

PROTEIN GENERATION WITH EMBEDDING LEARNING FOR MOTIF DIVERSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

A fundamental challenge in protein design is the trade-off between generating structural diversity while preserving motif biological function. Current state-of-the-art methods, such as partial diffusion in RFdiffusion, often fail to resolve this trade-off: small perturbations yield motifs nearly identical to the native structure, whereas larger perturbations violate the geometric constraints necessary for biological function. We introduce Protein Generation with Embedding Learning (PGEL), a general framework that learns high-dimensional embeddings encoding sequence and structural features of a target motif in the representation space of a diffusion model’s frozen denoiser, and then enhances motif diversity by introducing controlled perturbations in the embedding space. PGEL is thus able to loosen geometric constraints while satisfying typical design metrics, leading to more diverse yet viable structures. We demonstrate PGEL on **five** representative cases: a monomer, a protein-protein interface, a cancer-related transcription factor complex, **an antibody-antigen complex and an enzyme**. In all cases, PGEL achieves greater structural diversity, better designability, and improved self-consistency, as compared to partial diffusion. Our results establish PGEL as a general strategy for embedding-driven protein generation allowing for systematic, viable diversification of functional motifs.

1 INTRODUCTION

Designing proteins that achieve precise biological functions while allowing for structural diversity has long been a central goal in computational protein design. Recent advances in structure prediction models like AlphaFold (Jumper et al., 2021; Abramson et al., 2024), RoseTTAFold (Baek et al., 2021), ESMFold (Lin et al., 2023) and Boltz (Wohlwend et al., 2024; Passaro et al., 2025) have revolutionized protein generative models and paved the way for improved diffusion models in protein design. Among them, RFdiffusion (Watson et al., 2023), which results from fine-tuning RoseTTAFold, has shown strong performance in both unconditional and conditional generation.

Yet, targeted local modification remains a challenge. A common approach is *partial diffusion* in RFdiffusion, in which a native or designed structure undergoes only a few denoising steps to induce diversification (Watson et al., 2023; Vázquez Torres et al., 2024). However, this method faces a fundamental diversity-fidelity trade-off: small structural perturbations keep near-native conformations, but lack diversity, while larger perturbations induce excessive geometric drift that disrupts functional features (Lin et al., 2024). Overcoming this limitation requires rethinking how diffusion models can introduce controlled variation while still anchoring designs to essential geometric constraints.

A promising direction comes from recent advances in conditional image generation. Models such as Stable Diffusion and Latent Diffusion Models (LDMs) generate images from noise guided by text prompts (Ho et al., 2020; Rombach et al., 2022). Beyond standard prompting, textual inversion learns new prompt embeddings to represent unseen visual concepts (Gal et al., 2022; Jin et al., 2024). Once learned, these embeddings can be diversified to generate outputs that preserve the original concept while exploring novel variations. Here we adopt this embedding-centric view in the context of protein generation.

We present Protein Generation with Embedding Learning (PGEL), a general framework representing the first adaptation of textual inversion principles to protein diffusion models. PGEL introduces

054 two key approaches with broad applicability: (1) learning high-dimensional embeddings that cap-
055 ture the sequence and structural characteristics of target protein regions of interest, thus shifting the
056 paradigm from coordinate-space to embedding-space perturbations, and (2) relaxing evolutionary
057 and structural constraints by masking embeddings. Although we present our work here using RFdif-
058 fusion’s representation space, our method is general and readily adaptable to other protein diffusion
059 models, and can thus leverage the rich representational capacity of pre-trained diffusion models
060 without expensive retraining or fine-tuning.

061 **In this work we formalize a new design task, *motif diversification*: given an experimentally charac-**
062 **terized structure with a functional motif embedded in a larger scaffold, we aim to generate alternative**
063 **motif conformations and poses while keeping the rest of the protein fixed in real space. Diversifica-**
064 **tion is defined at the backbone level, meaning that motif residues are allowed to move and rearrange**
065 **relative to each other and to the scaffold, but without pre-specifying side-chain identities or chem-**
066 **istry. Instead, we require that downstream sequence design and structure prediction can recover**
067 **sequences that realize each diversified backbone. Biologically, this corresponds to exploring fami-**
068 **lies of alternative structural realizations of an underlying functional motif (*e.g.*, a binding epitope or**
069 **active site) that remain compatible with the scaffold, enabling the enhancement of properties such**
070 **as affinity, specificity, or stability while preserving the overall protein architecture. Some existing**
071 **approaches address related challenges, but differ in scope and implementation: structure inpainting**
072 **methods (*e.g.*, masked region generation) fully marginalize a region by masking and regenerating**
073 **it *de novo*, discarding the specific native geometry (Zhang et al., 2023), whereas flexible backbone**
074 **loop remodeling in Rosetta (KIC/Next-Generation KIC) samples local conformations under explicit**
075 **geometric and energetic restraints to achieve high-fidelity but relatively localized exploration (Man-**
076 **dell et al., 2009; Stein & Kortemme, 2013; Leman et al., 2020). Hence these tools do not explicitly**
077 **target controlled exploration of a *neighborhood* around an existing functional motif while keeping a**
surrounding scaffold nearly fixed.

078 We, therefore, compare chiefly to partial diffusion in RFdiffusion, the prevailing stochastic baseline
079 for local variation which has been recently applied in therapeutically relevant design settings, includ-
080 ing *de novo* creation of high-affinity peptide binders and venom toxin neutralizers (Vázquez Torres
081 et al., 2024; 2025). Across five representative scenarios involving a monomeric protein (calmod-
082 ulin), a protein–protein binding site (barstar–barnase), a p53 binder within the p53-MDM2 complex,
083 **an antibody bound to Alzheimer’s disease’s amyloid beta peptide and the adenylate kinase enzyme,**
084 PGEL (1000 samples) produces more designable structures (motif pLDDT ≥ 70 , scRMSD $\leq 1\text{\AA}$,
085 mRMSD $\leq 2\text{\AA}$) than partial diffusion. PGEL also yields more structurally diverse TM-score clus-
086 ters distinguishable from native, and shows better self-consistency after inverse folding and refolding
087 (meeting mRMSD and pAE thresholds), while maintaining predicted binding affinities comparable
088 to native and exceeding those obtained with partial diffusion. Our results support embedding learn-
089 ing combined with masking as a general, efficient strategy for systematic motif diversification.

091 2 BACKGROUND

093 **Functional motifs.** Conditional generation around functional residues, often framed as *motif scaf-*
094 *olding*, has been a focal point for recent protein design methods. In that setting, the motif and
095 scaffold are defined as disjoint subsets with the scaffold varied while the motif geometry is pre-
096 served. Approaches like RFdiffusion (where the motif coordinates are fixed), the Monte Carlo-
097 based Twisted Diffusion Sampler (Wu et al., 2023) applied to FrameDiff, and Genie2 (Lin et al.,
098 2024) have made progress on this task, though performance remains task-dependent and can yield
099 few or no backbones meeting success criteria in specific cases. In our motif diversification task, the
100 scaffold is held fixed and the motif is diversified to explore multiple, function-preserving geometric
101 realizations, enabling improvements in *e.g.*, affinity, specificity or stability, while maintaining the
102 broader structural context.

103 **Protein embeddings.** The limited availability of structural data motivated the development of mod-
104 els that transform sequences into sequence embeddings that encode structural information. These
105 embeddings have been employed for various tasks such as property prediction using a Gaussian Pro-
106 cess regression model (Yang et al., 2018) and residue-residue contact prediction via a Bidirectional
107 Long Short-Term Memory (BiLSTM) architecture (Bepler & Berger, 2019). Transformers (Vaswani
et al., 2017) have been used in generating sequence embeddings, including for antibody-specific

108 applications like paratope prediction (Leem et al., 2022). Transfer learning has also been shown
 109 to significantly improve performance across architectures by enabling the use of pre-trained em-
 110 beddings that capture fundamental sequence-structure relationships (Detlefsen et al., 2022). Other
 111 approaches explicitly include structural information (Ali et al., 2024), such as contact maps-derived
 112 embeddings, and have shown enhanced performance in particular downstream tasks such as struc-
 113 ture similarity assessment (Kandathil et al., 2025), structure searching (Greener & Jamali, 2024),
 114 property prediction (Blaabjerg et al., 2024; Danner et al., 2025), and domain classification (Lau
 115 et al., 2024). Similarly, protein function annotation and local flexibility prediction have benefited
 116 from Graph Convolutional Networks, which combine structure-derived graphs to propagate contex-
 117 tual signals from protein sequence embeddings obtained with pre-trained models (Gligorijević et al.,
 118 2021; Michalewicz et al., 2025).

