

# TOWARDS A GENERATIVE PROTEIN EVOLUTION MACHINE WITH DPLM-EVO

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Proteins are shaped by evolution under biophysical and functional constraints. Protein language models can learn rich evolutionary constraints, and discrete diffusion-based PLMs (DPLMs) are promising for both understanding and generation. However, existing DPLMs rely on masking-based diffusion, which is a loose proxy for evolution, and difficult to model the edit operations that drive sequence change in nature: substitutions and insertions/deletions (indels). In this paper, we present DPLM-EVO, an evolutionary discrete diffusion protein language model that explicitly predicts substitution, insertion, and deletion actions during denoising. To make indel-aware generation tractable, we introduce a latent alignment formulation that supports variable-length sequences. To make substitution corruption informative, we propose a contextual evolutionary noising kernel that generates biologically plausible mutations. On ProteinGym, DPLM-EVO achieves state-of-the-art mutation effect prediction in the single-sequence setting, and it enables variable-length generation and post-editing via explicit edit trajectories.

## 1 INTRODUCTION

Today’s rapidly growing sequence databases archive the results of protein evolution over millions of years, capturing both conserved patterns and extensive natural variation across families. For protein engineering, the practical goal is often not only to generate “protein-like” sequences, but also to leverage this evolutionary information to (i) predict the functional impact of mutations and (ii) propose variants that preserve the structure while improving or reprogramming function.

Protein language models (PLMs) trained on large protein sequence corpus have become a dominant paradigm for learning such evolutionary regularities (Lin et al., 2022; Hayes et al., 2024; Nijkamp et al., 2022; Wang et al., 2024b;a). By modeling the statistics of natural sequence variation, PLMs enable diverse applications including sequence-only, zero-shot mutation effect prediction (Meier et al., 2021) and protein structure prediction (Lin et al., 2022). In many real design workflows, however, the problem is inherently *edit-based*: engineers start from a natural scaffold and iteratively introduce substitutions and indels to modify loops, linkers, termini, or binding interfaces while preserving the overall fold and key functional sites.

Diffusion-based protein language models (Sohl-Dickstein et al., 2015; Austin et al., 2021; Hoogeboom et al., 2021; Zheng et al., 2023a; Sahoo et al., 2024; Shi et al., 2024; Nie et al., 2025) provide a powerful bidirectional modeling paradigm, but most existing methods use *masked diffusion*: the forward process hides tokens and generation recovers them. Masking is a loose proxy for evolution, which changes sequences through *substitutions and indels* and often requires variable-length edits for remodeling loops, linkers, and motifs. This mismatch makes it difficult to express realistic evolutionary trajectories and to post-edit existing proteins when length changes are needed.

We propose DPLM-EVO, an evolutionary discrete diffusion framework that predicts substitution, insertion, and deletion actions during denoising (Fig. 1). Our key design is to introduce a latent alignment space for diffusion, coupled to the observed sequence space, which makes indel-aware modeling tractable. We further introduce a contextual evolutionary noising kernel for substitution corruption to provide informative, biologically plausible mutation noise.

In short, our contributions are: (i) introducing an evolutionary discrete diffusion with explicit substitution, insertion, and deletion predictions; (ii) decoupling a latent alignment formulation for length-adaptive modeling; and (iii) biologically informed noising for substitution. Empirically, DPLM-EVO achieves state-of-the-art single-sequence mutation effect prediction on ProteinGym and supports variable-length generation and post-editing via explicit edit trajectories.

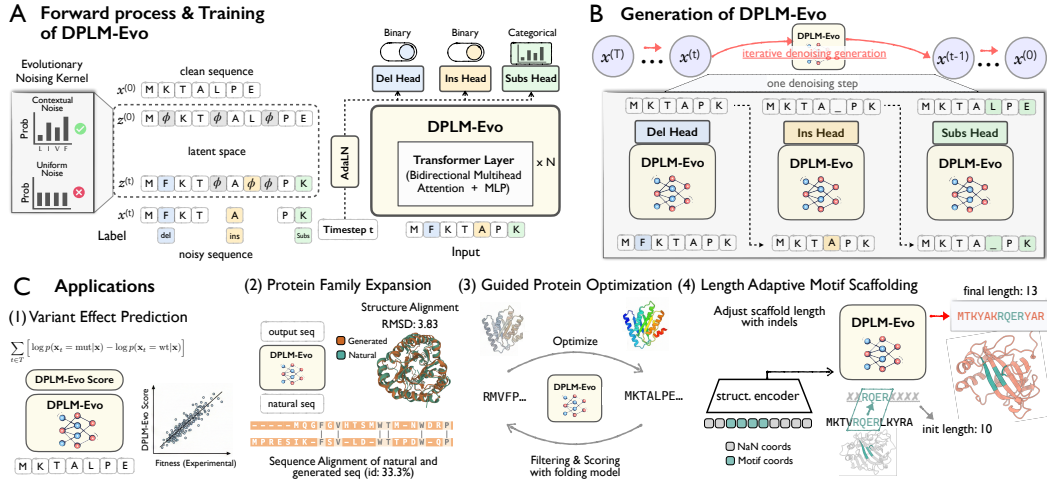


Figure 1: (A) DPLM-EVO models evolution with explicit substitution, insertion, and deletion predictions via a latent alignment space that supports variable length. (B) Sampling iteratively applies these edits to generate sequences. (C) Applications include mutation effect prediction, simulated evolution, optimization, and conditional motif scaffolding.

## 2 PRELIMINARIES

We briefly review diffusion protein language models (DPLMs) under the masked diffusion framework, which underpins our evolutionary extension.

**DPLM with masked diffusion.** The *masked* discrete diffusion framework (Sahoo et al., 2024; Liu et al., 2025; Nie et al., 2025) is characterized by a forward and backward Markov process. Let  $\text{Cat}(\mathbf{x}; \mathbf{p})$  be a categorical distribution on protein sequence  $\mathbf{x}$  parameterized by a probability vector  $\mathbf{p}$  over the vocabulary  $\mathcal{A}$  of  $K$  amino acids. The forward process of masked diffusion is governed by

$$q(\mathbf{x}^{(t)}|\mathbf{x}^{(t-1)}) = \text{Cat}(\mathbf{x}^{(t)}; \bar{\alpha}_t \delta_{\mathbf{x}^{(t-1)}} + (1 - \bar{\alpha}_t) \pi_{\text{mask}}),$$

which gradually perturb the data  $\mathbf{x}_0 \sim q(\mathbf{x})$  into an absorbing state  $\mathbf{x}^{(T)} \sim \pi_{\text{mask}}$ . The learned *backward* process  $p_\theta(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})$  reversely denoises  $\mathbf{x}^{(T)}$  towards the data distribution  $\mathbf{x}$ .

The learning objective can be simplified into weighted cross-entropies, resembling masked language modeling at arbitrary noise levels:

$$\begin{aligned} \mathcal{L}_t &= \mathbb{E}_{q(\mathbf{x})} \text{KL}[q(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)}, \mathbf{x}) \| p_\theta(\mathbf{x}^{(t-1)}|\mathbf{x}^{(t)})] \\ &= \mathbb{E}_{q(\mathbf{x})} \left[ \lambda_t \sum_{1 \leq i \leq L} \mathbb{I}_{x_t^{(i)} \neq x_0^{(i)}} \cdot \log p_\theta(x_0^{(i)}|\mathbf{x}^{(t)}) \right], \end{aligned}$$

where  $\lambda_t$  is a weighting coefficient induced by the noising schedule. For inference, DPLM samples from an all-mask initialization and iteratively denoises in a *mask-predict* manner.

**Motivation.** The masked diffusion formulation of makes it difficult to (i) represent the elementary evolutionary edits that biologists and engineers apply in practice—*substitutions, insertions, and deletions*—and (ii) support flexible, variable-length trajectories during generation and post-editing. To address this mismatch, we extend the DPLM framework to explicitly model evolutionary edit operations in next section.

## 3 METHODOLOGY

Standard discrete diffusion is typically fixed-length and therefore cannot explicitly represent insertions and deletions. We extend masked discrete diffusion to an evolutionary framework by coupling a variable-length *observed sequence space* to a *latent alignment space* (Havasi et al., 2025; Graves et al., 2006). The alignment is served as a latent variable for the forward noising process between observed clean and noisy sequence, enabling variable length modeling. The model predicts a substitution distribution together with insertion and deletion probabilities, thereby supporting explicit edit trajectories and variable-length sampling.

### 3.1 AN EVOLUTIONARY DISCRETE DIFFUSION FRAMEWORK

**Accommodating length-adaptive *indel* modeling with latent alignment.** Let  $\mathcal{V}$  represent the amino acid vocabulary. To handle variable-length sequences, we draw inspiration from latent-alignment methods (Graves et al., 2006; Havasi et al., 2025) and distinguish between two spaces:

- **Observed Space  $\mathcal{X}$** : Sequences  $\mathbf{x} \in \mathcal{V}^L$  with  $\mathcal{V} = \mathcal{A} \cup \{\mathbf{m}\}$ , where  $\mathbf{m}$  is a mask token.
- **Latent Alignment Space  $\mathcal{Z}$** : Sequences of length  $2L$  over  $\mathcal{V}^+ = \mathcal{V} \cup \{\phi\}$ , where  $\phi$  is a gap token.

