IMPROVING INVERSE FOLDING FOR PEPTIDE DE SIGN WITH DIVERSITY-REGULARIZED DIRECT PREF ERENCE OPTIMIZATION

Anonymous authors Paper under double-blind review

ABSTRACT

Inverse folding models play an important role in structure-based design by predicting amino acid sequences that fold into desired reference structures. Models like ProteinMPNN, a message-passing encoder-decoder model, are trained to reliably produce new sequences from a reference structure. However, when applied to peptides, these models are prone to generating repetitive sequences that do not fold into the reference structure. To address this, we fine-tune ProteinMPNN to produce diverse and structurally consistent peptide sequences via Direct Preference Optimization (DPO). We derive two enhancements to DPO: online diversity regularization and domain-specific priors. Additionally, we develop a new understanding on improving diversity in decoder models. When conditioned on Open-Fold generated structures, our fine-tuned models achieve state-of-the-art structural similarity scores, improving base ProteinMPNN by at least 8%. Compared to standard DPO, our regularized method achieves up to 20% higher sequence diversity with no loss in structural similarity score.

1 INTRODUCTION

Engineering biopolymers that fold into desired 3D structures, a computational challenge known as inverse protein folding problem, has broad applications in drug discovery and material science (Yang et al., 2023; Dill et al., 2008; Abascal & Regan, 2018). Several approaches for inverse folding have been adopted over the past decades, from molecular dynamics simulations to machine learning approaches (Dauparas et al., 2022b; Shanker et al., 2023; Hsu et al., 2022a; Yi et al., 2023; Correa, 1990). In the standard machine learning approach, a molecular backbone chain serves as input, and a model generates sequences that adopt folding topologies compatible with the reference backbone. Sequences do not necessarily share sequence homology, as multiple diverse sequences can fold into similar structures (Hsu et al., 2022a; Yue & Dill, 1992; Godzik et al., 1993).

Peptides, which are small biopolymers comprising 2-50 residues, are interesting targets for inverse 040 folding given their role in diverse biological functions, acting as hormones, neurotransmitters, sig-041 nalling molecules, or nanostructures assemblers (Chockalingam et al., 2007; Torres et al., 2019; 042 Copolovici et al., 2014; Ulijn & Smith, 2008). Only about 225,000 protein structures have been 043 experimentally determined¹ and made available via the Protein Data Bank (PDB). Training inverse-044 folding machine learning models in a supervised fashion is a challenging task, due to the complexity of the problem and the limited amount of experimental data. The challenge is aggravated in the peptide domain as fewer than 3.5% PDB structures contain 50 residues or less. In fact, applying 046 the SCOP classification filter in the PDB to display structures labelled as "Peptide" reveals only 509 047 entries, circa 0.2% of all experimentally determined structures available. 048

In addition to the lack of data, sequences are subject to composition bias. The incidence of certain amino acids may differ depending on the sequence length, as longer proteins have more options for accommodating multiple secondary structures and folding loops (Tiessen et al., 2012). Popular models like ProteinMPNN (Dauparas et al., 2022b), PiFold (Gao et al., 2022) and ESM-IF1 (Hsu

005 006

007

012 013

014

015

016

017

018

019

021

025

026 027 028

029

⁰⁵³

¹Updated figures available at https://www.rcsb.org/stats/growth/growth-released-structures.

et al., 2022b) are primarily trained on data derived from longer proteins, leading to poor performance for peptide design tasks. Additionally, shorter sequences fold into simpler structures. Indeed, Milner-White & Russell (2008) argue that short peptides are notoriously "structureless" and tend to flicker between conformations. For example, a structural conformation of a single alpha helix or beta sheet – or a combination of two or three of them – is not necessarily stable and can fluctuate.

Existing inverse folding models optimize sequence residue-recovery accuracy and structural similarity via template modeling (TM) score (Zhang & Skolnick, 2004). However, they often suffer from low sampling diversity (Gao et al., 2023b). Ideally, the inverse folding model generates maximally diverse candidate sequences, as additional design filters, such as synthesizability and thermal stability, reduce the number of potential hits downstream of the design process.

064 To address these issues, we apply Direct Preference Optimization (DPO) (Rafailov et al., 2023), a 065 fine-tuning method, to improve inverse-folding methods for peptide design. To the authors' knowl-066 edge, we are the first to apply DPO to this task. We propose several enhancements to DPO to address 067 specific problems that arise in inverse folding. Particularly, we forward fold generated sequences 068 and derive an online regularized algorithm for optimizing structural similarity to the reference and 069 sequence diversity simultaneously. We show empirically that this algorithm targets the differential 070 entropy in log-probability space. Furthermore, we present a simple reward scaling approach to in-071 corporate scalar reward information, showing that reward scaling adaptively selects a KL divergence penalty (Kullback, S. and Leibler, R. A., 1951) to improve performance on harder structures. 072

073 074

2 PRELIMINARIES

075 076

Inverse folding is the problem of inferring the sequence of amino acids y that fold into a 077 given protein structure. This protein structure is represented by a set of coordinates x = $(x_i, y_i, z_i) | i = 1, \dots, n \in \mathbb{R}^{3n}$, where each (x_i, y_i, z_i) represents the 3D position of a backbone 079 atom and n the number of atoms. The inverse folding problem is underconstrained; there may be 080 many solutions y that fold into structures similar to x (Koehl & Levitt, 1999). Prior research is pri-081 marily concerned with recovering the "ground-truth" sequence y_x from the experimental (reference) structure (Dauparas et al., 2022a; Hsu et al., 2022b; Jing et al., 2021b;a). Recently, forward folding 083 models like AlphaFold (Jumper et al., 2021), ESMFold (Lin et al., 2022), and OpenFold (Ahdritz 084 et al., 2022), have made it possible to estimate the structural similarity between generated sequences 085 and the reference structure. In this work, we focus on measuring structural similarity of generated 086 sequences to the reference structure via the self-consistency TM-score (sc-TM) (Gao et al., 2023b).

087 **ProteinMPNN** (Dauparas et al., 2022a) is a popular inverse-folding method that produces full pro-880 tein sequences from backbone features (distances and orientations between backbone atoms, back-089 bone dihedral angles, accounting for a virtual C β atom). ProteinMPNN is a 6-layer encoder-decoder 090 message-passing neural network based on the 'Structured Transformer' (Ingraham et al., 2019). Un-091 like autoregressive methods like ESM-IF1 (Hsu et al., 2022b), ProteinMPNN decodes residues in a 092 random decoding order, as opposed to a fixed left-to-right order from N-terminus to the C-terminus. 093 ProteinMPNN is trained on examples from the Protein Data Bank (PDB) to determine the most likely residues for a given protein backbone. On native protein backbones, it achieves a sequence 094 recovery of 52.4%, compared to 32.9% for Rosetta (Leman et al., 2020). In this work we build upon 095 ProteinMPNN by proposing methods for adapting it to peptide design. 096

097 **DPO for inverse folding**. Direct Preference Optimization (DPO) is a popular method for aligning 098 Large Language Models to a dataset of human-produced preference assignments, that discriminate amongst the responses grouped by prompt (Rafailov et al., 2023). We adapt DPO to fine-tune 099 inverse-folding models by replacing human preference labels on generated sentences with TM-score 100 rankings on generated sequences. Specifically, we generate a dataset of sequences conditioned on a 101 set of reference structures, and score each sequence with the TM-score computed between its pre-102 dicted structure and the reference structure. These scores define preference pairs over the generated 103 sequences, for each reference structure. 104

Let r(x, y) be the TM-score between structure x and the predicted fold of sequence y, and π_{ref} be a conditional distribution over y given x. Given a dataset D mapping structures ("prompts") to K generated sequences per structure ("responses"), DPO is derived within the KL-constrained reward-maximization objective (Ouyang et al., 2022; Rafailov et al., 2023):

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)}[r(x,y)] - \beta \mathbb{D}_{\mathrm{KL}}[\pi_{\theta}(y \mid x) \mid\mid \pi_{\mathrm{ref}}(y \mid x)]$$
(1)

where β controls the trade-off between reward maximization and deviation from the base model. The solution to Equation 1 is well-known (Ouyang et al., 2022; Rafailov et al., 2023):

$$\pi_r = \frac{1}{Z(x)} \pi_{\text{ref}}(y \mid x) \exp\left(\frac{1}{\beta} r(x, y)\right)$$
(2)

where Z(x) is the partition function to normalize π_r . Therefore, for any policy π_r , there is a corresponding reward function (DPO "implicit reward") for which π_r is optimal:

$$\mathbf{r}(x,y) = \beta \log \frac{\pi_r(y|x)}{\pi_{\mathrm{ref}}(y|x)} + \beta \log Z(x).$$
(3)

Substituting Equation 3 into the Bradley-Terry preference model (Bradley & Terry, 1952) and optimizing the policy under a maximum likelihood objective produces the DPO loss, typically solved via gradient descent:

$$\mathcal{L}_{\text{DPO}}\left(\pi_{\theta}; \pi_{\text{ref}}\right) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}}\left[\log \sigma \left(\beta \log \frac{\pi_{\theta}\left(y_w \mid x\right)}{\pi_{\text{ref}}\left(y_w \mid x\right)} - \beta \log \frac{\pi_{\theta}\left(y_l \mid x\right)}{\pi_{\text{ref}}\left(y_l \mid x\right)}\right)\right] \quad (4)$$

where y_w represents the sequence that is preferred or considered better, and y_l represents the sequence that is less preferred or considered worse according to TM-score ranking.