119 **Diffusion models for proteins.** Earlier works adapted Denoising Diffusion Probabilistic Mod-
 120 els (DDPMs) to protein design by conditioning on local structural elements or coarse fold con-
 121 straints (Wu et al., 2024; Anand & Achim, 2022; Trippe et al., 2023; Luo et al., 2022) yet, while
 122 encouraging, they produced few sequences that refolded to target backbones. RFdiffusion subse-
 123 quently emerged as the diffusion approach that reliably yields designable structures and sequences
 124 that recover the intended geometry. In RFdiffusion, a highly accurate protein structure prediction
 125 method (RoseTTAFold (Baek et al., 2021)) is fine-tuned to undo random perturbations of atomic
 126 coordinates introduced via 3D Gaussian noise (*i.e.*, to denoise). RFdiffusion can be constrained to
 127 specific binding targets, or symmetry specifications, and once trained it can be viewed as a *frozen*
 128 *denoiser*. RoseTTAFold/AlphaFold-style models (including RFdiffusion) learn so-called *state* and
 129 *pair* embeddings (related to per-residue and residue-residue properties of the protein structure, re-
 130 spectively) and MSA embeddings related to multiple sequence alignment (Jumper et al., 2021).

131 **Textual inversion.** Gal et al. (2022) builds on LDMs (Rombach et al., 2022), a specific class of
 132 DDPMs, to perform textual inversion. In the context of text-to-image models, let x represent an
 133 image, s a text prompt, ϵ_θ a pre-trained denoising network, and ε an image encoder. LDMs aim to
 134 minimize the following loss:

$$\mathcal{L}_{\text{LDM}} := \mathbb{E}_{z \sim \varepsilon(x), s, \epsilon \sim \mathcal{N}(0,1), t} [\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(s))\|_2^2] \quad (1)$$

135 Here, $c_\theta(s)$ represents a pre-trained text encoder that conditions the denoiser ϵ_θ based on the text
 136 prompt s , and z_t is a noised version of the image embedding z at timestep t . The goal of textual
 137 inversion is to learn a new text embedding v_* corresponding to a particular concept s_* such that it
 138 minimizes the LDM loss (equation 1). This means conditioning ϵ_θ on v_* so the generated image
 139 \tilde{x} closely resembles the original image x . Neutral prompts, such as “A photo of s_* ” or “A portrait
 140 of s_* ” are used while keeping ϵ_θ and c_θ frozen. Multi-Concept Prompt Learning (Jin et al., 2024)
 141 extends this idea to handle multiple concepts by incorporating three regularization techniques: at-
 142 tention masking, bind adjective, and prompts contrastive loss.
 143
 144

145 3 METHODS

146 We now present our method, *Protein Generation with Embedding Learning (PGEL)*, and describe
 147 how we learn the embedding representation of a motif in Section 3.1. In Section 3.2, we propose an
 148 approach to increase motif diversity, and Section 3.3 details the evaluation metrics.

149 3.1 PROTEIN GENERATION WITH EMBEDDING LEARNING (PGEL)

150 We generalize the notion of textual inversion with LDMs to proteins, treating the structure as anal-
 151 ogous to an image, and the sequence as analogous to a text prompt. Let R_* be a region of interest,
 152 or *motif*, defined as a continuous or discontinuous set of L_* amino acids within a protein. The motif
 153 has structure x_* and sequence s_* , where the coordinates of x_* are obtained from an experimental
 154 Protein Data Bank (PDB) entry, and the sequence s_* is *masked* when passed as an input to PGEL,
 155 *i.e.*, the amino acid range of the motif is specified, but not its exact composition.
 156
 157

158 PGEL learns a representation of R_* in embedding space, which we denote as v_* . The remainder of
 159 the protein constitutes the *scaffold*, with structure x_c and sequence s_c of length L_c , from which the
 160 protein LDM frozen ENCODER computes an embedding representation v_c .
 161

The procedure (see Figure 1 and Algorithm 1) starts by building a noised protein structure in which the scaffold coordinates are retained while the motif coordinates are subjected to T rounds of Gaussian noise injection, following Trippe et al. (2023). At each timestep t , the protein LDM frozen DENOISER predicts a denoised motif structure $\hat{x}_*^{(0)}$, conditioned jointly on the learnable motif embedding v_* and the fixed embedding v_c . These embeddings include state, pair and MSA embeddings. Then, by using structure $x^{(t)}$ and the intermediate structure $[x_c, \hat{x}_*^{(0)}]$, a reverse diffusion step REVERSESTEP, which does not contain any learnable parameters, yields $x^{(t-1)}$ (see Algorithm 3). In practice, we employ pre-trained building blocks of RFDiffusion for both the ENCODER and DENOISER, though alternative models could be substituted if desired.

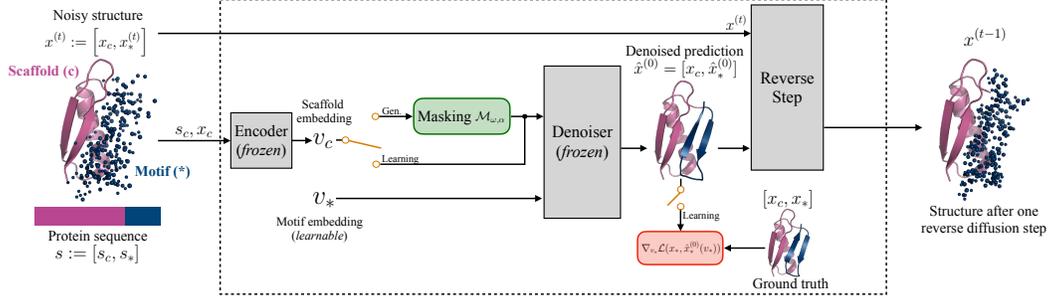


Figure 1: Outline of the PGEL learning and generation procedures during one reverse diffusion step.

Embedding optimization. The embedding v_* is learned by minimizing:

$$\mathcal{L} = \mathcal{L}_{\text{MSE}} + \lambda_{\text{DM}} \mathcal{L}_{\text{DM}} + \lambda_{\text{torsion}} \mathcal{L}_{\text{torsion}} \quad (2)$$

This loss function is composed of three terms, described hereafter, which compare different features of the ground truth structure x_* and the predicted structure $\hat{x}_*^{(0)}(v_*)$ of the motif with the coefficients $\lambda_{\text{DM}}, \lambda_{\text{torsion}} \in \mathbb{R}_{\geq 0}$ controlling the relative weight of the terms.

Data fidelity term (backbone atoms). For each motif residue $i \in R_*$ we consider the $A = 4$ backbone atoms (nitrogen N, α -carbon C_α , carbon C, oxygen O). Let $\hat{x}_{i,a}^{(0)} \in \mathbb{R}^3$ denote the predicted position of atom a in residue i and $x_{i,a} \in \mathbb{R}^3$ its ground truth counterpart. We then compute the mean squared error (MSE) between the backbone atoms of the ground truth and predicted motif:

$$\mathcal{L}_{\text{MSE}}(x_*, \hat{x}_*^{(0)}(v_*)) = \frac{1}{AL_*} \sum_{i \in R_*} \sum_{a \in \{N, C_\alpha, C, O\}} \left\| \hat{x}_{i,a}^{(0)}(v_*) - x_{i,a} \right\|^2 \quad (3)$$

Distance matrix between α -carbons. Let $\hat{x}_{i,C_\alpha}^{(0)} \in \mathbb{R}^3$ denote the predicted position of the α -carbon atom in residue i , and $x_{i,C_\alpha} \in \mathbb{R}^3$ its ground truth counterpart. We define the following loss term based on α -carbon Distance Matrices (DM), inspired by the distrogram notion (Senior et al., 2020):

$$\mathcal{L}_{\text{DM}}(x_*, \hat{x}_*^{(0)}(v_*)) = \frac{1}{L_*^2} \sum_{i \in R_*} \sum_{j \in R_*} \left(\left\| \hat{x}_{i,C_\alpha}^{(0)}(v_*) - \hat{x}_{j,C_\alpha}^{(0)}(v_*) \right\| - \|x_{i,C_\alpha} - x_{j,C_\alpha}\| \right)^2 \quad (4)$$

In contrast to \mathcal{L}_{MSE} , \mathcal{L}_{DM} is invariant under rigid motions (translations and rotations), thus encouraging global shape consistency.