We define a deterministic *collapse function*  $\Gamma^{-1}(\mathbf{z}) : \mathcal{Z} \rightarrow \mathcal{X}$  that maps a latent alignment  $\mathbf{z}$  to an observed sequence  $\mathbf{x}$  by removing all  $\phi$  tokens, i.e.,  $\Gamma^{-1}(\mathbf{z}) = [\mathbf{z}^{(j)} \mid \mathbf{z}^{(j)} \neq \phi]_{j=1}^N$ . Conversely,  $\Gamma(\mathbf{x})$  denotes the set of all latent alignments  $\mathbf{z}$  that collapse to  $\mathbf{x}$ , obtained by inserting exactly  $L$  gap tokens  $\phi$  into  $\mathbf{x}$  at arbitrary positions, for example  $[A, B, C] \mapsto [A, \phi, \phi, B, \phi, C]$ . This expanded latent canvas, which is strictly longer than observed  $\mathbf{x}$ , enables length-changing generation through explicit gap token transitions.

Given an observed protein  $\mathbf{x}_0$ , we treat its alignment  $\mathbf{z}_0$  as a latent variable, the training objective is to maximize the evidence lower bound (ELBO) of the log-likelihood:

$$\log p_\theta(\mathbf{x}_0) = \log \sum_{\mathbf{z}_0 \in \Gamma(\mathbf{x}_0)} p_\theta(\mathbf{z}_0) \geq \mathbb{E}_{\mathbf{z}_0 \in \Gamma(\mathbf{x}_0)} \left[ \mathbb{E}_{\mathbf{z}_t \sim q_t(\mathbf{z}_t | \mathbf{z}_0)} [\log p_\theta(\mathbf{z}_0 | \mathbf{z}_t)] \right]. \quad (1)$$

**Forward noising process for sequence with holistic edit operations.** Unlike predominant masked diffusion models that use absorbing-state (mask) noise, we introduce a new noising prior  $\pi(\mathbf{z}_0)$  that respects all possible sequence edit operations.

$$\begin{array}{ccccc} \mathbf{x}_0 & \xrightarrow{\Gamma} & \mathbf{z}_0 & \xrightarrow{\bar{\alpha}_t} & \mathbf{z}_t & \xrightarrow{\Gamma^{-1}} & \mathbf{x}_t \\ & & \mathbf{Q}_{\text{noise}} \downarrow & \nearrow 1-\alpha_t & & & \\ & & \pi(\mathbf{z}_0) & & & & \end{array}$$

Specifically, the forward transition  $q(\mathbf{z}_t | \mathbf{z}_0)$  is defined as an interpolant of the original data and the noise distribution:  $q_t(\mathbf{z}_t | \mathbf{z}_0) = \bar{\alpha}_t \delta_{\mathbf{z}_0} + (1 - \bar{\alpha}_t) \pi(\mathbf{z}_0)$ , where  $\bar{\alpha}_t$  is the noise schedule. The noise distribution  $\pi(\mathbf{z}_0)$  depends on the initial state  $\mathbf{z}_0$  via a transition matrix  $\mathbf{Q}_{\text{noise}} \in \mathbb{R}^{(K+2) \times (K+2)}$ , defining transitions between  $K$  amino acids,  $\mathbf{m}$  and  $\phi$ :

$$\pi(\mathbf{z}_0) = \text{Cat}(\cdot \mid \mathbf{Q}_{\text{noise}} \mathbf{z}_0).$$

To explicitly control the ratio of substitutions, insertions, and deletions, we parameterize  $\mathbf{Q}_{\text{noise}}$  with a deletion ratio  $\omega_{\text{del}}$  and an insertion ratio  $\omega_{\text{ins}}$ :

$$\mathbf{Q}_{\text{noise}} = \left( \begin{array}{c|c|c} z \in \mathcal{A} & z = \mathbf{m} & z = \phi \\ \hline (1 - \omega_{\text{del}})(1 - \rho_{\text{mask}}) \mathbf{U}_K & \mathbf{0}_K & \omega_{\text{ins}} \frac{1}{K} \mathbf{1}_K \\ (1 - \omega_{\text{del}}) \rho_{\text{mask}} \mathbf{1}_K^\top & 1 & 0 \\ \hline \omega_{\text{del}} \mathbf{1}_K^\top & 0 & 1 - \omega_{\text{ins}} \end{array} \right),$$

where  $\mathbf{1}_K$  denotes a  $K$ -dime column vector of all ones, and  $\mathbf{U}_K = \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top$  denotes the uniform transition matrix of size  $K \times K$  over amino acids. Intuitively, this implies:

- If  $z_0^{(j)} \in \mathcal{A}$  (amino acid): It is either substituted (with probability  $1 - \omega_{\text{del}}$ ) by another amino acid or mask token (up to  $\rho_{\text{mask}}$ , or deleted (becomes  $\phi$ , with probability  $\omega_{\text{del}}$ ).
- If  $z_0^{(j)} = \phi$  (gap): It either becomes an inserted amino acid, with probability  $\omega_{\text{ins}}$  or remains a gap (with probability  $1 - \omega_{\text{ins}}$ ).

**Simulating biological sequence mutations as evolutionary noising via contextualized on-policy substitution.** Uniform substitution noise ignores biophysical constraints and is often uninformative for learning evolutionary structure. We therefore use a *contextual evolutionary noising kernel* that draws substitution noise from the model’s on-policy predictions after a warmup stage, yielding more plausible and informative corruptions. We provide the details in §A.1.

**Connections to existing discrete diffusion.** By setting  $(\omega_{\text{del}}, \omega_{\text{ins}}, \rho_{\text{mask}})$  in  $\mathbf{Q}_{\text{noise}}$ , our model recovers masked diffusion and uniform diffusion as special cases, and enables initializing from pretrained masked diffusion models. We provide a discussion in Appendix A.2.

### 3.2 TRAINING

*Overall Objective.* We optimize a weighted sum of substitution, deletion, and insertion losses over observed tokens (collapsed from the latent alignment), with weights  $\lambda_{t-1}$  and  $\gamma$  controlling the denoising schedule and edit balance:

$$\mathcal{L}_t = \mathbb{E}_{\mathbf{x}_0, \mathbf{z}_0, \mathbf{z}_t} \left[ \sum_{k=1}^{|\Gamma^{-1}(\mathbf{z}_t)|} \lambda_{t-1} (\gamma_{\text{sub}} \mathcal{L}_{\text{sub}}^{(k)} + \gamma_{\text{del}} \mathcal{L}_{\text{del}}^{(k)} + \gamma_{\text{ins}} \mathcal{L}_{\text{ins}}^{(k)}) \right].$$

The introduced decomposed losses enables precise control over the model’s propensity for different evolutionary operations. We use three heads on the observed sequence  $\mathbf{x}_t = \Gamma^{-1}(\mathbf{z}_t)$ : a substitution head  $p_\theta^{\text{sub}}$  (multiclass over  $\mathcal{V}$ ), and deletion/insertion heads  $p_\theta^{\text{del}}, p_\theta^{\text{ins}}$  (binary). We provide the formulation of full decomposed objectives and the practical implementation in Appendix A.3.

### 3.3 GENERATION OF DPLM-EVO

DPLM-EVO follows standard iterative denoising for discrete diffusion, while enabling variable-length trajectories through explicit insertion and deletion actions in the observed sequence space. At each step, we apply deletion/insertion decisions from the indel heads, then fill substituted and newly inserted residues using the substitution head. Appendix B provides the full sampling procedure.

## 4 EXPERIMENTS

In this section, we evaluate DPLM-Evo across various understanding and generative tasks. We assess variant effect prediction to validate the model’s understanding of protein evolution, and evaluate generation capabilities through unconditional generation and motif scaffolding. Full experimental setups, additional analyses, and case studies are provided in Appendix C.

### 4.1 VARIANT EFFECT PREDICTION

**Setup.** Modeling the effect of sequence variation on function is fundamental for understanding and designing proteins.

DPLM-Evo predicts variant effects using protein sequence only, without supervision from experimental data. Unlike the common masked-residue scoring pipeline, DPLM-Evo is a substitution-based model that natively scores variants without masking by evaluating the substitution distribution at the mutated site(s). Please refer to Appendix C.1 for more details.

**Results.** On the ProteinGym DMS substitution zero-shot benchmark (Notin et al., 2023), DPLM-EVO achieves the highest correlation score among all the single sequence foundation models for variant effect prediction (Fig.2A). We attribute this strong correlation to the model’s evolutionary pretraining, which fundamentally enables it to learn mutation preferences from natural proteins, effectively capturing the constraints imposed by natural selection.

Moreover, explicitly aligning with evolutionary kernel further unlocks the potential of DPLM-EVO in mutation effect prediction. This makes the scores correlate more closely with natural mutations (Appendix C.1). Illustrated in Fig.2B, this alignment yields further enhancements, outperforming SaProt (Su et al., 2023) that takes additional structure and TrancepEVE L with supplementary MSA.

### 4.2 UNCONDITIONAL PROTEIN SEQUENCE GENERATION

DPLM-EVO performs unconditional generation by applying substitution, deletion, and insertion, enabling variable-length generation that more closely mirrors natural evolutionary trajectories. Full experimental settings can refer to Appendix C.2.

Fig.3 demonstrates the evaluation results of unconditional generation in various perspectives. (1) *Overall performance*: Fig. 3A shows DPLM-EVO achieves consistent high foldability across length as measured by ESMFold pLDDT. Fig.3B-C demonstrates the output lengths remain concentrated near their initial values without excessive expansion or collapse. This indicates that insertion and deletion prediction are invoked conservatively, resulting in a refinement process that prioritizes substitutions over drastic length changes. (2) *Diversity and Reduced Mode Collapse*: Compared with DPLM based on masked diffusion, Fig.3D-F shows DPLM-EVO achieves comparable foldability while possessing greater generation diversity and higher sequence entropy, indicating fewer repetition patterns and alleviating the mode collapse issue. (3) *Effect of Evolutionary Kernel*: Training with the contextual evolutionary noising kernel substantially outperforms uniform noising, as shown in Fig.3D. This indicates that biologically grounded corruptions encourages DPLM-EVO to learn more

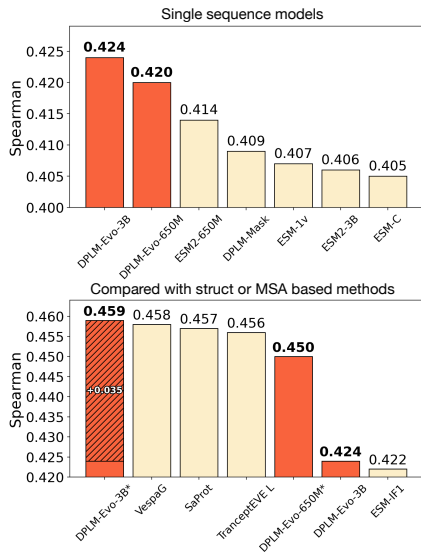


Figure 2: Results of ProteinGym zero-shot DMS substitution benchmark. \* represents the model which explicitly aligns the output distribution with MSA-based mutation effect prediction model.