3 PREFERENCE OPTIMIZATION FOR PEPTIDE DESIGN

r

In this section, we consider designing fine-tuning methods well-suited for peptide design. While DPO in its original formulation is useful for fine-tuning in biology (Park et al., 2023; Widatalla et al., 2024; Mistani & Mysore, 2024), we consider how it may be improved to tackle specific problems arising in inverse-folding for peptides (i.e., poor generation diversity and lack of peptide data in initial training). First, we derive a diversity-optimized DPO loss by incorporating an online diversity penalty directly into the top-level Reinforcement Learning objective. Next, we propose an ad-hoc modification to the DPO reward that incorporates scalar TM-scores instead of solely preference pairs, and address the distribution shift between peptides and longer-length proteins.

3.1 ONLINE DIVERSITY OPTIMIZATION

To encourage diversity in generated sequences while simultaneously maximizing reward (TM-score), we consider a modified DPO objective incorporating an auxiliary diversity reward:

 $= \max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} \left[r(x, y) - \beta \log \left(\frac{\pi_{\theta}(y \mid x)}{\pi_{\text{ref}}(y \mid x)} \right) - \alpha \mathbb{E}_{y' \sim \pi_{\theta}(y' \mid x)} \left[\gamma(y, y') \right] \right]$

(5)

147 148 149

108

110 111

112

118

119 120 121

122 123

124

125

131

132 133

134 135

136

137

138

139

140

141

142

143 144

145 146

where α controls the strength of diversity regularization, γ is a pairwise distance between sequences, and $\Gamma(\pi_{\theta})$ is the diversity of policy π_{θ} under the distance γ , i.e. $\Gamma(\pi) = \mathbb{E}_{y,y' \sim \pi} [\gamma(y,y')]$. In practice, we let γ be the fraction of pairwise different tokens in equal-length sequences y and y'. This is equivalent to the diversity metric defined in (Gao et al., 2023b).

 $\max_{x \sim D, y \sim \pi_{\theta}(y|x)} [r(x,y)] - \beta \mathbb{D}_{\mathrm{KL}} [\pi_{\theta}(y \mid x) \mid \mid \pi_{\mathrm{ref}}(y \mid x)] + \alpha \Gamma_{\gamma}(\pi_{\theta})$

This is similar in form to the auxiliary reward objective proposed by Park et al. (2024) and Zhou et al. (2024), but the auxiliary reward (diversity) depends on the policy π_{θ} . As a result, the standard analytic solution to Eq. 1 (Ziebart et al., 2008; Ouyang et al., 2022; Rafailov et al., 2023) is no longer valid. Instead, we consider an approximate objective, leveraging the fact that the DPO loss is solved via iterative gradient descent. Let $\tilde{\pi} = \pi_{\theta}^{(t-1)}$ be an approximation of $\pi_{\theta}^{(t)}$. That is, we will use the policy from a previous iteration to estimate diversity, while updating the policy in the current iteration. Then, let $\tilde{r}(x, y) = r(x, y) - \alpha \mathbb{E}_{y' \sim \tilde{\pi}(y|x)} [\gamma(y, y')]$. We can approximate Eq. 5 as:

 $\max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} \left[\tilde{r}(x, y) - \beta \log \left(\frac{\pi_{\theta}(y \mid x)}{\pi_{\text{ref}}(y \mid x)} \right) \right].$ (6)

The rest of the derivation follows from Rafailov et al. (2023). We produce the diversity-regularized implicit reward after writing the \tilde{r} -optimal policy $\pi_{\tilde{r}} = \frac{1}{Z(x)} \pi_{ref}(y \mid x) \exp\left(\frac{1}{\beta}\tilde{r}(x,y)\right)$:

$$r(x,y) = \beta \log \frac{\pi_r(y|x)}{\pi_{\text{ref}}(y|x)} - \alpha \mathbb{E}_{y' \sim \tilde{\pi}(y' \mid x)} \left[\gamma(y,y')\right] + \beta \log Z(x) \tag{7}$$

The resulting MLE loss under the Bradley-Terry preference model is:

$$\mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \Big[\log \sigma \Big(\beta \log \frac{\pi_{\theta} (y_w \mid x)}{\pi_{\text{ref}} (y_w \mid x)} - \beta \log \frac{\pi_{\theta} (y_l \mid x)}{\pi_{\text{ref}} (y_l \mid x)} + \alpha \mathbb{E}_{y' \sim \tilde{\pi}(y' \mid x)} \left[\gamma(y_l, y') \right] \\ - \alpha \mathbb{E}_{y' \sim \tilde{\pi}(y' \mid x)} \left[\gamma(y_w, y') \right] \Big) \Big]$$
(8)

The α -weighted scalar penalties require sampling from the approximate policy $\tilde{\pi}$ to compute the expectation. In practice, we only update $\tilde{\pi}$ a few times during training to minimize the cost of sampling. See Appendix A.2 for a detailed description of the diversity-regularized algorithm.

Unlike prior methods incorporating auxiliary rewards into DPO (Park et al., 2024; Zhou et al., 2024), ours includes an online sampling term. Empirically, we show that this online term is crucial. In the middle part of Fig. 1, we show that diversity of samples from the static dataset (the auxiliary reward prescribed in an offline approach) does not correlate with the diversity of samples from the trained policy (our auxiliary reward). That is, to accurately estimate – and therefore encourage – diversity over the course of training, online sampling is required. By optimizing the approximate objective in Equation 6, we provide principled motivation for including this online sampling term.

1953.2 ENTROPY IN DIVERSITY OPTIMIZATION

197 Here, we consider the mechanism behind diversity regularization, showing that diversity improves due to randomness in the token decoding order during sampling. Particularly, we develop a new 199 understanding of ProteinMPNN based on random decoding order. ProteinMPNN performs randomorder token decoding during sampling and all forward passes. So, for a fixed policy π , sequence y, 200 structure x, we highlight that $\log \pi(y \mid x)$ is a random variable, depending on a distribution over 201 token decoding orders. Let $\ell_{(\pi,x,y)} = \log \pi(y \mid x)$ denote this random variable, and let $f_{\pi}(d)$ 202 be a function parameterized by π that takes decoding orders d to log-probabilities s. In practice, 203 f_{π} is a forward pass through ProteinMPNN, and d is a uniformly chosen permutation over indices 204 $\{1, 2, \dots, |y|\}$. For a decoding order $d \sim \mathcal{U}$ sampled from the uniform decoding order distribution 205 \mathcal{U} , we can compute a log-probability $s = f_{\pi}(d)$. For this $s, \ell_{(\pi,x,y)}(s)$ gives the probability, over all 206 possible d, that the log-probability of y|x under π is s. 207

Revisiting training with random log-probabilities. This new viewpoint allows us to re-evaluate the training of random-order decoder models, like ProteinMPNN. The implicit rewards in Eqs.3 and 7, from which a MLE loss is derived, hide the stochastic nature of $\log \pi(y|x)$. In the explicit decoding order notation, we can write Eq. 3 as

212

166

167 168 169

170

171 172 173

174 175

181 182 183

185

186

187

213 214

$$r(x,y) = \beta \left[\log \left(f_{\pi_r}(d) \right) - \log \left(f_{\pi_{\text{ref}}}(d) \right) \right] + \beta \log Z(x), \quad d \sim \mathcal{U}$$
(9)

with a similar substitution for Eq. 7. This applies to SFT as well, as the loss function has a $\log \pi(y|x)$ term to which we can apply the same substitution. Note that Eq. 9 is computed over a single sample

216 d. That is, training random-order decoders usually involves computing losses with a single-sample 217 estimate of $\log \pi(y|x) = \mathbb{E}_{d \sim \mathcal{U}} \left[\ell_{(\pi, x, y)}(f_{\pi}(d)) \right]$. While not explored here, a simple approach to 218 improving the quality of this estimate is to compute log-probabilities over multiple decoding orders, 219 though this would linearly scale training time as well.

Explaining diversity with differential entropy. This view also offers insight into how diversity increases without token entropy increasing as well. Within the random log-probability view, there is another source of entropy to account for: the differential entropy \mathcal{H}_d of the distribution over log-probabilities, defined separately for each tuple (π, x, y) :

224 225

220

221

222

223

226 227 228

229

230

231

$$\mathcal{H}_d(\ell_{(\pi,x,y)}) = \mathbb{E}_{d \sim \mathcal{U}} \left[-\log \left[\ell_{(\pi,x,y)}(f_\pi(d)) \right] \right] = -\int_{-\infty}^0 \ell_{(\pi,x,y)}(s) \log \left[\ell_{(\pi,x,y)}(s) \right] ds \quad (10)$$

We claim that optimizing for diversity increases differential entropy in the continuous logprobability space, rather than increasing discrete entropy in token space. A proof sketch can be found in Appendix A.6, and empirical support for this theory is presented in Section 4.4.

