can be regarded as skipping several steps in DDPM but with close performance (see  
 611 Figure 5 in Section 4.2). DDPM is a special case of DDIM when skipping step  $k = 1$ .

---

**Algorithm 2** Sampling (DDPM)

---

- 1: Sample from uniformly prior  $\mathbf{X}_T \sim p(\mathbf{X}_T)$
  - 2: **for**  $t$  in  $\{T, T - 1, \dots, 1\}$  **do**
  - 3:     Predict  $\hat{p}(\mathbf{X}_0 | \mathbf{X}_t)$  by neural network  $\hat{p}(\mathbf{X}_0 | \mathbf{X}_t) = f_\theta(\mathbf{X}_t, t, \mathbf{E}, ss)$
  - 4:     Compute  $p_\theta(\mathbf{X}_{t-1} | \mathbf{X}_t) = \sum_{\hat{\mathbf{X}}^{\text{aa}}} q(\mathbf{X}_{t-1} | \mathbf{X}_t, \hat{\mathbf{X}}^{\text{aa}}) \hat{p}(\mathbf{X}^{\text{aa}} | \mathbf{X}_t)$
  - 5:     Sample  $\mathbf{X}_{t-1} \sim p_\theta(\mathbf{X}_{t-1} | \mathbf{X}_t)$
  - 6: **end for**
  - 7: Sample  $\mathbf{X}^{\text{aa}} \sim p_\theta(\mathbf{X}^{\text{aa}} | \mathbf{X}_1)$
- 

---

**Algorithm 3** Sampling (DDIM)

---

- 1: Sample from uniformly prior  $\mathbf{X}_T \sim p(\mathbf{X}_T)$
  - 2: **for**  $t$  in  $\{T, T - k, \dots, 1\}$  **do**
  - 3:     Predict  $\hat{p}(\mathbf{X}_0 | \mathbf{X}_t)$  by neural network  $\hat{p}(\mathbf{X}_0 | \mathbf{X}_t) = f_\theta(\mathbf{X}_t, t, \mathbf{E}, ss)$
  - 4:     Compute  $p_\theta(\mathbf{X}_{t-k} | \mathbf{X}_t) = \sum_{\hat{\mathbf{X}}^{\text{aa}}} q(\mathbf{X}_{t-k} | \mathbf{X}_t, \hat{\mathbf{X}}^{\text{aa}}) \hat{p}(\mathbf{X}^{\text{aa}} | \mathbf{X}_t)$
  - 5:     Sample  $\mathbf{X}_{t-k} \sim p_\theta(\mathbf{X}_{t-k} | \mathbf{X}_t)$
  - 6: **end for**
  - 7: Sample  $\mathbf{X}^{\text{aa}} \sim p_\theta(\mathbf{X}^{\text{aa}} | \mathbf{X}_1)$
- 

612 **E Inverse Folding Performance on TS50 and T500**

613 In addition to the CATH dataset, we also evaluated our model using the **TS50** and **T500** datasets.  
 614 These datasets were introduced by DenseCPD [31], encompassing 9888 structures for training, and  
 615 two distinct test datasets comprising 50 (TS50) and 500 (T500) test datasets, respectively. The

Table 2: Recovery rate performance of **TS50** and **T500** on zero-shot models.

Model	TS50		T500	
	Perplexity ↓	Recovery ↑	Perplexity ↓	Recovery ↑
STRUCTGNN [17]	5.40	43.89	4.98	45.69
GRAPHTRANS [17]	5.60	42.20	5.16	44.66
GVP [19]	4.71	44.14	4.20	49.14
GCA [41]	5.09	47.02	4.72	47.74
ALPHADESIGN [8]	5.25	48.36	4.93	49.23
PROTEINMPNN [5]	3.93	54.43	3.53	58.08
PIFOLD [9]	3.86	<b>58.72</b>	3.44	60.42
GRADE-IF(ours)	<b>3.71</b>	56.32	<b>3.23</b>	<b>61.22</b>

616 same preprocessing steps applied to the **CATH** dataset were utilized here. The denoising network  
 617 comprises six sequentially arranged EGNN blocks, each boasting a hidden dimension of 256. Our  
 618 model’s performance, outlined in Table 2, achieved an accuracy of 61.22% on **T500**, and 56.32% on  
 619 **TS50**, respectively.