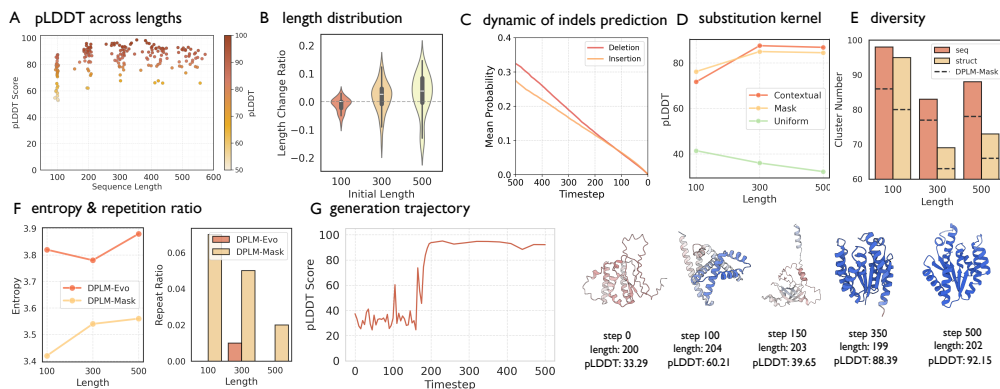


Figure 3: *Evaluation of unconditional sequence generation.* (A) Unconditional generation from length 100 to 500 evaluated by pLDDT. (B) Length distribution of DPLM-EVO generations from fixed initial lengths. (C) Insertion/deletion head predicted probability under different timesteps. (D) Ablation on different substitution kernels. (E) Sequence and structure diversity of DPLM-EVO compared with DPLM-Mask in different lengths. (F) Entropy and repetition comparison between DPLM-EVO and DPLM-Mask. (G) Demonstration of the generation trajectory.

evolutionarily plausible substitution predictions, yielding higher-quality samples at generation time. Further details and analyses about unconditional generation can be found in Appendix C.2.

### 4.3 LENGTH-ADAPTIVE SCAFFOLDING OF FUNCTIONAL MOTIFS

**Setup.** Motif scaffolding aims to generate a protein scaffold for a given functional motif.

We evaluate DPLM-EVO in *zero-shot* and *continued fine-tuning* settings. For finetuning, DPLM-EVO incorporates structural constraints for motif structure features, as illustrated in Fig.1C(4). Further details on settings and evaluation metrics are provided in the Appendix C.3 for more details.

**Results.** In the zero-shot setting, DPLM-EVO achieves higher overall success rate than EvoDiff and DPLM-Mask (Fig.4); we attribute this to the capability for dynamic scaffolding length adjustment and evolutionarily plausible mutations provided by the substitution head. In contrast, fixed-length sequence models require manually scaffold length enumeration and cannot revise length once an unsuitable initialization is chosen. Continued finetuning brings further improvements, highlighting the importance of multimodal conditioning (Appendix C.3).

### 4.4 CASE STUDIES

**In-silico sequence family expansion.** As shown in Fig. 5 and Appendix C.4, DPLM-EVO is able to generate diverse yet structurally similar relatives of a given protein via unconstrained post-editing. This suggest that DPLM-EVO performs unconstrained sequence optimization that preserves a shared structural scaffold while producing diverse sequences, with sequence identity mostly below 50%.

**Directed evolution of GFP.** Fig. 6 and Appendix C.5 demonstrates that DPLM-EVO improves the structural stability score (pTM) over iterations while keeping the chromophore site RMSD below a strict threshold, indicating it can leverage priors from evolution-scale sequences to optimize sequences under hard structural constraints.

## 5 CONCLUSION

In this work, we present DPLM-EVO, a unified framework for evolutionary discrete diffusion that explicitly models substitutions, insertions, and deletions. We decouple the latent space from observed sequences, enabling efficient indel-aware diffusion. We further enhance the substitution learning through a contextual evolutionary noising kernel. Extensive empirical evaluations demonstrate that DPLM-EVO not only achieves state-of-the-art performance in mutation effect prediction on ProteinGym but also opens new avenues for variable-length protein generation and optimization, bridging the gap between deep generative modeling and evolutionary biology.

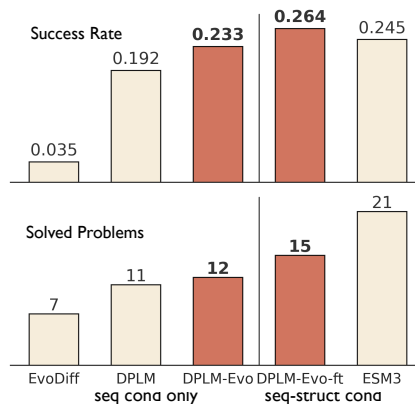


Figure 4: *Evaluation of motif scaffolding task.*

## REFERENCES

- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. In *Advances in Neural Information Processing Systems*, volume 34, pp. 17981–17993, 2021.
- Ethan Baron, Alan N Amin, Ruben Weitzman, Debora Marks, and Andrew Gordon Wilson. A diffusion model to shrink proteins while maintaining their function. *arXiv preprint arXiv:2511.07390*, 2025.
- Armin Behjati, Fatemeh Zare-Mirakabad, Seyed Shahriar Arab, and Abbas Nowzari-Dalini. Protein sequence profile prediction using protalbert transformer. *Computational Biology and Chemistry*, 99:107717, 2022.
- Nadav Brandes, Dan Ofer, Yam Peleg, Nadav Rappoport, and Michal Linial. Proteinbert: a universal deep-learning model of protein sequence and function. *Bioinformatics*, 38(8):2102–2110, 2022.
- Sebastian Deorowicz, Agnieszka Debudaj-Grabysz, and Adam Gudyś. Famsa: Fast and accurate multiple sequence alignment of huge protein families. *Scientific reports*, 6(1):33964, 2016.
- Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(10):7112–7127, 2021.
- ESM Team. Esm cambrian: Revealing the mysteries of proteins with unsupervised learning, 2024. URL <https://evolutionaryscale.ai/blog/esm-cambrian>.
- Shansan Gong, Shivam Agarwal, Yizhe Zhang, Jiacheng Ye, Lin Zheng, Mukai Li, Chenxin An, Peilin Zhao, Wei Bi, Jiawei Han, Hao Peng, and Lingpeng Kong. Scaling diffusion language models via adaptation from autoregressive models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=j1tSLYKwg8>.
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376, 2006.
- Jiatao Gu and Xiang Kong. Fully non-autoregressive neural machine translation: Tricks of the trade. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 120–133, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.11. URL <https://aclanthology.org/2021.findings-acl.11>.
- Jiatao Gu, Changhan Wang, and Junbo Zhao. Levenshtein transformer. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, volume 32, pp. 11179–11189, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/675f9820626f5bc0afb47b57890b466e-Abstract.html>.
- Yong Guo, Yin Zheng, Mingkui Tan, Qi Chen, Jian Chen, Peilin Zhao, and Junzhou Huang. Nat: Neural architecture transformer for accurate and compact architectures. *Advances in Neural Information Processing Systems*, 32, 2019.
- Marton Havasi, Brian Karrer, Itai Gat, and Ricky TQ Chen. Edit flows: Flow matching with edit operations. *arXiv preprint arXiv:2506.09018*, 2025.
- Tomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years of evolution with a language model. *bioRxiv*, pp. 2024–07, 2024.
- Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusionbert: Improving generative masked language models with diffusion models. 2023.

- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf>.
- Emiel Hoogetboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- Cheng-Yen Hsieh, Xinyou Wang, Daiheng Zhang, Dongyu Xue, Fei Ye, Shujian Huang, Zaixiang Zheng, and Quanquan Gu. Elucidating the design space of multimodal protein language models. In *Forty-second International Conference on Machine Learning*, 2025.
- 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. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 8946–8970. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/hsu22a.html>.
- Elodie Laine, Yasaman Karami, and Alessandra Carbone. Gemme: a simple and fast global epistatic model predicting mutational effects. *Molecular biology and evolution*, 36(11):2604–2619, 2019.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *BioRxiv*, 2022.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- Yangzhou Liu, Yue Cao, Hao Li, Gen Luo, Zhe Chen, Weiyun Wang, Xiaobo Liang, Biqing Qi, Lijun Wu, Changyao Tian, et al. Sequential diffusion language models. *arXiv preprint arXiv:2509.24007*, 2025.
- Céline Marquet, Julius Schlenzok, Marina Abakarova, Burkhard Rost, and Elodie Laine. Expert-guided protein language models enable accurate and blazingly fast fitness prediction. *Bioinformatics*, 40(11):btac621, 11 2024. ISSN 1367-4811. doi: 10.1093/bioinformatics/btac621. URL <https://doi.org/10.1093/bioinformatics/btac621>.
- Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language models enable zero-shot prediction of the effects of mutations on protein function. In *Advances in Neural Information Processing Systems*, pp. 29287–29303, 2021.
- Ananthan Nambiar, Maeve Heflin, Simon Liu, Sergei Maslov, Mark Hopkins, and Anna Ritz. Transforming the language of life: transformer neural networks for protein prediction tasks. In *Proceedings of the 11th ACM international conference on bioinformatics, computational biology and health informatics*, pp. 1–8, 2020.
- Shen Nie, Fengqi Zhu, Chao Du, Tianyu Pang, Qian Liu, Guangtao Zeng, Min Lin, and Chongxuan Li. Scaling up masked diffusion models on text. In *The Thirteenth International Conference on Learning Representations*, 2024.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025.
- Erik Nijkamp, Jeffrey Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2: exploring the boundaries of protein language models. *arXiv preprint arXiv:2206.13517*, 2022.
- Pascal Notin, Aaron W Kollasch, Daniel Ritter, Lood Van Niekerk, Steffan Paul, Han Spinner, Nathan J Rollins, Ada Shaw, Ruben Weitzman, Jonathan Frazer, et al. Proteingym: Large-scale benchmarks for protein fitness prediction and design. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