232 233 234

235

3.3 LEVERAGING DOMAIN-SPECIFIC PRIORS

Aligning train set and base model. Supervised fine-tuning (SFT) is a standard part of the DPO
 pipeline (Rafailov et al., 2023), where a gradient-based optimizer maximizes the log-probability of
 a target sequence prior to applying DPO. The left part of Fig. 1 shows that the distribution of to kens sampled from base ProteinMPNN and the peptide training dataset differs significantly. SFT
 assuages the token distribution problem by improving alignment with training distribution, and en sures consistency with previous research, where DPO has been applied to related biological tasks
 (Park et al., 2023; Widatalla et al., 2024; Mistani & Mysore, 2024).

Incorporating scalar rewards. Given multiple responses from a single prompt, DPO is derived assuming access to pairwise preferences. However, since we measure the reward of a sequence generation by the TM-score between the original structure and the generated sequence's predicted structure, we have access to a total ordering over generations through scalar rewards for each response. Widatalla et al. (2024) derives weighted DPO to incorporate scalar rewards, but it does not substantially outperform standard DPO.

249 We consider a simple ad-hoc method to incorporate these scalar scores. Given a structure x, K250 generated sequences $y_k | x$, and K corresponding TM-scores $r_k | x$, we scale the log-probabilities of 251 π_{ref} by the average of $r_k | x$. To see the effect of this scaling, consider the modified implicit reward:

$$r_{\text{scale}}(x,y) = \beta \left[\log \pi_{\theta}(y \mid x) - R(x) \cdot \log \pi_{\text{ref}}(y \mid x) \right] + \beta \log Z(x)$$

= $\beta R(x) \left[\log \pi_{\theta}(y \mid x) - \log \pi_{\text{ref}}(y \mid x) \right] + (\beta - \beta R(x)) \log \pi_{\theta}(y \mid x) + \beta \log Z(x)$
= $\left[\frac{\beta R(x)}{Reweighted standard reward} + \frac{\beta - \beta R(x)}{Reweighted standard reward} + \frac{\beta R(x$

258 259

260 When R(x) is large, the first term (mirroring the standard DPO reward, Eq. 3) has larger weight, and 261 the second term (a penalty on the entropy $\mathcal{H}(\pi_{\theta}) = -\mathbb{E} \left[\log \pi_{\theta}(y \mid x) \right]$) has smaller weight. Since 262 the weight on the DPO reward controls the strength of KL divergence regularization (Rafailov et al., 263 2023), large R(x) corresponds to less aggressive optimization and a lower-entropy policy.

264 When R(x) is near 1.0, the data-generating policy (i.e., base ProteinMPNN) already produces high-265 quality sequences on structure x. Therefore, we perform less aggressive optimization when R(x) is 266 large, and vice versa. Effectively, per structure, this reward scaling adapts the KL divergence penalty 267 to the strength of the dataset-generating policy. The right part of Fig. 1 shows that this reward scaling 268 method indeed helps ProteinMPNN accurately recover the total ordering over generated sequences. 269 Scaling log-probabilities improves their ability to rank sequences by TM-score, which provides a 269 more accurate reference model for DPO to start with.



Figure 1: Motivation for DPO design choices. **Left.** Frequency of amino acid across ProteinMPNN generations conditioned on the peptide train set, vs. frequency over the peptide sequences. **Middle.** When conditioned on the same structure, the diversity of sequences generated by base ProteinMPNN does not correlate with the diversity of sequences generated by fine-tuned ProteinMPNN. **Right.** Distribution of rank correlation coefficient between model log-probabilities and TM-score.

4 **RESULTS**

279

280

281

282

283 284

285 286

287

288

289

290

291

292 293 Here we present results on benchmarks across a suite of inverse folding models evaluated on peptide design tasks, as well as an exploration into the behavior of the two proposed algorithm enhancements (diversity regularization and reward scaling). We find that diversity regularization is effective in improving sampling diversity, and provide justification as to how this improvement happens. Additionally, we find reward scaling produces a small improvement in TM-score, enabling fine-tuned ProteinMPNN to reach SOTA performance in some situations.

4.1 EXPERIMENTAL OVERVIEW

295 Datasets. We fine-tune trained ProteinMPNN from Dauparas et al. (2022a) on a set of 211,402 dedu-296 plicated peptide structures with length up to 50 amino acids, derived from the ColabFold database 297 (Mirdita et al., 2022) by filtering for predicted local distance difference test (pLDDT) > 80. Each 298 structure was used to generate 4 candidate sequences using pretrained ProteinMPNN with T = 0.1. 299 Each sequence, including the true reference sequence from the structure, was folded with OpenFold (Ahdritz et al., 2022). TM-scores were computed to create a ranking over generated sequences for 300 every structure prompt. Each structure therefore contributed $\binom{4}{2} = 6$ chosen-rejected pairs for DPO, 301 for a total of 1,268,412 training pairs. Details of the folding process can be found in Appendix A.1. 302

303 Benchmarks. We consider two non-overlapping benchmarks. First, we take 50 sequences from 304 the OpenFold set, enforcing a sequence identity threshold of 0.4 from the train set and filtering for 305 sequence length L < 50. Next, we take the CATH 4.3 test split, as used in Hsu et al. (2022b), filter 306 for sequence identity threshold of 0.4 and $L \leq 60$, and use the resulting 173 peptides as a second benchmark. The OpenFold set contains structures resolved via OpenFold, while the CATH split has 307 structures from the PDB. For evaluation, we sample 4 sequences per benchmark structure at T = 0. 308 We compute diversity, TM-score, and recovery metrics as in Gao et al. (2023b). Following Section 309 3.1, we define diversity as the average fraction of non-identical amino acids, computed pairwise 310 across all sampled sequences for the same structure. 311

Hyperparameters. Supervised fine-tuning was run for 2 epochs with Adam through the PyTorch implementation (Paszke et al., 2019; Kingma & Ba, 2014), with structures' true sequence as the target. DPO was run for 20 epochs with default Adam hyperparameters on pairs of generated sequences. We vary β based on the method used to facilitate a fair comparison across similar KL budgets. For more details, see Appendix A.1.

Sweeps. We run DPO with diversity penalties $\alpha = \{0.0, 0.1, 0.2, 0.5\}$ and DPO with/without reward scaling, all on the same training set. All training is done with per-GPU batch size of 32 on a single node with 8xA100 80GB GPUs. Each run takes about 8 hours wall-clock time.

320

322

321 4.2 BENCHMARK SWEEPS

We benchmark our trained models against standard inverse-folding methods (Jing et al., 2021a;b; Hsu et al., 2022b; Dauparas et al., 2022a) across both the OpenFold and CATH benchmarks, filtered 324 for peptide-length proteins. We consider TM-score to be the most important metric, as unlike se-325 quence recovery, it can more accurately reflect the quality of generated sequences at high diversities. 326 As shown in Table 1, on OpenFold structures, all DPO methods outperform base ProteinMPNN and 327 other models by at least 8%. Diversity-regularized DPO ($\alpha = 0.1$) achieves state-of-the-art (SOTA) diversity, improving base DPO by 20% and even exceeding the diversity of base ProteinMPNN. 328 This is notable, as fine-tuning tend to decrease generation diversity (Wang et al., 2023). Combining 329 reward scaling with diversity regularization does not have a strong beneficial effect, with similar 330 performance compared to the regularized method. 331

Table 1: Comparison of models on TM-score, sampled sequence diversity, and native sequence 333 recovery for the OpenFold benchmark. The best results are bolded, and the second-best results are 334 underlined. Means and standard errors are reported in each cell. 335

336							
337	Models	OpenFold benchmark $(n = 50)$					
338	Woulds	TM-Score ↑	Diversity ↑	Recovery ↑			
339	GVP-GNN (Jing et al., 2021b)	0.62 ± 0.02	0.19 ± 0.01	0.27 ± 0.01			
340	ProteinMPNN (Dauparas et al., 2022a)	0.62 ± 0.01	0.31 ± 0.01	0.23 ± 0.01			
341	ESM-IF1 (Hsu et al., 2022b)	0.61 ± 0.01	$\overline{0.00\pm0.00}$	0.31 ± 0.01			
342	ProteinMPNN + DPO	$\underline{0.66 \pm 0.02}$	0.27 ± 0.01	$\overline{\textbf{0.33}\pm\textbf{0.01}}$			
343	ProteinMPNN + DPO (scaled)	$\overline{\textbf{0.67}\pm\textbf{0.02}}$	0.28 ± 0.01	0.31 ± 0.01			
344	ProteinMPNN + DPO ($\alpha = 0.1$)	$\textbf{0.67} \pm \textbf{0.02}$	$\textbf{0.32} \pm \textbf{0.01}$	0.31 ± 0.01			
345	ProteinMPNN + DPO (scaled, $\alpha = 0.1$)	$\textbf{0.67} \pm \textbf{0.02}$	$\underline{0.31\pm0.01}$	$\underline{0.32\pm0.01}$			

347 On the CATH benchmark, DPO does not help much (Table 2). This is likely due to the fact that 348 our training method leveraged OpenFold structure predictions, while the CATH benchmark contains 349 experimentally resolved protein structures. However, diversity regularization and reward scaling allow fine-tuned ProteinMPNN to nearly reach the performance of ESM-IF1, which was trained 350 on experimentally resolved structures and is a much larger model. Diversity regularization is still 351 effective in promoting diversity, beating base ProteinMPNN and improving standard DPO by 7%. 352

Table 2: Comparison of models on TM-score, sampled sequence diversity, and native sequence recovery for the CATH 4.3 benchmark. The best results are bolded, and the second-best results are underlined. Means and standard errors are reported in each cell. Since all values are rounded to 2 decimals, standard errors are not actually zero.