Backbone torsion angles. Let $\hat{\phi}_i$ and $\hat{\psi}_i$ denote the predicted backbone torsion angles at residue i , computed from $\hat{x}_*^{(0)}(v_*)$, and let ϕ_i and ψ_i be the corresponding ground truth values (Ramachandran et al., 1963). We impose a constraint on angular torsions through a cosine-based loss term akin to that of AlphaFold (Jumper et al., 2021):

$$\mathcal{L}_{\text{torsion}}(x_*, \hat{x}_*^{(0)}(v_*)) = \frac{1}{L_* - 2} \sum_{i=2}^{L_*-1} \left[1 - \cos(\hat{\phi}_i(v_*) - \phi_i) + 1 - \cos(\hat{\psi}_i(v_*) - \psi_i) \right] \quad (5)$$

This term penalizes sterically implausible geometries, helping improve performance under the predicted local distance difference test (pLDDT). We do not include the third backbone angle ω as it is typically considered fixed at 180 degrees (Cutello et al., 2006).

Algorithm 1 PGEL – Embedding learning

Input: region of interest/motif R_* with masked sequence s_* and structure x_* , fixed scaffold with sequence s_c and structure x_c , pre-trained ENCODER and DENOISER.
Output: learned embedding v_* for region R_* .
initialize v_* with zeros.
while not converged **do**
 Build noised structure $x^{(T)} := [x_c, x_*^{(T)}]$ with associated sequence $s := [s_c, s_*]$.
 $v_c = \text{ENCODER}(s_c, x_c)$
 for $t = T$ **down to** 1 **do**
 $\hat{x}_*^{(0)} = \text{DENOISER}(v_c, v_*)$
 $x^{(t-1)} = \text{REVERSESTEP}(x^{(t)}, [x_c, \hat{x}_*^{(0)}])$
 Update v_* by taking a gradient step $\nabla_{v_*} \mathcal{L}(x_*, \hat{x}_*^{(0)}(v_*))$
 end for
end while
Return v_*

Once the embeddings are learned for a particular protein, we employ Algorithm 2 to generate novel proteins with diversified region of interest R_* (see Figure 1). We demonstrate the necessity of learning motif embeddings by comparing with simpler baselines in Appendix D.

Algorithm 2 PGEL – Generation with embedding masking

Input: region of interest/motif R_* with learned embeddings v_* , fixed scaffold with sequence s_c and structure x_c , pre-trained ENCODER and DENOISER.
Output: generated structure.
Build noised structure $x^{(T)} := [x_c, x_*^{(T)}]$.
Draw at random the sample masking type $\omega \sim \text{Ber}(\frac{1}{2})$ (row if 0, column if 1).
Sample masking rate $\alpha \sim \mathcal{U}[0, 1]$.
Define $\mathcal{M}_{\omega, \alpha}(\cdot)$ as a zero mask with type ω and rate α .
 $v_c = \text{ENCODER}(s_c, x_c)$
for $t = T$ **down to** 1 **do**
 $\hat{x}_*^{(0)} = \text{DENOISER}(\mathcal{M}_{\omega, \alpha}(v_c), v_*)$
 $x^{(t-1)} = \text{REVERSESTEP}(x^{(t)}, [x_c, \hat{x}_*^{(0)}])$
end for
Return $x^{(0)}$

3.2 ENHANCING THE DIVERSITY OF GENERATED MOTIFS

MSA embeddings in sequence-to-structure predictors contain evolutionary covariation information about residues, thereby capturing geometric constraints such as residue proximity. In RFdiffusion, however, such embeddings are derived solely from the input sequence $s := [s_c, s_*]$ rather than from a full stack of aligned sequences, and can be represented as a $d_{\text{MSA}} \times L$ matrix, where d_{MSA} is the depth of the MSA embeddings and $L := L_* + L_c$ the total protein length. With PGEL, we show that the diversity of generated structures can be increased by applying perturbations to the scaffold MSA embeddings $v_c^{\text{MSA}} \in \mathbb{R}^{d_{\text{MSA}} \times L_c}$. These embeddings couple through attention mechanisms with state and pair embeddings produced by an internal RFdiffusion encoder, and also interact with the learned motif embedding v_* , which provides an independent conditioning signal for the frozen denoiser (see Appendix A for more details).

Embedding masking. We studied the effect of applying zero masks, *i.e.* masks zeroing specific elements, to the scaffold MSA embeddings during generation (see Algorithm 2). *Row masking* corresponds to masking specific features for all residues, whereas *column masking* zeroes out all features of specific residues. Both strategies lift some constraints on inter-residue distances, and modulate which co-variation patterns remain accessible. We sample $\omega \sim \text{Ber}(\frac{1}{2})$ to choose the masking mode ($\omega = 0$ for row masking and $\omega = 1$ for column masking) and $\alpha \sim \mathcal{U}[0, 1]$, the masking rate, to set the fraction of rows or columns masked. This defines the operator $\mathcal{M}_{\omega, \alpha}(\cdot)$, which implements zero masking with type ω and rate α . As such, masking v_c relaxes the geometric

constraints of the generated motif (see Appendix B), while its omission results in no diversity, as shown in Appendix D. We also assessed the correlation between motif structural diversity and rate α in Appendix E.

3.3 EVALUATION METRICS

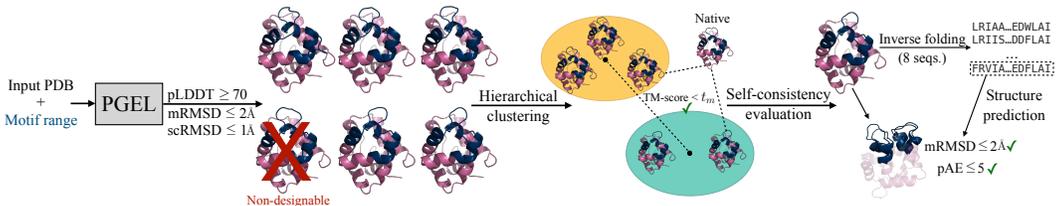


Figure 2: Summary of the evaluation metrics. PGEL takes as input a PDB entry and the amino acid range corresponding to the motif. From 1000 PGEL-generated backbones, designable candidates are filtered by root mean square deviation (RMSD) and pLDDT thresholds, and structural diversity is assessed via hierarchical clustering. Backbones are also required to be *distinguishable* from the native. Cluster representatives undergo self-consistency evaluation: sequences assigned to the designable backbones with ProteinMPNN are refolded, and at least one predicted structure must satisfy set mRMSD and predicted alignment error (pAE) conditions relative to the generated backbone.

Designability. To quantify designability in the motif diversification task, we first require a motif $pLDDT \geq 70$ as computed by an RFdiffusion internal block, following the threshold adopted for related tasks by Lin et al. (2024). We also require the scaffold RMSD to be $scRMSD \leq 1\text{\AA}$, to ensure that the residues surrounding the motif remain fixed. Finally, we set the threshold for the motif RMSD to $mRMSD \leq 2\text{\AA}$, to allow for structural diversification of the motif.

Diversity. To quantify structural diversity among generated proteins (Figure 2), we compute the pairwise TM-scores (Zhang & Skolnick, 2004) across all designable candidates, and employ hierarchical clustering (Lin et al., 2024) with several linkage thresholds t_m to group similar backbones under this score. Diversity is measured by the number of clusters. We evaluate also the TM-score with respect to the native motif: a cluster is considered *distinguishable* from native if, for TM-score threshold $t_m \in [0, 1]$, at least one cluster member exhibits lower similarity than t_m relative to the native. This analysis ensures that we capture the structural distinctiveness of the backbones.