- 378 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li.  
379 Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In  
380 *The Thirteenth International Conference on Learning Representations*, 2024.  
381
- 382 Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel,  
383 and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information*  
384 *processing systems*, 32, 2019.
- 385 Machel Reid, Vincent Josua Hellendoorn, and Graham Neubig. Diffuser: Diffusion via edit-based  
386 reconstruction. In *International Conference on Learning Representations*, 2022.  
387
- 388 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,  
389 Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function  
390 emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS*, 2019. doi:  
391 10.1101/622803. URL <https://www.biorxiv.org/content/10.1101/622803v4>.
- 392 Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu,  
393 Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language  
394 models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2024.  
395
- 396 Subham Sekhar Sahoo, Justin Deschenaux, Aaron Gokaslan, Guanghan Wang, Justin T Chiu, and  
397 Volodymyr Kuleshov. The diffusion duality. In *Forty-second International Conference on Machine*  
398 *Learning*.
- 399 Yair Schiff, Subham Sekhar Sahoo, Hao Phung, Guanghan Wang, Sam Boshar, Hugo Dalla-torre,  
400 Bernardo P de Almeida, Alexander M Rush, Thomas PIERROT, and Volodymyr Kuleshov. Simple  
401 guidance mechanisms for discrete diffusion models. In *The Thirteenth International Conference*  
402 *on Learning Representations*.  
403
- 404 Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and generalized  
405 masked diffusion for discrete data. *Advances in neural information processing systems*, 37:  
406 103131–103167, 2024.
- 407 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsu-  
408 pervised learning using nonequilibrium thermodynamics. In Francis Bach and David Blei  
409 (eds.), *International Conference on Machine Learning*, volume 37 of *Proceedings of Machine*  
410 *Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015. PMLR, PMLR. URL  
411 <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- 412 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Interna-*  
413 *tional Conference on Learning Representations*, 2020a.  
414
- 415 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
416 Poole. Score-based generative modeling through stochastic differential equations. In *International*  
417 *Conference on Learning Representations*, 2020b.
- 418 Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein  
419 language modeling with structure-aware vocabulary. *bioRxiv*, pp. 2023–10, 2023.  
420
- 421 Robert Verkuil, Ori Kabeli, Yilun Du, Basile IM Wicky, Lukas F Milles, Justas Dauparas, David  
422 Baker, Sergey Ovchinnikov, Tom Sercu, and Alexander Rives. Language models generalize beyond  
423 natural proteins. *bioRxiv*, pp. 2022–12, 2022.  
424
- 425 Dimitri von Rütte, Janis Fluri, Yuhui Ding, Antonio Orvieto, Bernhard Schölkopf, and Thomas  
426 Hofmann. Generalized interpolating discrete diffusion. In *Forty-second International Conference*  
427 *on Machine Learning*, 2025.
- 428 Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Dplm-2: A  
429 multimodal diffusion protein language model. *arXiv preprint arXiv:2410.13782*, 2024a.  
430
- 431 Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Diffusion  
language models are versatile protein learners. *arXiv preprint arXiv:2402.18567*, 2024b.

- 432 Zirui Wu, Lin Zheng, Zhihui Xie, Jiacheng Ye, Jiahui Gao, Yansong Feng, Zhenguo Li, Victoria W.,  
433 Guorui Zhou, and Lingpeng Kong. Dreamon: Diffusion language models for code infilling beyond  
434 fixed-size canvas, 2025. URL <https://hkunlp.github.io/blog/2025/dreamon>.
- 435 Yijia Xiao, Jiezhong Qiu, Ziang Li, Chang-Yu Hsieh, and Jie Tang. Modeling protein using large-scale  
436 pretrain language model. *arXiv preprint arXiv:2108.07435*, 2021.
- 437 Kevin K Yang, Alex X Lu, and Nicolo Fusi. Convolutions are competitive with transformers for  
438 protein sequence pretraining. *bioRxiv*, pp. 2022–05, 2022.
- 439 Jiacheng Ye, Zhihui Xie, Lin Zheng, Jiahui Gao, Zirui Wu, Xin Jiang, Zhenguo Li, and Lingpeng  
440 Kong. Dream 7b: Diffusion large language models. *arXiv preprint arXiv:2508.15487*, 2025.
- 441 Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Quanquan Gu. Diffusion language models  
442 can perform many tasks with scaling and instruction-finetuning. *arXiv preprint arXiv:2308.12219*,  
443 2023a.
- 444 Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Mingxuan Wang. Dinoiser: Diffused  
445 conditional sequence learning by manipulating noises. *arXiv preprint arXiv:2302.10025*, 2023b.
- 446 Lin Zheng, Jianbo Yuan, Lei Yu, and Lingpeng Kong. A reparameterized discrete diffusion model for  
447 text generation. *arXiv preprint arXiv:2302.05737*, 2023a.
- 448 Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei YE, and Quanquan Gu. Structure-informed  
449 language models are protein designers. In *International Conference on Machine Learning*, 2023b.

## 454 A METHODOLOGY DETAILS.

### 455 A.1 SUBSTITUTION LEARNING WITH CONTEXTUAL EVOLUTIONARY NOISE

456 In this section, we discuss the noise kernel for the substitution head. Standard multinomial diffusion  
457 models typically employ un-informative uniform noise for corruption, which ignores the biophysical  
458 constraints of the protein fitness landscape and leads to inefficient training. To address this, we  
459 propose a *contextual evolutionary noising kernel* that instantiates the substitution noise sampled  
460 from a contextualized distribution. The contribution is obtained using the model’s own on-policy  
461 denoising predictions, where the noise position is first noised with  $\mathbf{Q}_{\text{noise}}$ , then predicted. After a  
462 warmup stage trained with the original non-informative  $\mathbf{Q}_{\text{noise}}$ , this model’s on-policy substitution  
463 prediction provides biologically plausible and evolutionarily reasonable mutation as noise, which is  
464 more informative than random noise. Meanwhile, it encourages the model to capture evolutionary  
465 and homologous dependencies between amino acids and sequences. We will introduce this in detail  
466 in the following part.

#### 467 A.1.1 MOTIVATION: FROM UNIFORM NOISE TO BIOLOGICAL MANIFOLDS

468 Standard discrete diffusion models typically employ an uniform noise kernel (Hoogeboom et al.,  
469 2021). In the context of protein engineering, this implies that a mutation from a hydrophobic  
470 residue (e.g., Leucine) to a charged residue (e.g., Arginine) is as probable as a mutation to another  
471 hydrophobic residue (e.g., Isoleucine). This assumption fundamentally contradicts the biophysical  
472 constraints of the protein fitness landscape. Training with uniform noise suffers from significant  
473 inefficiency: the model spends a large portion of training correcting biologically obvious errors  
474 (e.g., restoring a hydrophobic core disrupted by charged noise) rather than learning evolutionary  
475 dependencies. To address this, we propose a **contextual evolutionary noising kernel**. Instead of  
476 corrupting data towards random chaos, we utilize the model’s own prediction capability to generate  
477 noise that lies on the protein manifold. This not only provides more informative noisy tokens, which  
478 is helpful for denoising, but also encourages the model to capture the dependencies between the  
479 wild-type sequence and plausible homologous variants generated by model prediction.

#### 480 A.1.2 FORMALIZATION: CONSTRUCT THE CONFIDENCE-AWARE KERNEL WITH MASK PREDICTION

481 We leverage mask token  $\mathbf{m}$  (distinct from the gap token  $\phi$  used for deletions) to represent unknown  
482 semantic identity in the protein sequence. Let  $p_\theta$  denote the neural network parameterized by  $\theta$ . At  
483 time step  $t$ , we construct a contextual distribution for an amino acid at index  $j$  in latent alignment  
484 sequence  $\mathbf{z}$  (where  $\mathbf{z}_0^{(j)} \in \mathcal{V}$ ) by masking the input at that position:

$$485 p_{\text{ctx}}(\cdot | \mathbf{z}_0^{(j)}) = \text{Softmax} \left( p_\theta(\Gamma^{-1}(\mathbf{z}_0^{\setminus j} \cup \{\mathbf{m}\}))^{(j)} \right) \quad (2)$$

where  $\mathbf{z}_0^{\setminus j} \cup \{\mathbf{m}\}$  represents the latent sequence with the  $j$ -th token replaced by  $\mathbf{m}$ . Specifically, contextual noise is sampled from  $p_{\text{ctx}}$ , then the model is trained to denoise based on contextual noise.