259							
350	Models	CATH 4.3 benchmark $(n = 173)$					
360	Woucis	TM-Score ↑	Diversity ↑	Recovery \uparrow			
361	GVP-GNN (Jing et al., 2021b)	0.69 ± 0.01	0.21 ± 0.00	$\textbf{0.35} \pm \textbf{0.01}$			
362	ProteinMPNN (Dauparas et al., 2022a)	0.68 ± 0.01	0.30 ± 0.00	0.32 ± 0.01			
363	ESM-IF1 (Hsu et al., 2022b)	$\textbf{0.72} \pm \textbf{0.01}$	$\overline{0.00\pm0.00}$	$\underline{0.34 \pm 0.01}$			
364	ProteinMPNN + DPO	0.70 ± 0.01	0.28 ± 0.00	0.32 ± 0.01			
365	ProteinMPNN + DPO (scaled)	$\underline{0.71\pm0.01}$	0.29 ± 0.00	0.32 ± 0.01			
366	ProteinMPNN + DPO ($\alpha = 0.1$)	0.70 ± 0.01	$\textbf{0.31} \pm \textbf{0.01}$	0.32 ± 0.01			
367	ProteinMPNN + DPO (scaled, $\alpha = 0.1$)	$\underline{0.71\pm0.01}$	$\underline{0.30\pm0.00}$	0.32 ± 0.01			

369 We find that both reward scaling and diversity regularization are effective in improving TM-score 370 over base DPO, achieving SOTA performance on the OpenFold benchmark. Diversity regularization is particularly effective in improving diversity across both benchmarks. Despite CATH structures 371 being out-of-distribution from the train set, given that they were not produced by OpenFold, our 372 DPO enhancements allow ProteinMPNN to approach the performance of ESM-IF1. 373

374 375

376

368

332

346

353

354

355

356

357

4.3 PARETO FRONT WITH DIVERSITY REGULARIZATION

In this section, we consider the impact of diversity regularization on DPO. We train four DPO models 377 on top of fine-tuned ProteinMPNN ($\alpha = \{0.0, 0.1, 0.2, 0.5\}$), and generate sequences across a range



Figure 2: Exploring the effect of online diversity optimization. Left. Improvement in diversity across temperatures 0, 0.5, and 1.0, with and without random decoding order. Middle. Sampling distribution entropy (average negative log-probability over samples) over various α values. Right. Entropy of log-probability distribution over validation reference sequences across α sweep.



Figure 3: Pareto front and KL divergences. Left. Pareto front for various α values over temperature sweep. Middle. KL divergence from π_{ref} for various α values. $\alpha = 0$ has $\beta = 0.5$, all other have $\beta = 0.1$. Right. KL divergence for DPO with reward scaling ($\beta = 0.1$) and without ($\beta = 0.5$).

of temperatures on the OpenFold test split. As temperature increases, the diversity of generated sequences increases, but the average TM-score (reward) decreases. We consider this reward-diversity tradeoff on the left side of Figure 3. At sufficiently low α values, diversity-regularized DPO produces a new Pareto frontier. With $\alpha = 0.1$, regularized DPO consistently achieves higher reward at the same diversity, indicating this regularization provides a strictly favorable reward-diversity tradeoff compared to simply increasing temperature. At temperature 0.0, $\alpha = 0.1$ yields a 20% relative improvement in diversity along with a small increase in TM-score (1.5%) over standard DPO.

However, at high α values, diversity regularization hurts both diversity and reward. We see symptoms of this in the KL divergences between trained DPO policies and initial fine-tuned Protein-MPNN (middle of Fig. 3). While for $\alpha = \{0.0, 0.1, 0.2\}$, KL divergence trajectories are similar during training, DPO with $\alpha = 0.5$ produces much higher divergences (i.e., the trained policy deviates significantly more from base ProteinMPNN). In this case, aggressive diversity optimization led to a collapse in model capacity via excessive deviation from the initial policy.

418 419 420

421

4.4 DIVERSITY REGULARIZATION TARGETS DIFFERENTIAL ENTROPY

We now consider the explanation for diversity improvement presented in Section 3.2. Naively, we 422 would expect diversity regularization to directly increase the entropy of the sampling distribution 423 $\mathcal{H}(\pi) = -\mathbb{E}_{x \sim \mathcal{D}, y \sim \pi} [\log \pi(y|x)]$. In the middle part of Fig. 2, we show that this is not the case. 424 Apart from $\beta = 0.5$ (an outlier as discussed in Section 4.3), increasing diversity does not increase 425 entropy, which seems contradictory. However, this is expected under the differential entropy for-426 mulation in Section 3.2, which argued that diversity optimization increases differential entropy in 427 the continuous log-probability space. To test this, we draw 128 samples from $\ell_{(\pi,x,y)}$ for all (x,y)428 structure-sequence test pairs, i.e., compute log-probabilities with 128 different decoding orders. In 429 the right part of Fig. 2, we show that the estimated \mathcal{H}_d increases with larger α , as expected under this theory. This agrees with our intuition, since \mathcal{H}_d controls variability in model log-probabilities, 430 which decide the next token during sampling. That is, larger \mathcal{H}_d leads to greater log-probability 431 variability, leading to greater diversity.

402

403 404

387

388

389



Figure 4: Sequence recovery with diversity optimization. Left. Stronger diversity regularization seems to hurt sequence recovery, though all fine-tuned models improve over ProteinMPNN. Middle. Diversity does not hurt the correlation between log-probabilities and TM-score. Right. Best-of-N sampling allows diverse models to achieve sequence recoveries comparable to standard models.

443 444 445

441

442

446 Effects at higher temperatures. In the left half of Fig. 2, we conduct a small ablation study to isolate the effect of random decoding order across temperatures. We find that at low temperatures, 447 random decoding order is necessary for diversity. This is because the next token is chosen determin-448 istically at T = 0, so entropy in the log-probabilities is the sole contributor to diversity. However, 449 at higher temperatures, the influence of using random decoding orders is less pronounced, i.e., tem-450 perature has a larger effect on next-token sampling compared to an increasing \mathcal{H}_d . Therefore, this 451 analysis applies only for random decoding order models at T = 0. For example, left-to-right autore-452 gressive methods have fixed decoding orders, meaning $\ell_{(\pi,y,x)}$ collapses around a single scalar and 453 the differential entropy is zero. At T = 0, this means that we cannot improve diversity by boosting 454 \mathcal{H}_d , and at higher T, \mathcal{H}_d does not empirically affect diversity much anyways. 455

Increased token entropy does not help diversity. We also try optimizing for discrete token-level entropy, and achieve around a 15% increase in this entropy. However, in line with the differential entropy formulation, this method does not improve sample diversity at temperature 0. See Appendix A.3 for more details and a complete derivation of the token-entropy regularized algorithm.

460 461 4.5 SEQUENCE RECOVERY WITH DIVERSITY-REGULARIZED MODELS

In this section, we consider how diversity optimization affects ProteinMPNN's ability to recover native sequences from structures, and accurately predict the quality of generated sequences.

First, we explore whether sequence recovery (i.e., the ability to recover the conditioning structure's 465 sequence) is still possible with diversity-regularized fine-tuning. Intuitively, it may seem that models 466 with higher sampling diversity have lower native sequence recovery rate. In the left half of Fig. 4, 467 it seems like this is the case, with sequence recovery rates decreasing as α increases across all 468 temperatures. However, for small alpha ($\alpha = 0.1$), the drop in sequence recovery compared to 469 standard DPO is small. Additionally, this recovery gap disappears if we allow for more compute 470 during inference time. In the right side of Fig. 4, we take N samples from each model and compute 471 the best-of-N sequence recovery rate. As N grows, the gap between standard DPO recovery rate 472 and $\alpha = 0.1$ recovery rate shrinks to nearly zero. Therefore, optimizing for diversity does not 473 significantly impact sequence recovery rate, particularly as more sequences are sampled.

Gao et al. (2023b) claims that log-probabilities may correlate with sequence quality, e.g. TM-score.
Given that diversity regularization increases log-probability entropy as shown in Section 4.4, one might expect this correlation to be less strong under more diverse models. As shown in the middle of Fig. 4, this is not the case. Indeed, across all models, log-probabilities do not correlate with TM-score.