Self-consistency. For the distinguishable backbones we use the procedure in Trippe et al. (2023) based on inverse folding to assess self-consistency between generated and predicted structures. Specifically, we use ProteinMPNN with default parameters (Dauparas et al., 2022) to assign 8 plausible sequences to each backbone, followed by a sequence-to-structure model, here AlphaFold3 (Abramson et al., 2024), to predict 8 full proteins. A designed backbone is deemed self-consistent if it satisfies for at least one of the 8 predicted structures: $mRMSD \leq 2\text{\AA}$ and $pAE \leq 5$ (Figure 2). Previous studies included this procedure under the designability assessment (Watson et al., 2023; Lin et al., 2024). However, this is computationally expensive and, when prioritizing diversity, often inefficient: many backbones either fail the initial scRMSD or mRMSD filters or exhibit negligible structural diversity. In motif diversification, diversity among generated proteins and distinguishability with respect to the native are decisive. We therefore invert the pipeline to enforce these criteria first and reserve the costly self-consistency evaluation only for diverse candidates.

Binding affinity. For protein-protein complexes, we run PRODIGY (Vangone & Bonvin, 2015; Xue et al., 2016) to estimate the binding affinity ΔG expressed in kcal/mol, with larger $|\Delta G|$ values indicating stronger binding. It is desirable that new designs present binding affinity values comparable to, or larger in magnitude than, those of the native complex. Note that learning is not optimized to enhance binding affinity; rather, this serves as an *a posteriori* assessment.

4 EXPERIMENTS

We focus on five representative test cases: (1) Calmodulin, a monomer that plays a pivotal role in regulating the activity of nearly 100 diverse target enzymes and structural proteins (Fallon & Quio-

Table 1: Comparison of partial diffusion in RFdiffusion and PGEL. The number of self-consistent clusters (diversity) is computed at TM-score threshold $t_m = 0.6$.

Example number	Designability					Diversity				
	(No. of viable structures out of 1000)					(No. of self-consistent clusters)				
	1	2	3	4	5	1	2	3	4	5
Partial diffusion	411	331	252	948	798	0	6	1	2	8
PGEL	1000	990	802	950	823	2	10	6	4	15

cho, 2003); (2) the barstar-barnase complex, in which the binding interface of barstar was diversified to probe its interaction with the extracellular ribonuclease barnase (Caro et al., 2023); (3) the cancer-related transcription factor p53 bound to its negative regulator MDM2 (Klein & Vassilev, 2004; Li et al., 2010); (4) a TAP01 family antibody in complex with an amyloid beta peptide, which is related to Alzheimer’s disease; (5) a mutated version of the adenylate kinase enzyme.

4.1 PROTOCOL

We established a protocol to systematically compare our method with RFdiffusion’s partial diffusion using the metrics introduced in Section 3.3. For partial diffusion, we generated 1000 protein backbones by uniformly sampling the number of diffusion timesteps, $T \sim \mathcal{U}\{2, 3, \dots, 49\}$, as $T = 50$ corresponds to the full diffusion process in RFdiffusion. In this way, we cover a spectrum of structural perturbations ranging from near-native backbones to unrelated ones. For PGEL, we performed the learning of v_* with Stochastic Gradient Descent with learning rate $l_r = 4 \times 10^{-4}$ and momentum $p = 0.9$, $\lambda_{\text{DM}} = 0.01$ and $\lambda_{\text{torsion}} = 0.05$ (Algorithm 1), and we then generated 1000 protein backbones (Algorithm 2).

For both sets of 1000 generated structures, we evaluated designability and, among those deemed designable, we computed TM-scores between all generated motifs and with respect to the native structure. We then plotted the number of clusters as a function of the TM-score. For structures that were designable and diverse according to a typical TM-score threshold $t_m = 0.6$ (Lin et al., 2024), we performed inverse folding through ProteinMPNN to generate compatible sequences, followed by AlphaFold3 inference to assess whether the predicted sequences refolded into the intended backbones, fulfilling the self-consistency requirement defined in Section 3.3.

4.2 EXAMPLE 1: MONOMER

Calmodulin (PDB entry: 1PRW), a monomeric protein containing a double EF-hand motif spanning residues 16-35 and 52-71 (structured multi-motif), was considered as the representative test case for single-chain proteins. All the 1000 backbones candidates generated by PGEL resulted to be designable, well exceeding the 411 obtained by partial diffusion (Table 1). All of the backbones generated by partial diffusion satisfied the pLDDT constraint, consistent with the fact that RFdiffusion’s training favors high-confidence local structures, but 589 of them failed to meet the expected motif RMSD threshold. In these cases, the added noise during diffusion excessively perturbed the initial backbone, leading to conformations that no longer preserved the intended geometry of the EF-hand motif.

We then evaluated the structural diversity of the designable backbones, recording the number of clusters as a function of the TM-score threshold (Figure 3A). PGEL consistently produced a higher number of clusters across thresholds, demonstrating that embedding perturbations through masking introduce greater variability in backbone conformations. On the other hand, partial diffusion yielded structures too similar to the native backbone, and hence not distinguishable from it.

We carried out the self-consistency assessment at $t_m = 0.6$, as per protocol. For PGEL, the two clusters had backbones that successfully refolded into the intended conformations after sequence design and AlphaFold3 inference. Figure 3B illustrates these two successful cases, along with examples of backbones generated by partial diffusion that either did not satisfy the mRMSD metric condition or the distinguishability from native.

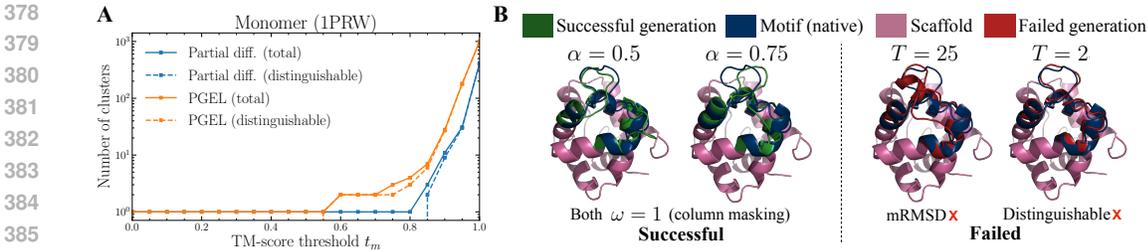


Figure 3: Results for example 1. (A) Number of clusters, both total and distinguishable from native, as a function of the TM-score threshold t_m for PGEL and partial diffusion. (B) *Left*: two successful PGEL designs at $t_m = 0.6$ using column masking with rates $\alpha = 0.5$ and $\alpha = 0.75$. *Right*: two partial diffusion failed backbones at $t_m = 0.6$, one obtained with $T = 25$ timesteps that violates the motif RMSD constraint, and one with $T = 2$ timesteps that is not distinguishable from the native.

4.3 EXAMPLE 2: BINDING SITE

Barstar is a small protein that binds the active site of barnase to prevent the latter from breaking RNA. This toxin-antitoxin pair (PDB entry: 7MRX) has been increasingly exploited in cancer therapy for targeted cytotoxicity (Kalinin et al., 2023). As target for motif diversification, we took the barstar’s binding interface region comprising residues 25 to 46, which can be classed as a structured single-motif (see Watson et al. (2023)).

Of the 1000 backbones generated with PGEL, 990 were classified as designable, nearly tripling the 331 obtained with partial diffusion. Correspondingly, Figure 4A demonstrates that PGEL consistently outperforms partial diffusion across the entire range of $t_m \in [0, 1]$, with pronounced differences observed at $t_m > 0.9$ and within $0.45 < t_m < 0.55$, around canonical TM-score thresholds. At $t_m = 0.6$, PGEL yielded 15 structural clusters compared to 13 for partial diffusion, which were reduced to 10 and 6, respectively, after self-consistency checks (Table 1).

In Figure 4A, we also display an overlay version of the native motif and 10 representative motifs derived from these clusters, highlighting the sequence variability both among generated barstar binding interfaces and relative to the native PDB structure. When predicting *in silico* the binding affinity of the generated complexes with PRODIGY, two of the generated structures exhibited binding affinities higher than the native complex (see Figure 5A and Table 5), while the remaining eight retained at least 80% of the original affinity ΔG_{native} . In contrast, only one complex generated by partial diffusion had a binding affinity value comparable to that of the native ($\Delta G_{\text{design}} > 0.9\Delta G_{\text{native}}$).