### A.1.3 IMPLEMENTATION: ON-POLICY CONFIDENCE-AWARE KERNEL FOR MORE EFFECTIVE TRAINING AND GENERATION

We leverage  $p_{\theta}^{\text{subs}}$  of DPLM-EVO for the mask prediction to construct  $p_{\text{ctx}}$ . However, during training, the model is trained to denoise based on the contextual noise, which consists of the standard amino acids. As training progresses, DPLM-EVO will lose its mask prediction capabilities. To prevent catastrophic forgetting of the masking prediction, we construct a mixing noising kernel for the contextual noise and the mask state. Therefore, DPLM-EVO can learn the mask prediction explicitly during training. Specifically, we employ a *confidence-aware* gating mechanism to construct the noising kernel:

$$\pi(v | z_0^{(j)}) = \mathbb{I}(c_j > \tau) \cdot p_{\text{ctx}}(v | z_0^{(j)}) + \mathbb{I}(c_j \leq \tau) \cdot \delta_{\mathbf{m}}(v) \quad (3)$$

Here,  $c_j = \max_{v \in \mathcal{V}} p_{\text{ctx}}(v | z_0^{(j)})$  is the confidence score of the model’s prediction, and  $\tau$  is a dynamic threshold. If  $c_j \leq \tau$ , this indicates the low confidence prediction that the model is uncertain, which represents insufficiently valuable or evolutionarily relevant information. Therefore, we fallback to the mask token  $\mathbf{m}$  for these positions. This reinforces the fundamental masked prediction objective, ensuring the model remains robust to missing information and avoid forgetting about mask prediction. Crucially, this process is **on-policy**: the noise is generated by the model state  $\theta$  at the current training step. To enhance the quality of contextual noise at the early training stage and prevent the training instability, we initialize the model parameters from a pre-trained MLM-based pLM or an absorbing discrete diffusion-based pLM.

### A.1.4 BIOLOGICAL INTERPRETATION: TRAVERSING THE FITNESS LANDSCAPE

By replacing the static uniform kernel with our dynamic contextual kernel, we reframe the diffusion training process as a traversal on the fitness landscape:

**(1) Denoising as Error Correction (Lethal Mutations):** When the contextual noise  $p_{\text{ctx}}$  generates a residue that violates structural constraints (e.g., steric clashes), the training objective is to enable the model to identify and correct these erroneous mutations that are evolutionarily unacceptable.

**(2) Denoising as Homology Modeling (Neutral Mutations):** The contextual noise sampled from  $p_{\text{ctx}}$  may also be biologically valid substitutions. In this regime, the ground truth  $\mathbf{x}_0$  represents a specific instance (wild type), while the noise  $\mathbf{x}_t$  represents a high-fitness neighbor (homologue). The denoising loss encourages the model to learn the *evolutionary dependencies* between the original functional sites  $\mathbf{x}_0$  and variable sites  $\mathbf{x}_t$ .

### A.1.5 THE LEARNABLE PRIOR

We discuss the sampling prior of the contextual noising kernel when  $t = T$ . In standard multinomial diffusion, the sampling prior is a uniform distribution  $\mathcal{U}(\mathcal{V})$ , which is a poor approximation of natural proteins. In our framework, the effective prior at  $T$  steps becomes the model’s prediction given a fully masked sequence:

$$p(z_T) \approx p_{\theta}(\cdot | \mathbf{m}^L) \quad (4)$$

This **learnable prior** captures the natural background frequencies of amino acids and global sequence statistics (e.g., length distributions and domain compositions) inherent in the pre-trained weights. Consequently, the reverse generation process initializes from a biologically informed distribution rather than random chaos, improving sampling efficiency and stability.

### A.1.6 OTHER ALTERNATIVE NOISING KERNEL: BLOSUM-INFORMED SUBSTITUTION

While the contextual evolutionary noising kernel provides dynamic, instance-specific noise, it incurs additional computational overhead due to the required forward pass of the model. For scenarios requiring high computational efficiency, we propose a static but biologically grounded alternative based on the BLOSUM substitution matrices.

Standard discrete diffusion typically uses a uniform transition matrix  $\mathbf{U}_K$ , which implies that all amino acid substitutions are equally probable. In contrast, the BLOSUM62 matrix encodes empirical substitution frequencies observed in homologous protein alignments. Let  $\mathbf{B} \in \mathbb{R}^{K \times K}$  be the BLOSUM62 scoring matrix, where  $\mathbf{B}_{ij}$  represents the log-odds score of substituting amino acid  $i$  with  $j$ .

We construct the static substitution noise matrix  $\mathbf{M}_{\text{BLOSUM}}$  by applying a row-wise Softmax operation over the scaled scores:

$$[\mathbf{M}_{\text{BLOSUM}}]_{ij} = \frac{\exp(\mathbf{B}_{ij}/\tau)}{\sum_{k=1}^K \exp(\mathbf{B}_{ik}/\tau)} \quad (5)$$

where  $\tau > 0$  is a temperature hyperparameter that controls the entropy of the noise distribution:

- As  $\tau \rightarrow \infty$ , the distribution approaches the uniform distribution  $\mathbf{U}_K$ .
- As  $\tau \rightarrow 0$ , the distribution collapses to the identity matrix (no mutation).
- At moderate  $\tau$ , the distribution favors physico-chemically conservative mutations (e.g.,  $L \leftrightarrow I$ ) over radical changes (e.g.,  $L \leftrightarrow K$ ).

This static kernel can be directly plugged into our proposed evolutionary discrete diffusion framework by replacing the uniform component in the noise transition matrix. Although less expressive than the contextual noising kernel, it still has large potential to outperform uniform noise by respecting fundamental biochemical properties.

## A.2 CONNECTIONS TO EXISTING DISCRETE DIFFUSION

Our model can recover existing discrete diffusion models by manipulating coefficients in  $\mathbf{Q}_{\text{noise}}$ . For example, when  $\omega_{\text{del}} = 0$  and  $\omega_{\text{ins}} = 0$ , we disable indels and our model becomes a fixed-length sequence diffusion model. In such circumstances,

- if  $\rho_{\text{mask}} = 1$ ,  $\mathbf{Q}_{\text{noise}}$  will always transits any token into  $\mathbf{m}$  and our model reduces to classical masked diffusion (Sahoo et al., 2024; Shi et al., 2024).
- if  $\rho_{\text{mask}} = 0$ ,  $\mathbf{Q}_{\text{noise}}$  will transits each token into another random token, and our model reduces to classical uniform diffusion (Austin et al., 2021; Schiff et al.).
- if  $\rho_{\text{mask}} \in (0, 1)$ ,  $\mathbf{Q}_{\text{noise}}$  will transits each token into either a random token or  $\mathbf{m}$ , and our model reduces to a generalized diffusion with mixed masked-uniform noising paths (Austin et al., 2021; von Rütte et al., 2025).

With these connections, we can also initialize our model from pretrained masked diffusion-based models, efficiently reprogramming classical discrete diffusion to enable the full spectrum of sequence edit operations. We note that there are also recent work on variable-length diffusion/flow models (Havasi et al., 2025; Wu et al., 2025) for text generation, and Baron et al. (2025) learning to shrink protein sequences. We extend our discussion to other related work in §D.

## A.3 DECOMPOSED OBJECTIVES AND PRACTICAL IMPLEMENTATIONS

**The decomposed training objectives.** To make the training tractable, we should solve a critical issue: the diffusion is defined on the latent sequence  $\mathbf{z}_t$ , but in practice the neural network  $f_\theta$  operates on the original sequence  $\mathbf{x}_t = \Gamma^{-1}(\mathbf{z}_t)$ , which is collapsed by  $\mathbf{z}_t$ . To bridge this gap, we define the *Index Mapping Function*  $\mathcal{I} : \{1, \dots, L_t\} \rightarrow \{1, \dots, N\}$  such that  $\mathcal{I}(k)$  is the index of the  $k$ -th non-gap token in the latent sequence  $\mathbf{z}_t$ . Then, we decompose the loss defined on the latent sequence  $\mathbf{z}$  into three mutually exclusive components defined in the observed space  $\mathbf{x}$ , i.e., substitution loss, deletion loss and insertion loss, based on the token category between  $\mathbf{z}_t$  (current noisy state) and  $\mathbf{z}_0$  (ground truth).

To more clearly decouple the prediction of the three operations, we leverage separate and independent heads,  $p_\theta^{\text{sub}}$ ,  $p_\theta^{\text{del}}$  and  $p_\theta^{\text{ins}}$  for the substitution, deletion and insertion prediction for each token in the original sequence  $\mathbf{x}_t$ . We define the loss for the  $k$ -th token of the input sequence  $\mathbf{x}_t$ :

(1). *Substitution Loss.* It is active only when the input and target token are both valid amino acids:

$$\mathcal{L}_{\text{sub}}^{(k)} = \mathbb{I}_{(\mathbf{z}_0^{(\mathcal{I}(k))} \in \mathcal{V})} \cdot \mathbb{I}_{(\mathbf{z}_t^{(\mathcal{I}(k))} \in \mathcal{V})} \cdot \mathbb{I}_{(\mathbf{z}_0^{(\mathcal{I}(k))} \neq \mathbf{z}_t^{(\mathcal{I}(k))})} \cdot \text{CE} \left( \mathbf{z}_0^{(\mathcal{I}(k))}, p_\theta^{\text{sub}}(\cdot | \mathbf{x}_t) \right).$$

(2). *Deletion Loss.* It encourages the model to predict  $\phi$  if the current token is noise when its target is a gap in  $\mathbf{z}_0$ :

$$\tilde{\mathcal{L}}_{\text{del}}^{(k)} = \mathbb{I}_{(\mathbf{z}_0^{(\mathcal{I}(k))} = \phi)} \cdot \mathbb{I}_{(\mathbf{z}_t^{(\mathcal{I}(k))} \in \mathcal{V})} \cdot \text{CE} \left( \mathbf{z}_0^{(\mathcal{I}(k))}, p_\theta^{\text{del}}(\cdot | \mathbf{x}_t) \right).$$

(3). *Insertion Loss.* Let  $v_{\text{next}}^{(k)}$  be the first non-gap token in  $\mathbf{z}_0$  between indices  $\mathcal{I}(k)$  and  $\mathcal{I}(k+1)$ . If no such token exists, i.e., there is no insertion needed between  $\mathbf{x}_t^k$  and  $\mathbf{x}_t^{k+1}$ , the  $v_{\text{next}}^{(k)}$  is  $\emptyset$ . The loss is calculated on the positions that need insertion for reconstruction:

$$\tilde{\mathcal{L}}_{\text{ins}}^{(k)} = \mathbb{I}_{(v_{\text{next}}^{(k)} \neq \emptyset)} \cdot \text{CE} \left( v_{\text{next}}^{(k)}, p_\theta^{\text{ins}}(\cdot | \mathbf{x}_t) \right).$$

**Practical considerations for  $\mathcal{L}_{\text{del}}$  and  $\mathcal{L}_{\text{ins}}$ .** In our preliminary experiments, we find that training with the original  $\mathcal{L}_{\text{del}}$  imposes a significant risk of mode collapse, while  $\mathcal{L}_{\text{ins}}$  leads to unstable training.