479 480

481

4.6 REWARD SCALING REINFORCES BIOLOGICAL PRIORS

Next, we consider the effect of reward scaling on DPO's behavior. On the left side of Fig. 5, the reward-diversity curves for DPO and reward-scaled DPO are computed across $T = \{0.0, 0.1, 0.2, ..., 1.0\}$. Reward-scaled DPO achieves consistently higher reward at the same diversity, indicating it is a Pareto improvement over standard DPO. Moreover, as in the right side of Fig. 3, reward-scaled DPO operates on a slightly smaller KL divergence budget compared to



Figure 5: Reward scaling improves DPO. Left. Reward-scaled DPO is a Pareto improvement over standard DPO over a temperature sweep. Middle. Left axis is the KL divergence between the token frequencies in the peptide train set vs. model samples (lower is better), right axis is the fraction of non-repeating tokens (higher is better). **Right.** TM-score improvement over base DPO. Lower TM-score buckets contain structures for which base ProteinMPNN generated low-quality sequences.

standard DPO. Despite maintaining a smaller deviation from pretrained ProteinMPNN, the policy trained with reward-scaled DPO is still strictly better than the policy trained with standard DPO.

504 The motivation for reward scaling, as presented in Eq. 11, is dynamic β selection based on the strength of the initial policy (or equivalently, the difficulty of the prompt), where β controls the KL 505 divergence from the reference model. On the right side of Fig. 5, we show the empirical effect 506 of reward scaling matches this intuition. For structure prompts whose ProteinMPNN-generated se-507 quences had low reward (i.e., where the base model performed poorly), reward-scaled DPO outper-508 forms standard DPO by around 5%. However, as the average TM-score of the generated sequences 509 increases, the performance gain from reward scaling drops to nearly 0. Therefore, reward-scaled 510 DPO targets hard examples where the data-generating policy is weak, agreeing with the motivation 511 in Section 3.3.

512 513

495

496

497

498

499 500 501

5 LIMITATIONS

514 515 516

517 While we illustrate an empirical connection between diversity optimization and differential entropy, 518 we do not establish a mathematical framework for how this optimization happens. Furthermore, 519 during the derivation of diversity-regularized DPO, we approximate π^* with the latest iteration of 520 gradient descent. Exploring the bounds on this approximation, and establishing a theory-first per-521 spective for the differential entropy framework, are promising directions for future research.

Since DPO and the variants proposed here are model-agnostic, they can be used to fine-tune any
inverse-folding model. Applying DPO to other inverse-folding models may produce even better
results compared to ProteinMPNN Gao et al. (2023b). For example, it may be possible to further
push the frontier of peptide sequence design by fine-tuning stronger base models like PiFold (Gao
et al., 2022) or KW-Design (Gao et al., 2023a).

527 528

6 CONCLUSION

529 530

531 We fine-tuned ProteinMPNN, a widely adopted inverse folding model, for diverse and structurally 532 consistent peptide sequence generation and proposed diversity-regularized DPO with an online sam-533 pling term to accurately estimate and encourage diversity. Domain-specific priors were also incorpo-534 rated into our methodology to account for peptides' residue distribution and dynamically control the 535 strength of KL divergence regularization. Our approach results in improvements on sampling diver-536 sity, sequence recovery and structural similarity of the generated peptide sequences. Furthermore, 537 we give additional intuition on the impact of diversity regularization on differential and discrete entropy. While our results are reported using ProteinMPNN as a base model for fine-tuning, our 538 proposed methods are agnostic to the inverse folding model, setting the grounds for future research in peptide design via fine-tuning.

540 REFERENCES

554

574

575

576

577

582

583

584 585

586

587 588

589

590

Nadia C. Abascal and Lynne Regan. The past, present and future of protein-based materials. *Open Biology*, 8(10), October 2018. ISSN 2046-2441. doi: 10.1098/rsob.180113. URL http://dx. doi.org/10.1098/rsob.180113.

Gustaf Ahdritz, Nazim Bouatta, Christina Floristean, Sachin Kadyan, Qinghui Xia, William 546 Gerecke, Timothy J O'Donnell, Daniel Berenberg, Ian Fisk, Niccolò Zanichelli, Bo Zhang, Arka-547 diusz Nowaczynski, Bei Wang, Marta M Stepniewska-Dziubinska, Shang Zhang, Adegoke Ojew-548 ole, Murat Efe Guney, Stella Biderman, Andrew M Watkins, Stephen Ra, Pablo Ribalta Lorenzo, 549 Lucas Nivon, Brian Weitzner, Yih-En Andrew Ban, Peter K Sorger, Emad Mostaque, Zhao Zhang, 550 Richard Bonneau, and Mohammed AlQuraishi. OpenFold: Retraining AlphaFold2 yields new 551 insights into its learning mechanisms and capacity for generalization. *bioRxiv*, 2022. doi: 552 10.1101/2022.11.20.517210. URL https://www.biorxiv.org/content/10.1101/ 2022.11.20.517210. 553

- Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. ISSN 00063444, 14643510. URL http://www.jstor.org/stable/2334029.
- Karuppiah Chockalingam, Mark Blenner, and Scott Banta. Design and application of stimulusresponsive peptide systems. *Protein Engineering, Design & Selection*, 20(4):155–161, 2007.
- Dana Maria Copolovici, Kent Langel, Elo Eriste, and Ulo Langel. Cell-penetrating peptides: design,
 synthesis, and applications. *ACS nano*, 8(3):1972–1994, 2014.
- Paul E. Correa. The building of protein structures from alpha-carbon coordinates. Proteins: Structure, Function, and Bioinformatics, 7(4):366–377, 1990. doi: https://doi.org/10.1002/ prot.340070408. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/ prot.340070408.
- J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky, A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan, B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, and D. Baker. Robust deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022a. doi: 10.1126/science.add2187. URL https://www.science.org/doi/abs/10. 1126/science.add2187.
 - Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles, Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning– based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022b.
- Ken A. Dill, S. Banu Ozkan, M. Scott Shell, and Thomas R. Weikl. The protein folding problem. *Annual Review of Biophysics*, 37(1):289–316, June 2008. ISSN 1936-1238. doi: 10.1146/annurev.biophys.37.092707.153558. URL http://dx.doi.org/10.1146/ annurev.biophys.37.092707.153558.
 - Zhangyang Gao, Cheng Tan, Pablo Chacón, and Stan Z Li. Pifold: Toward effective and efficient protein inverse folding. *arXiv preprint arXiv:2209.12643*, 2022.
 - Zhangyang Gao, Cheng Tan, and Stan Z. Li. Knowledge-design: Pushing the limit of protein design via knowledge refinement, 2023a. URL https://arxiv.org/abs/2305.15151.
 - Zhangyang Gao, Cheng Tan, Yijie Zhang, Xingran Chen, Lirong Wu, and Stan Z. Li. Proteininvbench: Benchmarking protein inverse folding on diverse tasks, models, and metrics. In *Thirtyseventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023b. URL https://openreview.net/forum?id=bqXduvuW5E.

592

593 A Godzik, A Kolinski, and J Skolnick. De novo and inverse folding predictions of protein structure and dynamics. J. Comput. Aided Mol. Des., 7(4):397–438, August 1993.