## 620 F Ablation Study

621 We conducted ablation studies to assess the impact of various factors on our model’s performance.  
 622 These elements encompassed the selection of the transition matrix (uniform versus BLOSUM),  
 623 the integration of secondary structure embeddings in the denoising procedure, and the function of  
 624 the equivariant neural network. As demonstrated in Figure 7, incorporating equivariance into the  
 625 denoising neural network substantially enhances the model’s performance. Given that the placement  
 626 of protein structures in space can be arbitrary, considering symmetry in the denoising neural network  
 627 helps to mitigate disturbances. Moreover, we found that including secondary structure as auxiliary  
 628 information lessens uncertainty and improves recovery. Lastly, utilizing the BLOSUM matrix as  
 629 the noise transition matrix boosted the recovery rate by 2%, highlighting the benefits of infusing  
 630 biological information into the diffusion and generative processes. This approach reduces sample  
 631 variance and substantially benefits overall model performance.

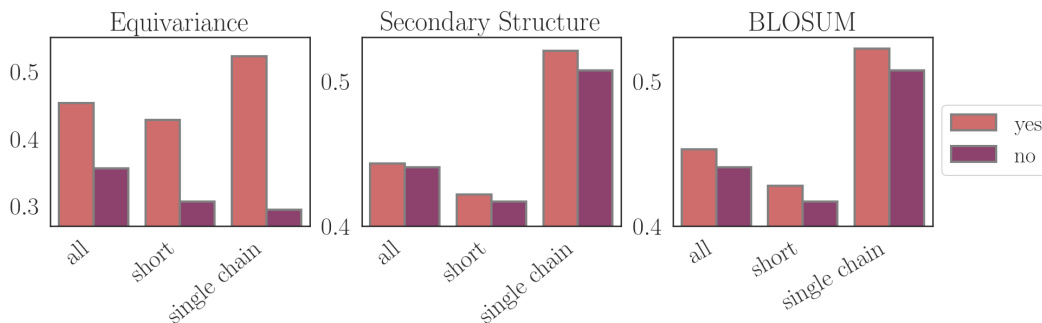


Figure 7: Recovery rate with the different selection of the transition matrix, whether considering equivariance and secondary structure.

632 In our sampling procedure, we accelerate the original DDPM sampling algorithm, which takes every  
 633 step in the reverse sampling process, by implementing the discrete DDIM as per Equation 6. This  
 634 discrete DDIM allows us to skip every  $k$  steps, resulting in a speed-up of the original DDPM by  
 635 a factor of  $k$ . We conducted an ablation study on the impact of speed and recovery rate by trying  
 636 different skip steps: 1, 2, 5, 10, 20, 25, 50, and 100. We compare the recovery rates achieved by these  
 637 different steps. Our results revealed that the recovery rate performance decays as the number of  
 638 skipped steps increases. The best performance is achieved when skipping a single step, resulting in a  
 639 recovery rate of 52.21%, but at a speed of 100 times slower than when skipping 100 steps, which  
 640 yields a recovery rate of 47.66%.