4.4 EXAMPLE 3: BINDER

The interaction between the transcription factor p53 and its negative regulator MDM2 is a key molecular process in cancer progression. Specifically, pharmacological disruption of the p53-MDM2 complex restores p53 activity and has been proven beneficial in cancer therapy (Hu et al., 2021).

Starting from PDB entry 1YCR, we addressed the motif diversification task by redesigning the complete p53 under the RMSD constraints described in Section 3.3. PGEL generated 802 designable backbones out of 1000 trials, with most non-designable cases attributable to low pLDDT confidence scores (Table 1). Partial diffusion produced only 252 designable backbones with considerably reduced structural diversity (a single cluster at TM-score threshold $t_m = 0.6$, see Figure 4B). PGEL, by comparison, gave six clusters at $t_m = 0.6$, all of which passed the self-consistency checks.

When assessing the binding affinity *a posteriori*, five out of six representatives of PGEL clusters exhibited lower affinity compared to the sole valid instance of partial diffusion (Figure 5B, Table 5). Notably, sequence SSMWELWQEIEGE (see Figure 4B), designed with PGEL in combination with ProteinMPNN, folded, as predicted by AlphaFold3, into a structure with a binding affinity comparable to that of the native structure, despite sharing only around 15% of sequence identity. This result highlights PGEL’s ability to generate backbones that can accommodate sequences unrelated to the native while refolding into structures that preserve function.

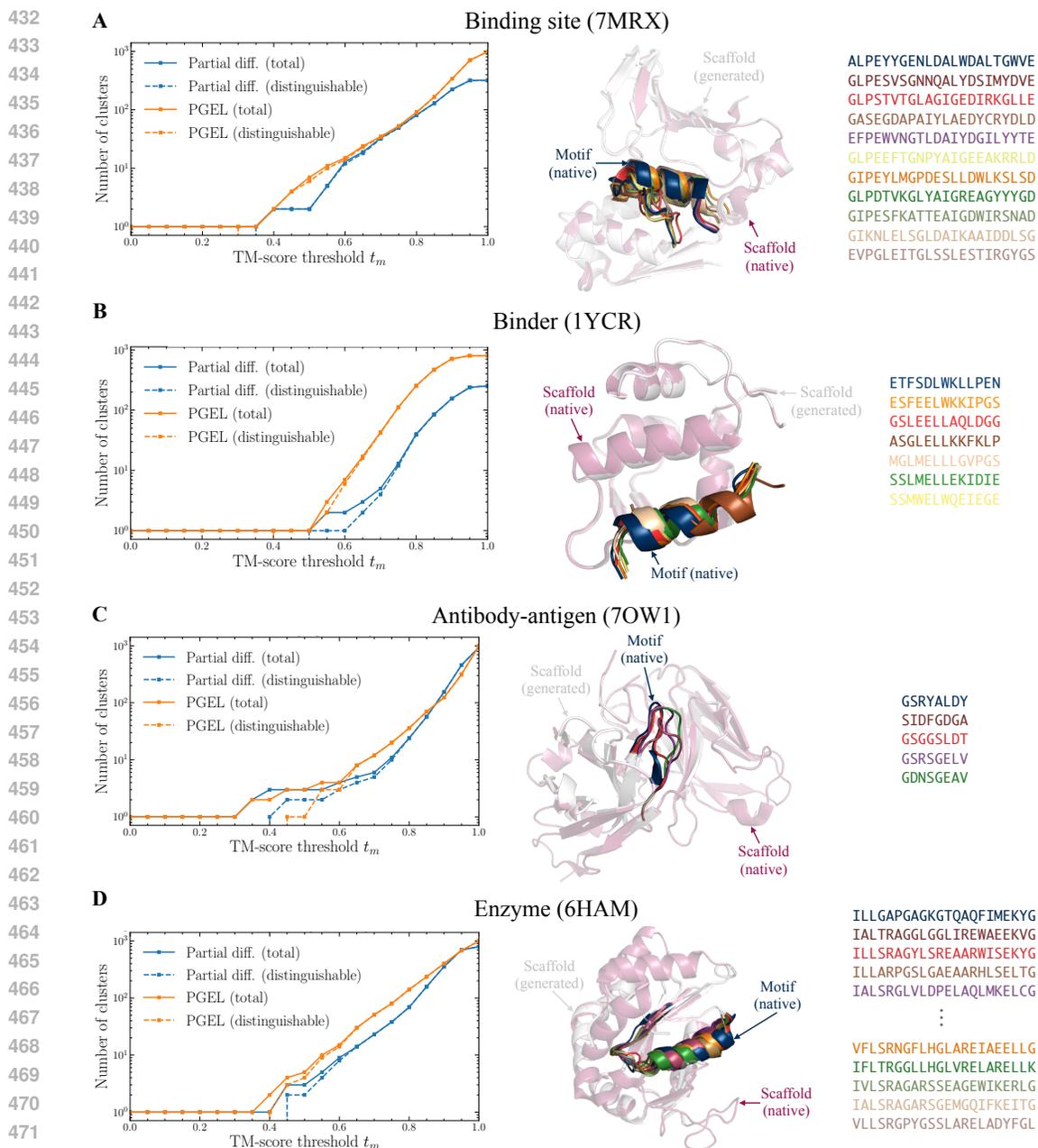


Figure 4: (A) Example 2: *Left*: Number of clusters identified PGEL and partial diffusion in the binding site example, both total and distinguishable from native, as a function of the TM-score threshold t_m . *Right*: generated binding site backbones (overlaid with native), alongside the native sequence and sequences that refold to self-consistent structures. (B,C,D) Same as A but for examples 3, 4 and 5, respectively.

4.5 EXAMPLE 4: ANTIBODY-ANTIGEN COMPLEX

The gradual accumulation of the peptide amyloid beta in neural tissue leads to the formation of insoluble aggregates, a pathological feature of Alzheimer’s disease (van Dyck, 2018). PDB entry 7OW1 contains a TAP01 family antibody bound to amyloid beta. We diversified the third complementarity-determining region of the antibody heavy chain (CDR-H3), an unstructured (loop) motif, which in this case comprises 8 residues and is commonly targeted to enhance antibody binding proper-

ties (Michalewicz et al., 2024). In terms of designability and diversity, both PGEL and partial diffusion produced 4 clusters that are distinguishable from the native structure. However, following self-consistency checks, only 2 out of the 4 clusters remained valid for partial diffusion.

4.6 EXAMPLE 5: ENZYME

Adenylate kinase (AdK) from *E. coli*, found in PDB entry 6HAM, is a multi-domain enzyme featuring a Walker A motif in its N-terminal region (Kantaev et al., 2018; Jafri et al., 2025). We explored the generation of AdK variants by diversifying the Walker A motif and its surrounding residues (positions 4-25; structured single-motif), which resulted in the highest number of self-consistent clusters among all five experiments: 15 for $t_m = 0.6$, compared to 8 with partial diffusion. We also identified a valid cluster at $t_m = 0.4$ containing sequences such as IALSRGLVLDPELAQLMKELCG while also satisfying the mRMSD constraint despite exhibiting only 27% sequence identity relative to the native.

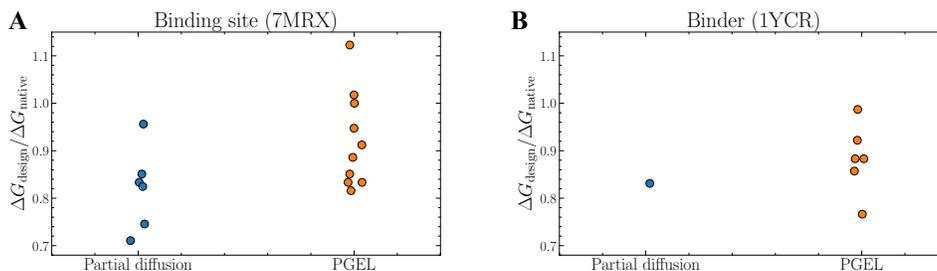


Figure 5: Ratios of the binding affinity predicted with PRODIGY of the native structure vs the AlphaFold3-predicted structures from backbones (one per cluster) generated by PGEL and partial diffusion for example 2 (A) and example 3 (B).