- ⁵⁹⁴ Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. *ICML*, April 2022a. doi: 10.1101/2022.04.10.487779. URL http://dx.doi.org/10.1101/2022.04.10.487779.
- Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. *bioRxiv*, 2022b. doi: 10.1101/2022.04.10.487779. URL https://www.biorxiv.org/content/early/2022/04/10/2022.04.10.487779.
- John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph based protein design. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox,
 and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Cur ran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/
 paper/2019/file/f3a4ff4839c56a5f460c88cce3666a2b-Paper.pdf.
- Bowen Jing, Stephan Eismann, Pratham N Soni, and Ron O Dror. Equivariant graph neural networks for 3d macromolecular structure. *arXiv preprint arXiv:2106.03843*, 2021a.
- Bowen Jing, Stephan Eismann, Patricia Suriana, Raphael John Lamarre Townshend, and Ron Dror.
 Learning from protein structure with geometric vector perceptrons. In International Conference on Learning Representations, 2021b. URL https://openreview.net/forum?id= 1YLJDvSx6J4.
- 615 Peter St. John, Dejun Lin, Polina Binder, Malcolm Greaves, Vega Shah, John St. John, Adrian Lange, 616 Patrick Hsu, Rajesh Illango, Arvind Ramanathan, Anima Anandkumar, David H Brookes, Akosua 617 Busia, Abhishaike Mahajan, Stephen Malina, Neha Prasad, Sam Sinai, Lindsay Edwards, Thomas 618 Gaudelet, Cristian Regep, Martin Steinegger, Burkhard Rost, Alexander Brace, Kyle Hippe, Luca 619 Naef, Keisuke Kamata, George Armstrong, Kevin Boyd, Zhonglin Cao, Han-Yi Chou, Simon Chu, Allan dos Santos Costa, Sajad Darabi, Eric Dawson, Kieran Didi, Cong Fu, Mario Geiger, 620 Michelle Gill, Darren Hsu, Gagan Kaushik, Maria Korshunova, Steven Kothen-Hill, Youhan 621 Lee, Meng Liu, Micha Livne, Zachary McClure, Jonathan Mitchell, Alireza Moradzadeh, Ohad 622 Mosafi, Youssef Nashed, Saee Paliwal, Yuxing Peng, Sara Rabhi, Farhad Ramezanghorbani, 623 Danny Reidenbach, Camir Ricketts, Brian Roland, Kushal Shah, Tyler Shimko, Hassan Sirelkha-624 tim, Savitha Srinivasan, Abraham C Stern, Dorota Toczydlowska, Srimukh Prasad Veccham, Nic-625 colò Alberto Elia Venanzi, Anton Vorontsov, Jared Wilber, Isabel Wilkinson, Wei Jing Wong, 626 Eva Xue, Cory Ye, Xin Yu, Yang Zhang, Guoqing Zhou, Becca Zandstein, Christian Dallago, 627 Bruno Trentini, Emine Kucukbenli, Saee Paliwal, Timur Rvachov, Eddie Calleja, Johnny Israeli, 628 Harry Clifford, Risto Haukioja, Nicholas Haemel, Kyle Tretina, Neha Tadimeti, and Anthony B 629 Costa. Bionemo framework: a modular, high-performance library for ai model development in 630 drug discovery, 2024. URL https://arxiv.org/abs/2411.10548.
- John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, 632 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, 633 Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-634 Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, 635 Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Se-636 bastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Push-637 meet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold. 638 Nature, 596(7873):583-589, 2021. ISSN 1476-4687. doi: 10.1038/s41586-021-03819-2. URL 639 https://doi.org/10.1038/s41586-021-03819-2.
- 641 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2014.
- Patrice Koehl and Michael Levitt. De novo protein design. i. in search of stability and specificity11edited by f. e. cohen. *Journal of Molecular Biology*, 293(5):1161–1181, 1999. ISSN 0022-2836. doi: https://doi.org/10.1006/jmbi.1999.3211. URL https://www.sciencedirect.com/science/article/pii/S0022283699932114.
- 646

640

631

647 Kullback, S. and Leibler, R. A. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86., 1951.

673

689

690

691

692

648 Julia Koehler Leman, Brian D Weitzner, Steven M Lewis, Jared Adolf-Bryfogle, Nawsad Alam, 649 Rebecca F Alford, Melanie Aprahamian, David Baker, Kyle A Barlow, Patrick Barth, Benjamin 650 Basanta, Brian J Bender, Kristin Blacklock, Jaume Bonet, Scott E Boyken, Phil Bradley, Chris 651 Bystroff, Patrick Conway, Seth Cooper, Bruno E Correia, Brian Coventry, Rhiju Das, René M 652 De Jong, Frank DiMaio, Lorna Dsilva, Roland Dunbrack, Alexander S Ford, Brandon Frenz, Darwin Y Fu, Caleb Geniesse, Lukasz Goldschmidt, Ragul Gowthaman, Jeffrey J Gray, Dominik 653 Gront, Sharon Guffy, Scott Horowitz, Po-Ssu Huang, Thomas Huber, Tim M Jacobs, Jeliazko R 654 Jeliazkov, David K Johnson, Kalli Kappel, John Karanicolas, Hamed Khakzad, Karen R Khar, 655 Sagar D Khare, Firas Khatib, Alisa Khramushin, Indigo C King, Robert Kleffner, Brian Koep-656 nick, Tanja Kortemme, Georg Kuenze, Brian Kuhlman, Daisuke Kuroda, Jason W Labonte, Ja-657 son K Lai, Gideon Lapidoth, Andrew Leaver-Fay, Steffen Lindert, Thomas Linsky, Nir London, 658 Joseph H Lubin, Sergey Lyskov, Jack Maguire, Lars Malmström, Enrique Marcos, Orly Marcu, 659 Nicholas A Marze, Jens Meiler, Rocco Moretti, Vikram Khipple Mulligan, Santrupti Nerli, Christoffer Norn, Shane Ó'Conchúir, Noah Ollikainen, Sergey Ovchinnikov, Michael S Pacella, 661 Xingjie Pan, Hahnbeom Park, Ryan E Pavlovicz, Manasi Pethe, Brian G Pierce, Kala Bharath 662 Pilla, Barak Raveh, P Douglas Renfrew, Shourya S Roy Burman, Aliza Rubenstein, Marion F 663 Sauer, Andreas Scheck, William Schief, Ora Schueler-Furman, Yuval Sedan, Alexander M Sevy, Nikolaos G Sgourakis, Lei Shi, Justin B Siegel, Daniel-Adriano Silva, Shannon Smith, Yifan 665 Song, Amelie Stein, Maria Szegedy, Frank D Teets, Summer B Thyme, Ray Yu-Ruei Wang, Andrew Watkins, Lior Zimmerman, and Richard Bonneau. Macromolecular modeling and design in 666 rosetta: recent methods and frameworks. Nat. Methods, 17(7):665–680, July 2020. 667

- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic level protein structure with a language model. July 2022. doi: 10.1101/2022.07.20.500902. URL http://dx.doi.org/10.1101/2022.07.20.500902.
- E James Milner-White and Michael J Russell. Predicting the conformations of peptides and proteins in early evolution. a review article submitted to biology direct. *Biol. Direct*, 3(1):3, January 2008.
- Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin Steinegger. Colabfold: making protein folding accessible to all. *Nature Methods*, 19(6):679– 682, June 2022. doi: 10.1038/s41592-022-01488-1. URL https://doi.org/10.1038/ s41592-022-01488-1.
- Pouria Mistani and Venkatesh Mysore. Preference optimization of protein language models as a multi-objective binder design paradigm, 2024. URL https://arxiv.org/abs/2403. 04187.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.
 URL https://arxiv.org/abs/2203.02155.
 - Ryan Park, Ryan Theisen, Navriti Sahni, Marcel Patek, Anna Cichońska, and Rayees Rahman. Preference optimization for molecular language models, 2023. URL https://arxiv.org/ abs/2310.12304.
- Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality
 in direct preference optimization, 2024. URL https://arxiv.org/abs/2403.19159.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. 2019.
- 701 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward

702 703 704	<pre>model. ArXiv, abs/2305.18290, 2023. URL https://api.semanticscholar.org/ CorpusID:258959321.</pre>
705 706 707 708 709	Varun R. Shanker, Theodora U.J. Bruun, Brian L. Hie, and Peter S. Kim. Inverse folding of protein complexes with a structure-informed language model enables unsupervised antibody evolution. December 2023. doi: 10.1101/2023.12.19.572475. URL http://dx.doi.org/10.1101/2023.12.19.572475.
710 711 712 713	Martin Steinegger and Johannes Soeding. Mmseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. <i>Nature Biotechnology</i> , 35(11):1026–1028, 2017. doi: 10. 1038/nbt.3988. URL https://doi.org/10.1038/nbt.3988.
714 715 716 717	Axel Tiessen, Paulino Pérez-Rodríguez, and Luis José Delaye-Arredondo. Mathematical modeling and comparison of protein size distribution in different plant, animal, fungal and microbial species reveals a negative correlation between protein size and protein number, thus providing insight into the evolution of proteomes. <i>BMC Res. Notes</i> , 5(1):85, February 2012.
718 719 720 721	Marcelo DT Torres, Shanmugapriya Sothiselvam, Timothy K Lu, and Cesar de la Fuente-Nunez. Peptide design principles for antimicrobial applications. <i>Journal of molecular biology</i> , 431(18): 3547–3567, 2019.
722 723 724	Rein V Ulijn and Andrew M Smith. Designing peptide based nanomaterials. <i>Chemical Society Reviews</i> , 37(4):664–675, 2008.
725 726 727 728	Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints, 2023. URL https://arxiv.org/abs/2309.16240.
729 730 731 732	 Talal Widatalla, Rafael Rafailov, and Brian Hie. Aligning protein generative models with experimental fitness via direct preference optimization. <i>bioRxiv</i>, 2024. doi: 10.1101/2024.05.20.595026. URL https://www.biorxiv.org/content/early/2024/05/21/2024.05.20.595026.
733 734 735 736 737	Zhenyu Yang, Xiaoxi Zeng, Yi Zhao, and Runsheng Chen. Alphafold2 and its applications in the fields of biology and medicine. <i>Signal Transduction and Targeted Therapy</i> , 8(1), March 2023. ISSN 2059-3635. doi: 10.1038/s41392-023-01381-z. URL http://dx.doi.org/10.1038/s41392-023-01381-z.
738 739 740	Kai Yi, Bingxin Zhou, Yiqing Shen, Pietro Liò, and Yu Guang Wang. Graph denoising diffusion for inverse protein folding. 2023.
741 742 743	Kaizhi Yue and Ken A Dill. Inverse protein folding problem: designing polymer sequences. <i>Proceedings of the National Academy of Sciences</i> , 89(9):4163–4167, 1992.
744 745 746 747 748	Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure template quality. <i>Proteins: Structure, Function, and Bioinformatics</i> , 57(4):702–710, 2004. ISSN 1097-0134. doi: 10.1002/prot.20264. URL http://bioinformatics.buffalo.edu/ TM-score. Copyright 2004 Wiley-Liss, Inc.
749 750 751 752	Zhanhui Zhou, Jie Liu, Jing Shao, Xiangyu Yue, Chao Yang, Wanli Ouyang, and Yu Qiao. Beyond one-preference-fits-all alignment: Multi-objective direct preference optimization, 2024. URL https://arxiv.org/abs/2310.03708.
753 754 755	Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey. Maximum entropy inverse reinforcement learning. In <i>Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3</i> , AAAI'08, pp. 1433–1438. AAAI Press, 2008. ISBN 9781577353683.