We find that only training  $\mathcal{L}_{\text{del}}$  to predict the  $\phi$  token poses a significant risk of mode collapse. Without exposure to tokens that should *not* be deleted, the head may converge to a degenerate solution that always predicts the gap token  $\phi$  for any input, resulting in excessive deletion during generation.

Therefore, we introduce negative samples that represents tokens that should be preserved for deletion training. Given that there are only two distinct prediction targets: either a gap token or retaining the original token, deletion is essentially a *binary classification* task. Therefore, we parameterize the deletion head as a binary prediction head, and define the binary deletion target  $y_k^{\text{del}} = \mathbb{I}(\mathbf{z}_0^{(\mathcal{I}^{(k)})} = \phi)$ . The  $\mathcal{L}_{\text{del}}$  can be formalized as the Binary Cross-Entropy (BCE) computed over all the positions:

$$\mathcal{L}_{\text{del}} = \sum_k^{|\Gamma^{(-1)}(\mathbf{z}_t)|} \text{BCE}(y_k^{\text{del}}, p_{\theta}^{\text{del}}(\cdot|\mathbf{x}_t)) \quad (6)$$

where  $\text{BCE}(y, p) = -[y \log p + (1 - y) \log(1 - p)]$ .

Meanwhile, we find that parameterize the insertion head as a binary classifier leads to more stable training. For each token position  $k$  in the observed sequence  $\mathbf{x}_t$ , the head predicts whether an additional token should be inserted immediately to the right of  $\mathbf{x}_t^{(k)}$ . We define the binary insertion target as  $y_k^{\text{ins}} = \mathbb{I}(v_{\text{next}}^{(k)} \neq \emptyset)$ , and train the head with a BCE objective:

$$\mathcal{L}_{\text{ins}} = \sum_k^{|\Gamma^{(-1)}(\mathbf{z}_t)|} \text{BCE}(y_k^{\text{ins}}, p_{\theta}^{\text{ins}}(\cdot|\mathbf{x}_t)). \quad (7)$$

When an insertion is triggered, we first insert a special mask token  $\mathbf{m}$  as a noisy placeholder at that location, and then reuse the substitution head for filling the masked position.

## B SAMPLING DETAILS

### B.1 OVERVIEW

The generation process of DPLM-EVO follows the standard iterative denoising paradigm of discrete diffusion models. Starting from a completely noisy sequence  $\mathbf{x}_T$ , which is initialized with mask tokens with length  $L$ , and sampled from the learned prior  $p_{\theta}(\cdot|\mathbf{m}^L)$ , the model iteratively refines the sequence by sampling from the reverse distribution  $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ . Formally, the reverse transition is derived by marginalizing over the predicted ground truth  $\hat{\mathbf{x}}_0$ :

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \sum_{\hat{\mathbf{x}}_0} q(\mathbf{x}_{t-1}|\mathbf{x}_t, \hat{\mathbf{x}}_0) p_{\theta}(\hat{\mathbf{x}}_0|\mathbf{x}_t) \quad (8)$$

In practice, since summing over the entire sequence space is intractable, we approximate this process that we first sample a  $\hat{\mathbf{x}}_0$  via  $p_{\theta}(\cdot|\mathbf{x}_t)$ , and then samples the previous state  $\mathbf{x}_{t-1}$  using the tractable posterior  $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \hat{\mathbf{x}}_0)$ . According to Zheng et al. (2023a), this can be implemented with a predict-then-renoise strategy: first update a set of indices in  $\mathbf{x}_t$  with  $\hat{\mathbf{x}}_0$ , then renoise the indices with low confidence from the noise distribution.

### B.2 EVOLUTIONARY SAMPLING WITH DELETION, INSERTION AND SUBSTITUTION

We maintain a noisy state  $\mathcal{N}_t$  during sampling, which tracks the indices of tokens that are considered "noisy" at the current step  $t$  and will be updated for the next step  $t - 1$ . The denoising process at each iteration step  $t$  proceeds through four steps.

**Step 1: Deletion prediction.** We first modify the sequence  $\mathbf{x}_t$  by removing noise. The deletion head is applied to the indices in the noisy set  $j \in \mathcal{N}_t$ . If the model predicts deletion (i.e.,  $p_{\theta}^{\text{del}}(\mathbf{x}_t^j) > \tau_{\text{del}}$ ), the token is removed from the sequence. The noisy set  $\mathcal{N}_t$  is updated to reflect the shifted indices of the remaining tokens.

**Step 2: Insertion prediction.** The insertion head scans the current noisy set. If an insertion is predicted at index  $j$  (i.e.,  $p_{\theta}^{\text{ins}}(\mathbf{x}_t^j) > \tau_{\text{ins}}$ ), a mask token  $\mathbf{m}$  is inserted into the sequence at position  $j + 1$ . Crucially, since these new tokens lack semantic content, their indices are immediately added to the noisy set  $\mathcal{N}_t$ .

**Step 3: Substitution prediction, along with the insertion content.** The substitution head makes prediction for all tokens, yielding a  $\hat{\mathbf{x}}_0$  sampled from  $p_{\theta}^{\text{subs}}(\cdot|\mathbf{x}_t)$  and the corresponding confidence score. We update the indices in the noisy set  $\mathcal{N}_t$  of  $\mathbf{x}_t$  with  $\hat{\mathbf{x}}_0$ , including the substitution prediction and inserted content prediction. Then, we update the noisy set for the next step, i.e.,  $\mathcal{N}_{t-1}$ , by selecting

**Algorithm 1** Evolutionary sampling with DPLM-Evo

---

**Require:** Trained prediction heads  $p_\theta^{\text{del}}$ ,  $p_\theta^{\text{ins}}$  and  $p_\theta^{\text{subs}}$ , Prior length  $L_{\text{init}}$ , Steps  $T$

- 1: **Initialize:**  $x_T \sim p_\theta(\cdot | \mu^{L_{\text{init}}})$ ,  $\mathcal{N}_T \leftarrow \{1, \dots, L_{\text{init}}\}$
- 2: **for**  $t = T$  **down to** 1 **do**
- 3:     // Step 1: Deletion
- 4:     Predict deletion:  $\mathcal{D} \leftarrow \{j \in \mathcal{N} \mid p_\theta^{\text{del}}(x_t^j) > \tau_{\text{del}}\}$
- 5:      $x_t \leftarrow \text{Delete}(x_t, \mathcal{D})$ , Update indices in  $\mathcal{N}_t$
- 6:     // Step 2: Insertion
- 7:     Predict insertion:  $\mathcal{I} \leftarrow \{j \in \mathcal{N} \mid p_\theta^{\text{ins}}(x_t^j) > \tau_{\text{ins}}\}$
- 8:      $x_t \leftarrow \text{Insert}(x_t, \mathcal{I}, \text{token} = \mathbf{m})$
- 9:      $\mathcal{N}_t \leftarrow \mathcal{N}_t \cup \text{Indices}(\text{InsertedToken})$
- 10:    // Step 3: Substitution
- 11:    Sample  $\hat{x}_0 \sim p_\theta^{\text{subs}}(\cdot | x_t)$
- 12:    **for**  $j$  in  $\mathcal{N}_t$  **do**
- 13:       $x_t^{(j)} \leftarrow \hat{x}_0^{(j)}$
- 14:    **end for**
- 15:     $k_t \leftarrow \text{Schedule}(t)$
- 16:     $\mathcal{N}_{t-1} \leftarrow \text{TopKLowestConfidence}(c, k_t)$
- 17:    // Step 4: Renoise
- 18:    Sample  $x_{t-1} \sim \pi_{\text{noise}}(\cdot | x_t)$  for indices in  $\mathcal{N}_{t-1}$ .
- 19: **end for**
- 20: **Return:**  $x_0$

---

the  $k_t\%$  tokens with the *lowest* confidence scores, where  $k_t\%$  follows a linear decay schedule from 100% to 0%.

**Step 4: Renoising.** Finally, we perform renoising for the indices in  $\mathcal{N}_{t-1}$ . This also prevents the model from collapsing into local optima. For every index  $j \in \mathcal{N}_{t-1}$ , we sample  $\mathbf{x}_{t-1}^{(j)}$  from the noise distribution. This noise distribution can be instantiated as the contextual evolutionary noising kernel (using the model’s own predictions) or the BLOSUM-based kernel, ensuring that the noise state remains biologically plausible, which is consistent with the training.

The full procedure is summarized in Algorithm 1.