Computational remarks. Across the five case studies, column masking was more effective than row masking (see Table 4, with nearly 80% of successful outcomes with $\omega = 1$).

The time required per timestep during the generation process is nearly indistinguishable between PGEL and partial diffusion: average timestep 0.7 s for partial diffusion and 0.71 s for PGEL on a single NVIDIA GeForce RTX 3090 GPU with 24GB of memory.

5 LIMITATIONS AND FUTURE WORK

PGEL inherits the biases and limitations of the underlying frozen RFdiffusion denoiser, including its training data distribution and architectural constraints. Moreover, learning embeddings requires additional optimization time, which ranged in our examples from 2 minutes (example 2) to 2 hours (example 3) on a single GPU, with a trade-off between speed and improved results (more details in Appendix F). Our experiments were limited to motifs of up to 40 residues, with practical limits of around 50 residues given available memory, though scaling to longer motifs should be feasible with larger hardware or engineering optimization. Beyond this, our evaluation is entirely *in silico* (pLDDT/RMSD/TM-score filtering, AlphaFold3 refolding, and PRODIGY ΔG) and thus predictive rather than experimental.

Future work will investigate alternative ways of perturbing embeddings, as this strategy for motif diversification remains largely unexplored, as well as different strategies for sampling the masking parameters ω and α . For instance, instead of sampling α uniformly between 0 and 1, one could bias it toward smaller masking rates (e.g., using a Poisson distribution with rate λ , where λ tunes how conservative or aggressive the masking is), thus providing finer control over structural perturbations. A more systematic mapping between PGEL’s ω and α and partial diffusion’s T would also clarify the relationship between the diversity-fidelity trade-off in both methods. Finally, experimental validation will be pursued in follow-up work.

Code availability. Upon publication, we will release code and configurations to facilitate reproducibility.

REFERENCES

- 540
541
542 Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf
543 Ronneberger, Lindsay Willmore, Andrew J. Ballard, Joshua Bambrick, Sebastian W. Boden-
544 stein, David A. Evans, Chia-Chun Hung, Michael O’Neill, David Reiman, Kathryn Tunyasuvu-
545 nakool, Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex
546 Bridgland, Alexey Cherepanov, Miles Congreve, Alexander I. Cowen-Rivers, Andrew Cowie,
547 Michael Figurnov, Fabian B. Fuchs, Hannah Gladman, Rishub Jain, Yousuf A. Khan, Caro-
548 line M. R. Low, Kuba Perlin, Anna Potapenko, Pascal Savy, Sukhdeep Singh, Adrian Stecula,
549 Ashok Thillaisundaram, Catherine Tong, Sergei Yakneen, Ellen D. Zhong, Michal Zielinski,
550 Augustin Židek, Victor Bapst, Pushmeet Kohli, Max Jaderberg, Demis Hassabis, and John M.
551 Jumper. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*,
552 630(8016):493–500, Jun 2024. ISSN 1476-4687. doi: 10.1038/s41586-024-07487-w. URL
553 <https://doi.org/10.1038/s41586-024-07487-w>.
- 554 Sarwan Ali, Prakash Chourasia, and Murray Patterson. When protein structure embedding meets
555 large language models. *Genes*, 15(1), 2024. ISSN 2073-4425. doi: 10.3390/genes15010025.
556 URL <https://www.mdpi.com/2073-4425/15/1/25>.
- 557 Namrata Anand and Tudor Achim. Protein structure and sequence generation with equivariant de-
558 noising diffusion probabilistic models. *arXiv preprint arXiv:2205.15019*, 2022.
559
- 560 Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, Sergey Ovchinnikov, Gyu Rie
561 Lee, Jue Wang, Qian Cong, Lisa N. Kinch, R. Dustin Schaeffer, Claudia Millán, Hahnbeom
562 Park, Carson Adams, Caleb R. Glassman, Andy DeGiovanni, Jose H. Pereira, Andria V. Ro-
563 drigues, Alberdina A. van Dijk, Ana C. Ebrecht, Diederik J. Opperman, Theo Sagmeister,
564 Christoph Buhlheller, Tea Pavkov-Keller, Manoj K. Rathinaswamy, Udit Dalwadi, Calvin K.
565 Yip, John E. Burke, K. Christopher Garcia, Nick V. Grishin, Paul D. Adams, Randy J. Read,
566 and David Baker. Accurate prediction of protein structures and interactions using a three-
567 track neural network. *Science*, 373(6557):871–876, 2021. doi: 10.1126/science.abj8754. URL
568 <https://www.science.org/doi/abs/10.1126/science.abj8754>.
- 569 Tristan Bepler and Bonnie Berger. Learning protein sequence embeddings using information from
570 structure. In *International Conference on Learning Representations, ICLR’19*, 2019.
571
- 572 Lasse M. Blaabjerg, Nicolas Jonsson, Wouter Boomsma, Amelie Stein, and Kresten Lindorff-
573 Larsen. Ssemb: A joint embedding of protein sequence and structure enables robust variant effect
574 predictions. *Nature Communications*, 15(1):9646, Nov 2024. ISSN 2041-1723. doi: 10.1038/
575 s41467-024-53982-z. URL <https://doi.org/10.1038/s41467-024-53982-z>.
- 576 José A. Caro, Kathleen G. Valentine, Taylor R. Cole, and A. Joshua Wand. Pressure, motion, and
577 conformational entropy in molecular recognition by proteins. *Biophysical Reports*, 3(1), Mar
578 2023. ISSN 2667-0747. doi: 10.1016/j.bpr.2022.100098. URL [https://doi.org/10.](https://doi.org/10.1016/j.bpr.2022.100098)
579 [1016/j.bpr.2022.100098](https://doi.org/10.1016/j.bpr.2022.100098).
- 580 Vincenzo Cutello, Giuseppe Narzisi, and Giuseppe Nicosia. A multi-objective evolution-
581 ary approach to the protein structure prediction problem. *Journal of The Royal So-*
582 *ciety Interface*, 3(6):139–151, 2006. doi: 10.1098/rsif.2005.0083. URL [https://](https://royalsocietypublishing.org/doi/abs/10.1098/rsif.2005.0083)
583 royalsocietypublishing.org/doi/abs/10.1098/rsif.2005.0083.
- 584
585 Martin Danner, Matthias Begemann, Miriam Elbracht, Ingo Kurth, and Jeremias Krause. Utiliz-
586 ing protein structure graph embeddings to predict the pathogenicity of missense variants. *NAR*
587 *Genomics and Bioinformatics*, 7(3):lqaf097, 07 2025. ISSN 2631-9268. doi: 10.1093/nargab/
588 lqaf097. URL <https://doi.org/10.1093/nargab/lqaf097>.
- 589 J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky,
590 A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan,
591 B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, and D. Baker. Robust
592 deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56,
593 2022. doi: 10.1126/science.add2187. URL [https://www.science.org/doi/abs/10.](https://www.science.org/doi/abs/10.1126/science.add2187)
[1126/science.add2187](https://www.science.org/doi/abs/10.1126/science.add2187).