Figure 6: Length, token, and pLDDT statistics on the peptide training set from ColabFold. Left. Distribution of amino acid tokens on the train set. Middle. Length of sequences in train set. Right. Histogram of pLDDT scores from ColabFold. This is after filtering for pLDDT > 80.

APPENDIX А

756

758

759

761

763

764

765

766

767 768 769

770

771 ADDITIONAL EXPERIMENTAL DETAILS A.1 772

773 Folding of peptides was done with the OpenFold module in Folding peptide sequences. 774 NVIDIA BioNeMo Framework (John et al., 2024), version 1.8, with default settings. We used 775 mmseqs2(Steinegger & Soeding, 2017) to generate multiple sequence alignments (MSAs), refer-776 encing against UniRef90, Small BFD and MGnify datasets, as the input. For template searches, 777 hhsearch was used with the PDB70 database. OpenFold inferencing was performed with a single set of weights, converted from an AlphaFold2 (Jumper et al., 2021) model checkpoint 778 params-model-4; typically AlphaFold2 is run with five checkpoints. Folding was run on 8 to 32 779 NVIDIA A100-SXM4-80GB GPUs, with the overall folding throughput around 1.4 seconds per sequence per GPU. We did not perform structural relaxation after folding. For model checkpoint down-781 load scripts, and database download scripts, see github.com/aqlaboratory/openfold. 782

783 **Choosing** β fairly. Since β is a proxy for specifying the amount of allowable deviation from the 784 base (reference) policy (Rafailov et al., 2023), we ensure fair comparison between DPO and its 785 variants by modifying β so that all methods operate on a similar same KL divergence budget. For the diversity-regularized and reward-scaled methods, we choose $\beta = 0.1$; for base DPO, we choose 786 $\beta = 0.5$. In the middle and left parts of Fig. 3, we show that these choices allow all models to 787 deviate from the base policy by around the same amount, with standard DPO still dominating the 788 other model's KL divergences. DPO with $\alpha = 0.5$ is an exception, but we consider this model to be 789 an outlier as described in Section 4.3. 790

Method evaluations. While recent inverse-folding methods like PiFold and KW-Design show 791 strong performance compared to ProteinMPNN (Gao et al., 2023b) and ESM-IF1 Hsu et al. (2022b), 792 we were unable to get their implementation working on the necessary timeline. As a result, we did 793 not include them in our benchmarks. 794

Train dataset statistics. In Fig. 6, we report the distribution of protein lengths, the token-level histograms, and the AlphaFold pLDDT scores for the train dataset referenced earlier. 796

797 798

799

A.2 DETAILS OF ONLINE DIVERSITY REGULARIZATION

800 In Algorithm 1 we present the full online diversity-regularized DPO algorithm. Note that γ computes 801 the pairwise diversity between a sequence y and N other sequences y', so it returns an N-length vector. $\tilde{\pi}$ is updated only once every K epochs, so in practice we only have to sample once every 802 K epochs. Between these updates, samples are cached and take up only a few megabytes of GPU 803 memory. The modifications compared to base DPO (Rafailov et al., 2023) are highlighted in blue. 804

805

806 A.3 ENTROPY-REGULARIZED DPO

807

Here, we consider optimizing for discrete entropy over the sequences sampled from ProteinMPNN. 808 The derivation is the same as in Section 3.1, but with the entropy $\mathcal{H}(\pi) = -\mathbb{E}_{y \sim \pi} \left[\log \pi(y) \right]$ as the 809 diversity penalty instead of $\Gamma(\pi)$. The objective is:

Algorithm 1 Diversity-regularized DPO

Require: Dataset $D = \{(x^{(i)}, y^{(i)}_w, y^{(i)}_l)\}$, base policy π_{ref} , KL deviation penalty β , diversity in-812 centive α , max steps N, sample frequency K, sequence distance γ , number of samples M 813 $\pi_{\theta}^{(0)} \leftarrow \pi_{\text{ref}}$ 814 $\tilde{\pi} \leftarrow \pi_{\mathrm{ref}}$ 815 for t = 1 to N do 816 if $t \mod \underset{\tilde{\pi} \leftarrow \pi_{\theta}^{(t-1)}}{K = 0} 0$ then 817 818 end if 819 $(x, y_w, y_l) \leftarrow \text{Minibatch}(D)$ \triangleright Batch size B, sequence length L 820 $r(y_w, x) \leftarrow \beta \log \frac{\pi_{\theta}^{(t-1)}(y_w|x)}{\pi_{\text{ref}}(y_w|x)}$ $r(y_l, x) \leftarrow \beta \log \frac{\pi_{\theta}^{(t-1)}(y_l|x)}{\pi_{\text{ref}}(y_l|x)}$ Chosen reward 821 822 ▷ Rejected reward 823 $s \leftarrow \text{Sample}(\tilde{\pi}, x, M)$ \triangleright Shape (B, M, L)824 $d(\tilde{\pi}, y_w, x) \leftarrow \text{Average}(\gamma(y_w, s))$ \triangleright Shape (B,)825 $d(\tilde{\pi}, y_l, x) \leftarrow \text{Average}(\gamma(y_l, s))$ \triangleright Shape (B,) $\pi_{\theta}^{(t)} \leftarrow \operatorname{argmin} - \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(r(y_w, x) - r(y_l, x) + \alpha d(\tilde{\pi}, y_w, x) - \alpha d(\tilde{\pi}, y_l, x) \right) \right]$ 827 end for 829 830 831 832 $\max_{x \sim D, y \sim \pi_{\theta}(y|x)} [r(x,y)] - \beta \mathbb{D}_{\mathrm{KL}} [\pi_{\theta}(y \mid x) \mid \mid \pi_{\mathrm{ref}}(y \mid x)] + \alpha \mathcal{H}(\pi_{\theta})$ 833 $= \max_{\pi_{\theta}} \mathbb{E}_{x \sim D, y \sim \pi_{\theta}(y|x)} \left[r(x, y) - \beta \log \left(\frac{\pi_{\theta}(y \mid x)}{\pi_{ref}(y \mid x)} \right) - \alpha \log \pi_{\theta}(y \mid x) \right]$ 834 (12)835 836 837 With the same approximation $\tilde{\pi}$ and modified reward $\tilde{r}(x,y) = r(x,y) - \alpha \log \tilde{\pi}(y \mid x)$: 838 839 840 $\mathcal{L}_{\text{DivPO}}\left(\pi_{\theta}; \pi_{\text{ref}}\right) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \Big[\log \sigma \Big(\beta \log \frac{\pi_{\theta} \left(y_w \mid x \right)}{\pi_{\text{ref}} \left(y_w \mid x \right)} - \beta \log \frac{\pi_{\theta} \left(y_l \mid x \right)}{\pi_{\text{ref}} \left(y_l \mid x \right)} + \frac{\pi_{\theta} \left(y_l \mid x \right)}{\pi_{\text{ref}} \left(y_l \mid x \right)} \Big]$ 841 842 $\alpha \log \tilde{\pi}(y_l \mid x) - \alpha \log \tilde{\pi}(y_w \mid x) \Big) \Big|$ 843 (13)844 845

Algorithm 2 Entropy-regularized DPO

846

847

Require: Dataset $D = \{(x^{(i)}, y^{(i)}_w, y^{(i)}_l)\}$, base policy π_{ref} , hyperparameters β, α , max steps $N, \tilde{\pi}$ 848 update frequency K 849 $\pi_{\theta}^{(0)} \leftarrow \pi_{\text{ref}}$ 850 $\tilde{\pi} \leftarrow \pi_{\text{ref}}$ 851 for t = 1 to N do 852 if $t \mod K = 0$ then $\tilde{\pi} \leftarrow \pi_{\theta}^{(t-1)}$ 853 854 end if 855 $(x, y_w, y_l) \leftarrow \text{Minibatch}(D)$ $r(y_w, x) \leftarrow \beta \log \frac{\pi_{\theta}^{(t-1)}(y_w|x)}{\pi_{\mathrm{ref}}(y_w|x)}$ $r(y_l, x) \leftarrow \beta \log \frac{\pi_{\theta}^{(t-1)}(y_l|x)}{\pi_{\mathrm{ref}}(y_l|x)}$ 856 Chosen reward 858 ▷ Rejected reward 859 $\hat{H}(y_w, x) \leftarrow -\alpha \log \tilde{\pi}(y_w \mid x)$ ▷ No gradient computation $\hat{H}(y_l, x) \leftarrow -\alpha \log \tilde{\pi}(y_l \mid x)$ ▷ No gradient computation 861 $\pi_{\theta}^{(t)} \leftarrow \operatorname{argmin}_{\circ} - \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(r(y_w, x) - r(y_l, x) + \hat{H}(y_w, x) - \hat{H}(y_l, x) \right) \right]$ 862 863 end for

The full algorithm is described in Algorithm 2, with deviations from base DPO (Rafailov et al., 2023)
highlighted in blue. We find that this algorithm does not produce diversity gains (Table 4) or TMscore gains (Table 3). After entropy regularization, sample entropy does increase by around 25%
compared to standard DPO; however, diversity does not improve at temperature 0. This supports
the differential entropy theory from Section 4.4, since explicitly increasing discrete entropy does not
help improve diversity.