## C FULL EXPERIMENTAL RESULTS

In this section, we evaluate DPLM-Evo across various understanding and generative tasks. First, we assess variant effect prediction to validate the model’s understanding of protein evolution. Subsequently, we examine the model’s generation capabilities, including unconditional generation (covering both substitution-only and full edit operations) and conditional motif scaffolding scenario. Finally, we demonstrate the potential application of DPLM-Evo in protein sequence optimization.

### C.1 VARIANT EFFECT PREDICTION

**Setup.** Modeling the effect of sequence variation on function is fundamental for understanding and designing proteins. DPLM-Evo predicts variant effects using protein sequence only, without supervision from experimental data. **Unlike** the common masked-residue scoring pipeline used by masked language models (i.e., masking the residue(s) of interest and reading out the logits at the masked positions), DPLM-Evo is a substitution-based model that natively scores variants **without masking**. Instead, we directly input the wild-type sequence and evaluate the model’s substitution distribution at the mutated site(s), which better matches the model design and avoids introducing an artificial mask token.

For a variant with mutation set  $T$ , similar to ESM-1v (Meier et al., 2021), we use a log-odds mutation score that compares the mutant residue to the wild-type residue:

$$\sum_{t \in T} \left[ \log p(\mathbf{x}_t = \text{mut} \mid \mathbf{x}) - \log p(\mathbf{x}_t = \text{wt} \mid \mathbf{x}) \right],$$

where  $\mathbf{x}$  denotes the wild-type sequence and  $p(\mathbf{x}_t = \cdot \mid \mathbf{x})$  is the substitution probability at position  $t$  predicted by DPLM-Evo conditioned on the unmodified wild-type context. In contrast, masked-

702 residue approaches typically score mutations via

$$703 \sum_{t \in T} \left[ \log p(\mathbf{x}_t = \text{mut} \mid \mathbf{x}_{\setminus T}) - \log p(\mathbf{x}_t = \text{wt} \mid \mathbf{x}_{\setminus T}) \right],$$

704 where  $\mathbf{x}_{\setminus T}$  indicates that the mutated positions are masked/removed from the input; DPLM-Evo does  
705 not require this step.

706 **Results.** We evaluate the performance on the ProteinGym DMS substitution zero-shot bench-  
707 mark (Notin et al., 2023) by calculating the correlation between DPLM-EVO’s score and experimental  
708 fitness score across all 217 DMS assays. As shown in Fig.2 upper figure, **(1) DPLM-EVO achieves**  
709 **the highest correlation score among all the single sequence foundation models for variant effect**  
710 **prediction in ProteinGym.** DPLM-EVO outperforms the ESM model series, including ESM-2 (Lin  
711 et al., 2022), ESM-C (ESM Team, 2024), ESM-1v (Meier et al., 2021), and DPLM (Wang et al.,  
712 2024b). According to Fig.2(b), DPLM-EVO even surpasses ESM-IF1 (Hsu et al., 2022), which  
713 utilizes additional structure information, despite DPLM-EVO using a single sequence. Crucially,  
714 we observe that scaling up the model leads to further improvements, as evidenced by the 3B model  
715 outperforming the 650M model. This scalability stands in contrast to ESM-2, which exhibits a  
716 performance regression as model size increases (with ESM-2 3B underperforming ESM-2 650M by  
717 approximately 0.01 in correlation). We attribute this strong correlation to the model’s evolutionary  
718 pretraining, which fundamentally enables it to learn mutation preferences from natural proteins,  
719 effectively capturing the constraints imposed by natural selection.

720 **(2) Explicitly aligning with evolutionary kernel further unlocks the potential of DPLM-EVO in**  
721 **mutation effect prediction.** We adopted the strategy proposed by VespaG Marquet et al. (2024) to  
722 explicitly align DPLM-EVO output distribution with GEMME (Laine et al., 2019), a state-of-the-art  
723 evolutionary-informed prediction model. Leveraging multiple sequence alignment information (De-  
724 rowicz et al., 2016), GEMME analyzes the evolutionary mutation sensitivity of individual sites,  
725 thereby providing a substitution distribution at each position. By aligning the substitution kernel  
726 of DPLM-EVO with that of GEMME, the scores correlate more closely with natural mutations.  
727 Illustrated in Fig.2B, this alignment yields further enhancements, outperforming SaProt (Su et al.,  
728 2023) that takes additional structure, TranceptEVE L with supplementary MSA, and the original  
729 VespaG method (based on ESM2-3B).

## 730 C.2 UNCONDITIONAL PROTEIN SEQUENCE GENERATION

731 **Setup.** We initialize DPLM-EVO from a pretrained DPLM-650M model (Wang et al., 2024b).  
732 Specifically, the backbone parameters (token embeddings and Transformer blocks) and the substitution  
733 head are initialized from DPLM, while the two binary operation heads for indel prediction (deletion  
734 and insertion) are randomly initialized. We train on the UniRef50 dataset for 100,000 steps, using  
735 2,000 warmup steps to a peak learning rate of  $10^{-4}$ , followed by linear decay to  $0.1 \times 10^{-4}$  by the  
736 end of training. The diffusion timestep is set to  $T = 500$ . For unconditional generation, we consider  
737 initial lengths  $L_{\text{init}} \in \{100, 200, 300, 400, 500\}$ .

738 **Results.** DPLM-EVO performs iterative denoising by jointly applying substitution, deletion, and  
739 insertion, enabling variable-length generation that more closely mirrors natural evolutionary trajec-  
740 tories. DPLM-EVO generation starts from corrupted sequences sampled from the learnable diffusion  
741 prior rather than all-mask initialization used in masked diffusion. Fig.3 demonstrates the evaluation  
742 results of unconditional generation in various perspectives: (1) **Diversity and Foldability:** Fig. 3A  
743 shows DPLM-EVO achieves consistent high foldability across length as measured by ESMFold  
744 pLDDT. Compared with DPLM based on masked diffusion, Fig.3D-E shows DPLM-EVO achieves  
745 comparable foldability while possessing greater generation diversity, reflected by a larger number  
746 of clusters in both sequence and structure. (2) **Reduced Mode Collapse:** DPLM-EVO produces  
747 higher sequence entropy than DPLM, as is shown in Fig.3F, indicating fewer repetition patterns and  
748 alleviating the mode collapse issue. (3) **Effect of Evolutionary Kernel:** Training with the *contextual*  
749 *evolutionary noising kernel* substantially outperforms uniform noising, as shown in Fig.3D. This  
750 indicates that biologically grounded corruptions encourages DPLM-EVO to learn more evolutionarily  
751 plausible substitution predictions, yielding higher-quality samples at generation time. (4) **Length**  
752 **Control:** Output lengths remain concentrated near their initial values without excessive expansion or  
753 collapse. The distribution is visualized in Fig.3B. This indicates that insertion and deletion prediction  
754 are invoked conservatively, resulting in a refinement process that prioritizes substitutions over drastic  
755 length changes.

To better understand how indel operations are scheduled over the diffusion trajectory, we probe  
the deletion and insertion heads across timesteps on a representative natural sequence (Fig.3C).

sample 50 candidates starting from the original sequence

Iter number	pLDDT		Clusters			
	mean/std	min/max	Seq Id 0.5	Seq Id 0.4	Seq Id 0.3	Struct TM 0.5
0 (original)	87.43	-/-	-	-	-	-
100	83.75/2.8	74.75/88.31	51	37	4	1

RMSD: 3.83      Seq Identity: 33.33%

Figure 5: *Unconstrained in-silico family expansion by DPLM-Evo.* 50 candidates are sampled from the original sequence show an overall high pLDDT, structure consistency and sequence diversity. One case is visualized to show structure alignment with low RMSD below 4 and sequence alignment with low identity at 33%.

The predicted indel probabilities decrease as the timestep decreases towards clean, suggesting the model primarily uses indels for coarse adjustments during high-noise stages, and gradually shifts to fine-grained refinements later. This behavior indicates that indel operations can be manipulated through timestep control, e.g., fixing the deletion timestep to 0 for insertion-only tasks.

### C.3 LENGTH-ADAPTIVE SCAFFOLDING OF FUNCTIONAL MOTIFS

**Setup.** Motif scaffolding aims to generate a protein scaffold for a given functional motif. We evaluate DPLM-EVO in *zero-shot* and *continued finetuning* settings. For finetuning, DPLM-EVO incorporates structural constraints for motif structure features, as illustrated in Fig.1C(4).

During generation, DPLM-EVO edits only the scaffold region and never modifies motif residues, allowing dynamic length adjustment to better accommodate the motif. In contrast, fixed-length sequence models require manually scaffold length enumeration and cannot revise length once an unsuitable initialization is chosen. For each motif instance, we sample 100 candidate scaffolds. Success is defined as pLDDT > 70, and motif RMSD < 1 Å.

**Results.** As shown in Fig.4, in the zero-shot setting, DPLM-EVO solves more motif problems than EvoDiff and DPLM-Mask, and achieves higher overall success rate (0.23). We attribute this to the capability for dynamic scaffolding length adjustment and evolutionarily plausible mutations provided by the substitution head. Continued finetuning brings further improvements, highlighting the importance of multimodal conditioning. Compared to multi-modal models like ESM-3 (Hayes et al., 2024), the finetuned DPLM-EVO achieves a higher overall success rate but resolves slightly less targets. We hypothesize this gap arises because DPLM-EVO only supports multimodal conditioning, without native end-to-end training for structural understanding. We leave multimodal evolutionary discrete diffusion modeling as an exciting direction for future work.

### C.4 CASE STUDY: IN-SILICO SEQUENCE FAMILY EXPANSION

**Setup.** To assess whether DPLM-EVO can generate diverse yet structurally consistent relatives of a given protein, we perform unconstrained post-editing starting from natural sequences. Specifically, we randomly select sequences from the CAMEO dataset and let DPLM-EVO refine them without imposing explicit functional constraints. We evaluate both structural preservation relative to the starting sequence and sequence diversification.