- 594 Nicki Skafte Detlefsen, Søren Hauberg, and Wouter Boomsma. Learning meaningful repre-
595 sentations of protein sequences. *Nature Communications*, 13(1):1914, Apr 2022. ISSN
596 2041-1723. doi: 10.1038/s41467-022-29443-w. URL [https://doi.org/10.1038/
597 s41467-022-29443-w](https://doi.org/10.1038/s41467-022-29443-w).
- 598 Jennifer L Fallon and Florante A Quioco. A closed compact structure of native ca²⁺-calmodulin.
599 *Structure*, 11(10):1303–1307, 2003. ISSN 0969-2126. doi: [https://doi.org/10.1016/j.str.
600 2003.09.004](https://doi.org/10.1016/j.str.2003.09.004). URL [https://www.sciencedirect.com/science/article/pii/
601 S0969212603002053](https://www.sciencedirect.com/science/article/pii/S0969212603002053).
- 602 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel
603 Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual
604 inversion, 2022. URL <https://arxiv.org/abs/2208.01618>.
- 606 Vladimir Gligorijević, P. Douglas Renfrew, Tomasz Kosciolk, Julia Koehler Leman, Daniel Beren-
607 berg, Tommi Vatanen, Chris Chandler, Bryn C. Taylor, Ian M. Fisk, Hera Vlamakis, Ramnik J.
608 Xavier, Rob Knight, Kyunghyun Cho, and Richard Bonneau. Structure-based protein function
609 prediction using graph convolutional networks. *Nature Communications*, 12(1):3168, May 2021.
610 ISSN 2041-1723. doi: 10.1038/s41467-021-23303-9. URL [https://doi.org/10.1038/
611 s41467-021-23303-9](https://doi.org/10.1038/s41467-021-23303-9).
- 612 Joe G Greener and Kiarash Jamali. Fast protein structure searching using structure graph embed-
613 dings. *Bioinformatics Advances*, 5(1):vbaf042, 03 2024. ISSN 2635-0041. doi: 10.1093/bioadv/
614 vbaf042. URL <https://doi.org/10.1093/bioadv/vbaf042>.
- 616 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in
617 neural information processing systems*, 33:6840–6851, 2020.
- 618 Jiahao Hu, Jiasheng Cao, Win Topatana, Sarun Juengpanich, Shijie Li, Bin Zhang, Jiliang Shen,
619 Liuxin Cai, Xiujun Cai, and Mingyu Chen. Targeting mutant p53 for cancer therapy: di-
620 rect and indirect strategies. *Journal of Hematology & Oncology*, 14(1):157, Sep 2021. ISSN
621 1756-8722. doi: 10.1186/s13045-021-01169-0. URL [https://doi.org/10.1186/
622 s13045-021-01169-0](https://doi.org/10.1186/s13045-021-01169-0).
- 623 Raza Jafri, Yash Raj, and Jacinta D’Souza. Changes in the adenylate kinase activity are proportional
624 to the adp/atp ratio upon resorption and regeneration of chlamydomonas reinhardtii flagella. *Cell
625 Biochemistry and Biophysics*, pp. 1–16, 07 2025. doi: 10.1007/s12013-025-01825-z.
- 626 Chen Jin, Ryutaro Tanno, Amrutha Saseendran, Tom Diethel, and Philip Alexander Teare. An
627 image is worth multiple words: Discovering object level concepts using multi-concept prompt
628 learning. In *Forty-first International Conference on Machine Learning*, 2024. URL [https://
629 //openreview.net/forum?id=F3x6uYILgL](https://openreview.net/forum?id=F3x6uYILgL).
- 630 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
631 Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, Alex Bridgland,
632 Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-
633 Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman,
634 Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Se-
635 bastian Bodenstern, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Push-
636 meet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with alphafold.
637 *Nature*, 596(7873):583–589, Aug 2021. ISSN 1476-4687. doi: 10.1038/s41586-021-03819-2.
638 URL <https://doi.org/10.1038/s41586-021-03819-2>.
- 639 Yogesh Kalakoti and Björn Wallner. Afsample2 predicts multiple conformations and ensembles
640 with alphafold2. *Communications Biology*, 8(1):373, Mar 2025. ISSN 2399-3642. doi: 10.1038/
641 s42003-025-07791-9. URL <https://doi.org/10.1038/s42003-025-07791-9>.
- 642 R. S. Kalinin, V. O. Shipunova, Y. P. Rubtsov, V. M. Ukrainskay, A. Schulga, E. V. Konovalova, D. V.
643 Volkov, I. A. Yaroshevich, A. M. Moysenovich, A. A. Belogurov, G. B. Telegin, A. S. Chernov,
644 M. A. Maschan, S. S. Terekhov, V. D. Knorre, E. Khurs, N. V. Gnuchev, A. G. Gabibov, and S. M.
645 Deyev. Barnase-barstar specific interaction regulates car-t cells cytotoxic activity toward malignancy.
646 *Doklady Biochemistry and Biophysics*, 508(1):17–20, Feb 2023. ISSN 1608-3091. doi: 10.
647 1134/S1607672922700041. URL <https://doi.org/10.1134/S1607672922700041>.

- 648 Shaun M Kandathil, Andy M Lau, Daniel W A Buchan, and David T Jones. Foldclass and merizo-
649 search: scalable structural similarity search for single- and multi-domain proteins using geometric
650 learning. *Bioinformatics*, 41(5):btaf277, 05 2025. ISSN 1367-4811. doi: 10.1093/bioinformatics/
651 btaf277. URL <https://doi.org/10.1093/bioinformatics/btaf277>.
- 652 Raisa Kantaev, Inbal Riven, Adi Goldenzweig, Yoav Barak, Orly Dym, Yoav Peleg, Shira Albeck,
653 Sarel J. Fleishman, and Gilad Haran. Manipulating the folding landscape of a multidomain pro-
654 tein. *The Journal of Physical Chemistry B*, 122(49):11030–11038, 2018. doi: 10.1021/acs.jpcc.
655 8b04834. URL <https://doi.org/10.1021/acs.jpcc.8b04834>. PMID: 30088929.
- 656
657 C. Klein and L. T. Vassilev. Targeting the p53–mdm2 interaction to treat cancer. *British Journal*
658 *of Cancer*, 91(8):1415–1419, Oct 2004. ISSN 1532-1827. doi: 10.1038/sj.bjc.6602164. URL
659 <https://doi.org/10.1038/sj.bjc.6602164>.
- 660 Andy M. Lau, Nicola Bordin, Shaun M. Kandathil, Ian Sillitoe, Vaishali P. Waman, Jude Wells,
661 Christine A. Orengo, and David T. Jones. Exploring structural diversity across the protein uni-
662 verse with the encyclopedia of domains. *Science*, 386(6721):eadq4946, 2024. doi: 10.1126/
663 science.adq4946. URL [https://www.science.org/doi/abs/10.1126/science.](https://www.science.org/doi/abs/10.1126/science.adq4946)
664 [adq4946](https://www.science.org/doi/abs/10.1126/science.adq4946).
- 665 Jinwoo Leem, Laura S. Mitchell, James H.R. Farmery, Justin Barton, and Jacob D. Galson. De-
666 ciphering the language of antibodies using self-supervised learning. *Patterns*, 3(7):100513,
667 2022. ISSN 2666-3899. doi: <https://doi.org/10.1016/j.patter.2022.100513>. URL <https://www.sciencedirect.com/science/article/pii/S2666389922001052>.
- 668
669 Julia Koehler Leman, Brian D Weitzner, Steven M Lewis, Jared Adolf-Bryfogle, Nawsad Alam, Re-
670becca F Alford, Melanie Aprahamian, David Baker, Kyle A Barlow, Patrick Barth, et al. Macro-
671 molecular modeling and design in rosetta: recent methods and frameworks. *Nature methods*, 17
672(7):665–680, 2020.
- 673 Chong Li, Marzena Pazgier, Changqing Li, Weirong Yuan, Min Liu, Gang Wei, Wei-Yue Lu,
674 and Wuyuan Lu. Systematic mutational analysis of peptide inhibition of the p53–mdm2/mdmx
675 interactions. *Journal of Molecular Biology*, 398(2):200–213, 2010. ISSN 0022-2836. doi:
676 <https://doi.org/10.1016/j.jmb.2010.03.005>. URL [https://www.sciencedirect.com/
677 science/article/pii/S0022283610002433](https://www.sciencedirect.com/science/article/pii/S0022283610002433).
- 678
679 Yeqing Lin, Minji Lee, Zhao Zhang, and Mohammed AlQuraishi. Out of many, one: Designing
680 and scaffolding proteins at the scale of the structural universe with genie 2, 2024. URL <https://arxiv.org/abs/2405.15489>.
- 681
682 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
683 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom
684 Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level
685 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/
686 science.ade2574. URL [https://www.science.org/doi/abs/10.1126/science.](https://www.science.org/doi/abs/10.1126/science.ade2574)
687 [ade2574](https://www.science.org/doi/abs/10.1126/science.ade2574).
- 688
689 Shitong Luo, Yufeng Su, Xingang Peng, Sheng Wang, Jian Peng, and Jianzhu Ma. Antigen-specific
690 antibody design and optimization with diffusion-based generative models for protein structures.
691 *Advances in Neural Information Processing Systems*, 35:9754–9767, 2022.
- 692 Daniel J. Mandell, Evangelos A. Coutsias, and Tanja Kortemme. Sub-angstrom accuracy in protein
693 loop reconstruction by robotics-inspired conformational sampling. *Nature Methods*, 6:551–552,
694 2009. doi: 10.1038/nmeth0809-551.
- 695
696 Kevin Michalewicz, Mauricio Barahona, and Barbara Bravi. Antipasti: Interpretable prediction of
697 antibody binding affinity exploiting normal modes and deep learning. *Structure*, 32(12):2422–
698 2434.e5, 2024. ISSN 0969-2126. doi: <https://doi.org/10.1016/j.str.2024.10.001>. URL <https://www.sciencedirect.com/science/article/pii/S0969212624004362>.
- 699
700 Kevin Michalewicz, Mauricio Barahona, and Barbara Bravi. Integrating protein sequence embed-
701 dings with structure via graph-based deep learning for the prediction of single-residue properties,
2025. URL <https://arxiv.org/abs/2502.17294>.