Interestingly, as shown in Table 4, at higher temperatures, entropy-regularized DPO does slightly
 increase diversity. This again aligns with the analysis in Section 4.4, as at higher temperatures, the
 randomness introduced by the token sampling distribution dominates the randomness introduced by
 log-probability variations due to random decoding orders.

Table 3: Comparison of TM-Score for ProteinMPNN, standard DPO, and DPO with entropy regularization across different temperatures (T). Same evaluation setting as the OpenFold benchmark.

Method	TM-Score									
Witthou	T=0.0	T=0.1	T=0.2	T=0.3	T=0.4	T=0.5	T=0.6	T=0.7	T=0.8	T=0.9
ProteinMPNN	0.62	0.61	0.62	0.61	0.62	0.61	0.60	0.57	0.56	0.56
Standard DPO	0.66	0.66	0.65	0.66	0.65	0.65	0.65	0.64	0.63	0.61
DPO ($\alpha = 0.1$)	0.66	0.67	0.66	0.66	0.65	0.64	0.63	0.63	0.61	0.60

Table 4: Comparison of diversity for ProteinMPNN, standard DPO, and DPO with entropy regularization across different temperatures (T). Same evaluation setting as the OpenFold benchmark.

Method	Diversity									
	T=0.0	T=0.1	T=0.2	T=0.3	T=0.4	T=0.5	T=0.6	T=0.7	T=0.8	T=0.9
ProteinMPNN	0.31	0.35	0.39	0.47	0.54	0.59	0.65	0.71	0.74	0.78
Standard DPO	0.27	0.29	0.35	0.42	0.50	0.56	0.62	0.67	0.71	0.75
DPO ($\alpha = 0.1$)	0.26	0.30	0.38	0.46	0.53	0.60	0.65	0.70	0.73	0.77

A.4 ADDITIONAL BENCHMARK RESULTS

In Tables 5 and 6, we report the same benchmark results as in Tables 2 and 1, but with N = 64 samples per structure, as well as at T = 0.1. Only ESM-IF and ProteinMPNN results are reported.

Table 5: Comparison of methods on the OpenFold benchmark.

Method	Score	Diversity	Recovery
ESM-IF1 (Hsu et al. (2022a))	0.616 ± 0.004	0.252 ± 0.002	$\textbf{0.305} \pm 0.002$
ProteinMPNN (Jing et al. (2021b))	0.552 ± 0.004	0.421 ± 0.002	0.155 ± 0.002
ProteinMPNN + DPO (scaled)	$\textbf{0.633} \pm 0.004$	$\overline{0.404} \pm 0.002$	0.192 ± 0.002
ProteinMPNN + DPO ($\alpha = 0.1$)	0.631 ± 0.004	$\textbf{0.440} \pm 0.002$	0.194 ± 0.002
ProteinMPNN + DPO	0.625 ± 0.004	0.389 ± 0.002	$\underline{0.212}\pm0.002$

A.5 SAMPLED SEQUENCE EXAMPLE STRUCTURES

In Figures 8 and 7, we present some sequence samples conditioned on structures from both the
 OpenFold and CATH 4.3 benchmarks. We select pairs where reward scaled DPO and diversity regularized DPO outperform base ProteinMPNN for visualization.

914 A.6 DIFFERENTIAL ENTROPY AND SAMPLE DIVERSITY

We will prove that increasing differential entropy in the log-probabilities increases the expected pairwise difference in next-token sampling. Assume sampling temperature T = 0. We prove this for the slightly weaker case considering only next-token sampling, not over the entire sequence.

Table 6: Comparison of methods on the CATH4.3 benchmark.

919	_			
920	Method	Score	Diversity	Recovery
921	ESM-IF1 (Hsu et al. (2022a))	$\textbf{0.719} \pm 0.002$	0.223 ± 0.001	$\textbf{0.341} \pm 0.001$
922	ProteinMPNN (Jing et al. (2021b))	0.662 ± 0.002	$\underline{0.316} \pm 0.001$	0.318 ± 0.001
923	ProteinMPNN + DPO	0.688 ± 0.002	$\overline{0.315}\pm0.001$	0.317 ± 0.001
924	ProteinMPNN + DPO (scaled)	$\underline{0.701} \pm 0.002$	0.308 ± 0.001	$\underline{0.320}\pm0.001$
925	ProteinMPNN + DPO ($\alpha = 0.1$)	0.695 ± 0.002	$\textbf{0.340} \pm 0.001$	0.317 ± 0.001

Suppose $Y_1 \dots Y_n$ are discrete random variables over the set of amino acid tokens, where Y_i is conditioned on $Y_{j < i}$. This represents the generative process of sampling token *i* conditioning on the previous sampled tokens. For simplicity, we omit reference to the conditioning backbone structure, the distribution Y_1 can be implicitly viewed as capturing this dependence.

Let $y_1
dots y_{i-1}$ be a partial realization of $Y_1
dots Y_n$ up to the (i-1)-th token. Consider two sequences A and B, where $A = y_1
dots y_{i-1}, y_A$ and $B = y_1
dots y_{i-1}, y_B$, where y_A, y_B are independent samples from $Y_i | y_1
dots y_{i-1}$. We omit the dependence of Y_i on $y_1
dots y_{i-1}$ in the notation from here on. Per the definition of diversity in Section 3.1, the pairwise diversity on the partial sequences $\gamma(A, B)$ is dependent only on $\mathbb{1}(y_A \neq y_B)$. Assuming the next token of A, B to be conditionally independent given $y_1
dots y_{i-1}$, we consider this quantity in expectation over Y_i ,

938 939

944 945 946

918

926

$$\mathbb{E}[\mathbb{1}(y_A \neq y_B)] = 1 - P(y_a = y_b)$$

= $1 - \sum_{y \in S} P(Y_i = y, Y_i = y \mid y_1 \dots y_{i-1})$
= $1 - \sum_{y \in S} P(Y_i = y \mid y_1 \dots y_{i-1})^2$ (14)

where S is the set of amino acid tokens. Maximizing Eq. 14 corresponds to minimizing the sum of squared probabilities, i.e. $\mathbb{E}[P(Y_i)]$. Noting that log is concave and applying Jensen's inequality, we can derive a lower bound on the diversity via the discrete Shannon entropy $\mathcal{H}(Y_i)$,

947

$$\log \mathbb{E}[P(Y_i)] \ge \sum_{x \in S} P(Y_i = y) \log P(Y_i = y) = -\mathcal{H}(Y_i)$$
$$\mathbb{E}[P(Y_i)] \ge \exp\{-\mathcal{H}(Y_i)\}$$
$$P(y_A \neq y_B) \ge 1 - \exp\{-\mathcal{H}(Y_i)\}$$
(15)

953 954 955

952

That is, if the entropy of Y_i increases, so does the expected diversity of the next token ². Now, since $T = 0, Y_i = \operatorname{argmax} \{\ell_1, \dots, \ell_{|S|}\}$ where ℓ_k is the log-probability of the k-th token. In the notation of Section 3.1, $\ell_k = \ell_{(\pi,x,\{y_1\dots,y_{i-1},S_k\})}$, i.e. a random variable defined as a transformation over uniform decoding orders given by the parameters of the MPNN.

Since ℓ_k are all mutually dependent with negative pairwise correlation (the sum of their exponentials must be one) and share the same support, increasing the differential entropy $\mathcal{H}_d(\ell_k)$ across all all |S| log-probabilities necessarily increases $\mathcal{H}(Y_i)$, as the uncertainties of the individual variables determine the uncertainty of their maximum. Therefore, larger differential entropies in the log-probabilities leads to larger next-token entropies, which lower-bounds the expected next-token diversity in independent samples.

966

⁹⁶⁷ 968

²Note this is a subtly different type of entropy compared to the one referenced in the middle part of 2. When T = 0, as in this proof, discrete sequence entropy is 0, since the expectation is over the non-existent stochasticity in the token sampling process. However, when T = 0, the quantity $\mathcal{H}(Y_i)$ is nonzero, since the expectation is over the stochasticity in decoding orders.