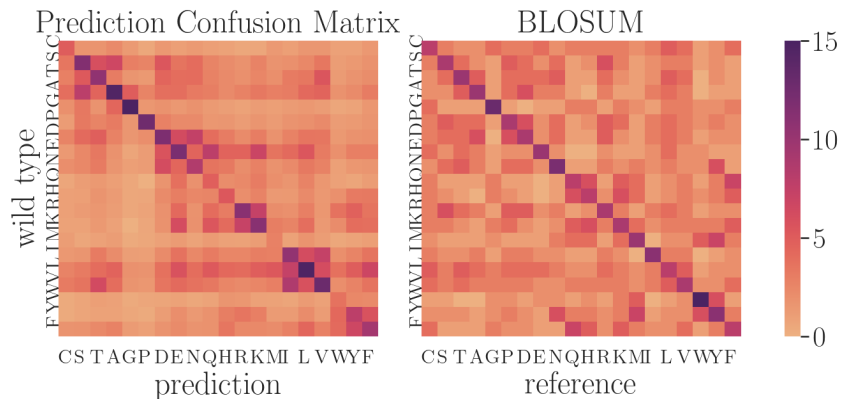


Figure 8: Comparison of the distribution of mutational prior (BLOSUM replacement matrix) and sampling results.

## 641 G Additional Folding Results

642 We further analyzed the generated sequences by comparing different protein folding predictions. We  
 643 consider the crystal structures of three native proteins with PDB IDs: **1ud9** (A chain), **2rem** (B chain),  
 644 **3drn** (B chain), which we randomly choose from **CATH** dataset. For each structure, we generated  
 645 three sequences from the diffusion model and used ALPHAFOLD 2 [20] to predict the respective  
 646 structures. As shown in Figure 9, these predictions (in purple) were then compared with the structures  
 647 of the native protein sequences (in nude). We can observe that the RMSD for all cases is lower than  
 648 the preparation accuracy of the wet experiment. The results demonstrate that our model-generated  
 649 sequences retain the core structure, indicating their fidelity to the original structures.

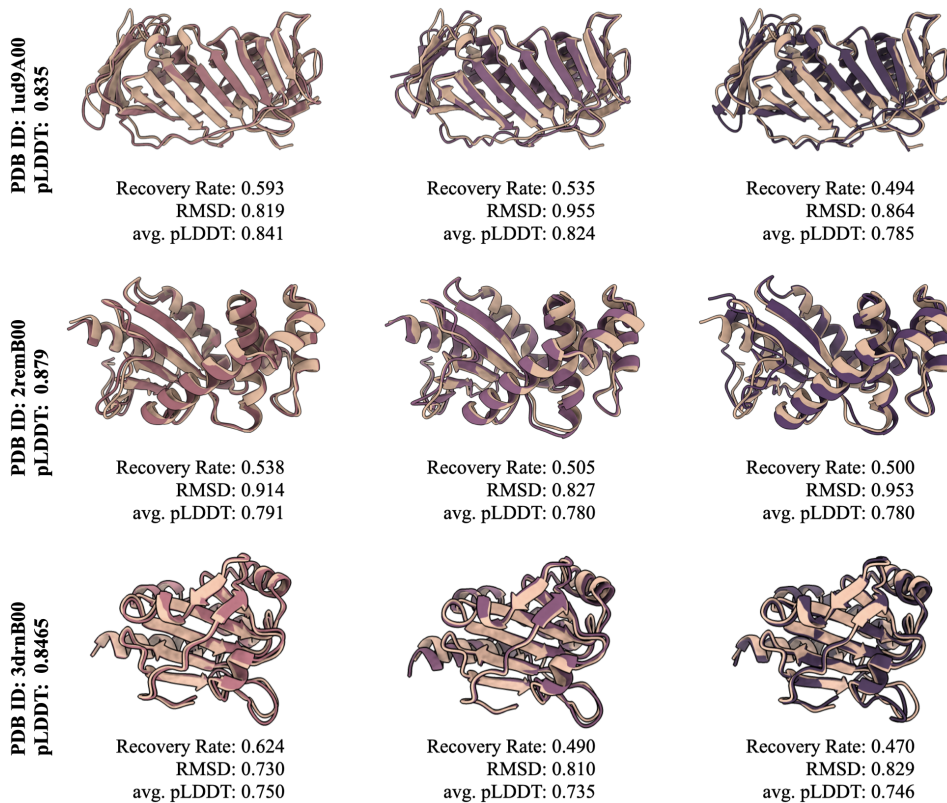


Figure 9: Folding comparison between native sequence and generated sequence