**Results.** DPLM-EVO generates diverse, yet structurally similar protein sequences in the unconstrained optimization setting. **Structural preservation:** We find that DPLM-EVO preserves structural plausibility (evaluated via comparing predicted structure to wild type structure) and at the same time introduces substantial edits. While it does not necessarily increase the initial pLDDT, it effectively explores the sequence space around a given fold without catastrophic structural degradation. **Sequence diversification and family expansion:** Meanwhile, DPLM-EVO modifies a large proportion of the initial sequence (with sequence identity mostly below 50%). Fig.5 shows a case where the highly modified sequence still aligns structurally with the original. These results suggest that DPLM-EVO performs unconstrained sequence optimization that preserves a shared structural scaffold while

810 producing diverse sequences. This implies that DPLM-EVO captures latent regularities of natural  
 811 proteins, including constraints related to fold and stability. In this way, the generated sequences can  
 812 be potentially viewed as in silico expanded homologs of the starting protein, holding the potential for  
 813 purely sequence based orpha protein understanding.

### 814 C.5 CASE STUDY: DIRECTED EVOLUTION OF GFP

815 **Setup.** We optimize the green fluorescent protein (GFP) via directed evolution using DPLM-EVO as  
 816 illustrated in Algorithm 2.

817 Starting from the template, we adopt the beam  
 818 search strategy to maintain a candidate set: in  
 819 each iteration, 10 optimized sequences are gener-  
 820 ated for each sequence in the candidate set.  
 821 In each step, we employ the Chai-1 model for  
 822 filtering and structure scoring to keep only the  
 823 top-scoring candidates retained for the next iter-  
 824 ation. Following ESM3 (Hayes et al., 2024),  
 825 The criteria for filter is that the template chromo-  
 826 phore site RMSD is less than 1.5Å, while  
 827 the scoring term is the pTM score produced by  
 828 Chai-1.

829 **Results.** Fig.6A depicts the trajectory of opti-  
 830 mization. We observed that as the iteration pro-  
 831 cesses, the pTM value gradually increases, while  
 832 the RMSD remained consistently below 1.5Å.  
 833 After 20 iterations, the pTM increased from  
 834 an initial 0.263 to 0.792 while a random mu-  
 835 tation baseline converges below 0.6. Fig.6B and  
 836 Fig.6C visualize the structures of the GFP before  
 837 and after optimization, respectively. Residues  
 838 are colored according to its pLDDT value, and  
 839 indicate a significant increase in stability. These  
 840 results prove that DPLM-EVO can leverage the  
 841 priors learned from evolution-scale protein se-  
 842 quences to optimize the GFP sequence toward  
 843 greater overall structural stability, while pre-  
 844 serving the structure of the chromophore site to  
 maintain its fluorescence effect.

## 845 D RELATED WORK

846 **Discrete diffusion model** are diffusion models operating in discrete state space (Sohl-Dickstein et al.,  
 847 2015; Austin et al., 2021). They noises samples with discrete transition probabilities and learn to  
 848 denoise them iteratively, in comparison to their continuous counterpart using continuous distribution  
 849 such as Gaussian kernels (Ho et al., 2020; Song et al., 2020a;b). Various transition kernels have been  
 850 explored, typically uniform transitions (Austin et al., 2021; Sahoo et al.) and masking (He et al., 2023;  
 851 Ye et al., 2023a; Zheng et al., 2023a; Sahoo et al., 2024; Shi et al., 2024; Nie et al., 2024). Among  
 852 the variants, masked diffusion has attracted the most recent interest for their simplicity (Zheng et al.,  
 853 2023a; Ou et al., 2024), scalability (Ye et al., 2023a; Nie et al., 2024), and empirical effectiveness,  
 854 gaining success as protein language models (Zheng et al., 2023b; Wang et al., 2024b;a; Hsieh et al.,  
 855 2025) and large language models (Ye et al., 2023a; Nie et al., 2025; Ye et al., 2025; Gong et al.,  
 2025).

856 **Flexible length generative model.** Pre-deciding the length of a answer to many problems can be hard  
 857 due to the uncertain computation budget and answer shapes required. While autoregressive models,  
 858 which iteratively produce one token each step, naturally provide flexible length generation, they  
 859 performs inferiorly in generating data without fixed-order structure such as protein sequences (Zheng  
 860 et al., 2023b). Non-autoregressive generative models (e.g., diffusion models), on the other hand, does  
 861 not require presumption on the orders but mostly require preset or predicted answer lengths (Guo et al.,  
 862 2019; Gu & Kong, 2021). To address this, researchers have explored enabling non-autoregressive  
 863 models to vary the output lengths by introducing indels operations in their predictions, which are  
 also referred to as edit-based models. As early attempt, Levenshtein Transformer (Gu et al., 2019)

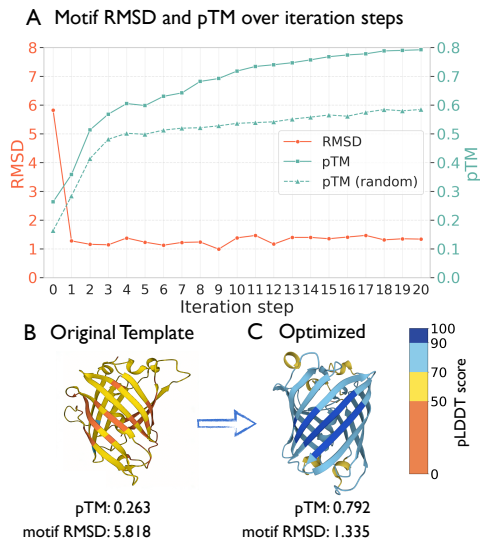


Figure 6: A case study of sequence optimization: directed evolution for GFP. (A) The motif RMSD and pTM scores are monitored through time step. DPLM-EVO maintains low motif RMSD while boosting pTM to over 0.8, compared to a random mutation baseline having pTM capped at 0.6. (B-C) Visualization of the optimization starting template and the final result.

864 studies non-autoregressive machine translations and show edit-based models perform on par with  
 865 autoregressive models while showing much lower sampling latency thanks to parallel generation.  
 866 Later progress revisit edit-based models under the formulation of diffusion models (Reid et al.,  
 867 2022). More recently, DreamOn (Wu et al., 2025) introduces a large diffusion language model with  
 868 indels as special tokens in vocabulary to vary the length during sampling, tailored for coding tasks.  
 869 EditFlow (Havasi et al., 2025), most related to our work, takes a flow-based model perspective to  
 870 construct indels training signals by aligning the noisy samples, which are interpolations between  
 871 clean samples and noises, with ground truth target. Although our work also extend fixed-length  
 872 diffusion to supports indels, we highlight the perspective of diffusion transition kernels and their  
 873 evolutionary significance.

874 **Protein language model.** Motivated by the success of large language models (LLMs), similar  
 875 practice has been extended to the development of protein language models. ESM-1b Rives et al.  
 876 (2019) utilizes self-supervised masked language modeling on 250 million protein sequences spanning  
 877 evolutionary diversity, later leading to the development of ESM-2 Lin et al. (2023) scales further.  
 878 ProtTrans Elnaggar et al. (2021), ProteinBERT Brandes et al. (2022), PRoBERTa Nambiar et al.  
 879 (2020), ProtAlbert Behjati et al. (2022), TAPE Rao et al. (2019), ProteinLM Xiao et al. (2021), and  
 880 CARP Yang et al. (2022) involve several other representative masked language modeling (MLM)  
 881 paradigm. These sequence-based PLMs perform competitively with classic methods that rely on mul-  
 882 tiple sequence alignments, indicating that PLMs have captured some of the evolutionary information  
 883 from sequences alone. In particular, these protein language models achieve powerful generalization  
 884 on various downstream tasks involving the secondary and tertiary structures. Recent findings further  
 885 showcase their capabilities in predicting protein functions Meier et al. (2021), structure folding Lin  
 886 et al. (2023), and de novo designs Verkuil et al. (2022). Beyond representation learning, DPLM (Wang  
 887 et al., 2024b) unlocks the generation capabilities of protein MLM through scalable discrete diffusion  
 888 training process (Ye et al., 2023b;a; Zheng et al., 2023b), enabling generating high-quality protein  
 889 sequences. Building upon DPLM, DPLM-2 (Wang et al., 2024a) and DPLM-2.1 (Hsieh et al., 2025)  
 890 further enhance the model with multimodal understanding and generation capabilities. The series lay  
 891 the foundation of pretrained models and diffusion algorithms for ours.

## 891 E GFP OPTIMIZATION PIPELINE

---

### 892 **Algorithm 2** GFP Directed Evolution by DPLM-EVO

---

- 893 1: **Initialization:** Starting from the GFP template sequence, add it to the candidate set  $\mathcal{C}$ .
  - 894 2: **Hyperparameters:**
  - 895 3:   max iteration  $T = 20$
  - 896 4:   search width  $w = 100$
  - 897 5:   beam size  $b = 10$
  - 898 6: **for**  $i = 1, \dots, T$  **do**
  - 899 7:   **Generate mutated sequences:** For each sequence in  $\mathcal{C}$ , generate  $w$  mutated sequences (one  
 900 position is mutated at a time).
  - 901 8:   Obtain a total of  $|\mathcal{C}| \times w$  samples.
  - 902 9:   **Filter candidates:** Filter the generated samples by  $\{\text{Common Filters}\}$ .
  - 903 10:   **Sort candidates:** Sort the filtered candidates according to  $\{\text{Common Score Terms}\}$ .
  - 904 11:   **Update candidate set:** Select the sequences with top  $b$  scores as the next-iteration candidate  
 905 set  $\mathcal{C}$ .
  - 906 12: **end for**
  - 907 13: **Return:** The final candidate set  $\mathcal{C}$
-