- 702 Saro Passaro, Gabriele Corso, Jeremy Wohlwend, Mateo Reveiz, Stephan Thaler, Vignesh Ram
703 Somnath, Noah Getz, Tally Portnoi, Julien Roy, Hannes Stark, David Kwabi-Addo, Dominique
704 Beaini, Tommi Jaakkola, and Regina Barzilay. Boltz-2: Towards accurate and efficient binding
705 affinity prediction. *bioRxiv*, 2025. doi: 10.1101/2025.06.14.659707.
- 706 Valerio Piomponi, Alberto Cazzaniga, and Francesca Cuturello. Evolutionary constraints guide
707 alphafold2 in predicting alternative conformations and inform rational mutation design. *Journal of*
708 *Chemical Information and Modeling*, 65(18):9459–9468, 2025. doi: 10.1021/acs.jcim.5c01090.
709 URL <https://doi.org/10.1021/acs.jcim.5c01090>. PMID: 40902999.
- 710 G.N. Ramachandran, C. Ramakrishnan, and V. Sasisekharan. Stereochemistry of polypeptide
711 chain configurations. *Journal of Molecular Biology*, 7(1):95–99, 1963. ISSN 0022-2836. doi:
712 [https://doi.org/10.1016/S0022-2836\(63\)80023-6](https://doi.org/10.1016/S0022-2836(63)80023-6). URL <https://www.sciencedirect.com/science/article/pii/S0022283663800236>.
- 713 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
714 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
715 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 716 Andrew W. Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green,
717 Chongli Qin, Augustin Židek, Alexander W. R. Nelson, Alex Bridgland, Hugo Penedones, Stig
718 Petersen, Karen Simonyan, Steve Crossan, Pushmeet Kohli, David T. Jones, David Silver, Ko-
719 ray Kavukcuoglu, and Demis Hassabis. Improved protein structure prediction using potentials
720 from deep learning. *Nature*, 577(7792):706–710, Jan 2020. ISSN 1476-4687. doi: 10.1038/
721 s41586-019-1923-7. URL <https://doi.org/10.1038/s41586-019-1923-7>.
- 722 Amelie Stein and Tanja Kortemme. Improvements to robotics-inspired conformational sampling in
723 rosetta. *PLoS ONE*, 8(5):e63090, 2013. doi: 10.1371/journal.pone.0063090.
- 724 Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and
725 Tommi S. Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-
726 scaffolding problem. In *The Eleventh International Conference on Learning Representations*,
727 2023. URL <https://openreview.net/forum?id=6TxBxqNME1Y>.
- 728 Christopher H. van Dyck. Anti-amyloid- β monoclonal antibodies for alzheimer’s disease: Pitfalls
729 and promise. *Biological Psychiatry*, 83(4):311–319, 2018. ISSN 0006-3223. doi: [https://doi.org/](https://doi.org/10.1016/j.biopsych.2017.08.010)
730 [10.1016/j.biopsych.2017.08.010](https://doi.org/10.1016/j.biopsych.2017.08.010). URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S0006322317318978)
731 [article/pii/S0006322317318978](https://www.sciencedirect.com/science/article/pii/S0006322317318978). Mechanisms of Alzheimer’s Disease and Treatment.
- 732 Anna Vangone and Alexandre MJJ Bonvin. Contacts-based prediction of binding affinity in pro-
733 tein–protein complexes. *eLife*, 4:e07454, jul 2015. ISSN 2050-084X. doi: 10.7554/eLife.07454.
734 URL <https://doi.org/10.7554/eLife.07454>.
- 735 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
736 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von
737 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-*
738 *vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
739 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
740 [file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 741 Susana Vázquez Torres, Philip J. Y. Leung, Preetham Venkatesh, Isaac D. Lutz, Fabian Hink, Huu-
742 Hien Huynh, Jessica Becker, Andy Hsien-Wei Yeh, David Juergens, Nathaniel R. Bennett, An-
743 drew N. Hoofnagle, Eric Huang, Michael J. MacCoss, Marc Expòsit, Gyu Rie Lee, Asim K.
744 Bera, Alex Kang, Joshmyn De La Cruz, Paul M. Levine, Xinting Li, Mila Lamb, Stacey R. Ger-
745 ben, Analisa Murray, Piper Heine, Elif Nihal Korkmaz, Jeff Nivala, Lance Stewart, Joseph L.
746 Watson, Joseph M. Rogers, and David Baker. De novo design of high-affinity binders of bioac-
747 tive helical peptides. *Nature*, 626(7998):435–442, Feb 2024. ISSN 1476-4687. doi: 10.1038/
748 s41586-023-06953-1. URL <https://doi.org/10.1038/s41586-023-06953-1>.
- 749 Susana Vázquez Torres, Melisa Benard Valle, Stephen P Mackessy, Stefanie K Menzies, Nicholas R
750 Casewell, Shirin Ahmadi, Nick J Burlet, Edin Muratspahić, Isaac Sappington, Max D Overath,
751 et al. De novo designed proteins neutralize lethal snake venom toxins. *Nature*, 639(8053):225–
752 231, 2025.

- 756 Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, Helen E.
757 Eisenach, Woody Ahern, Andrew J. Borst, Robert J. Ragotte, Lukas F. Milles, Basile I. M.
758 Wicky, Nikita Hanikel, Samuel J. Pellock, Alexis Courbet, William Sheffler, Jue Wang, Preetham
759 Venkatesh, Isaac Sappington, Susana Vázquez Torres, Anna Lauko, Valentin De Bortoli, Emile
760 Mathieu, Sergey Ovchinnikov, Regina Barzilay, Tommi S. Jaakkola, Frank DiMaio, Minkyung
761 Baek, and David Baker. De novo design of protein structure and function with rfdiffusion. *Nature*,
762 620(7976):1089–1100, Aug 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06415-8.
763 URL <https://doi.org/10.1038/s41586-023-06415-8>.
- 764 Jeremy Wohlwend, Gabriele Corso, Saro Passaro, Noah Getz, Mateo Reveiz, Ken Leidal, Wojtek
765 Swiderski, Liam Atkinson, Tally Portnoi, Itamar Chinn, Jacob Silterra, Tommi Jaakkola, and
766 Regina Barzilay. Boltz-1: Democratizing biomolecular interaction modeling. *bioRxiv*, 2024. doi:
767 10.1101/2024.11.19.624167.
- 768 Kevin E Wu, Kevin K Yang, Rianne van den Berg, Sarah Alamdari, James Y Zou, Alex X Lu, and
769 Ava P Amini. Protein structure generation via folding diffusion. *Nature communications*, 15(1):
770 1059, 2024.
- 771 Luhuan Wu, Brian L Trippe, Christian A. Naesseth, David M Blei, and John P Cunningham.
772 Practical and asymptotically exact conditional sampling in diffusion models. *arXiv preprint*
773 *arXiv:2306.17775*, 2023.
- 774 Li C. Xue, João Pglm Rodrigues, Panagiotis L. Kastritis, Alexandre Mjj Bonvin, and Anna Vangone.
775 Prodigy: a web server for predicting the binding affinity of protein–protein complexes. *Bioinform-*
776 *atics*, 32(23):3676–3678, 08 2016. ISSN 1367-4803. doi: 10.1093/bioinformatics/btw514.
777 URL <https://doi.org/10.1093/bioinformatics/btw514>.
- 778 Kevin K Yang, Zachary Wu, Claire N Bedbrook, and Frances H Arnold. Learned protein embed-
779 dings for machine learning. *Bioinformatics*, 34(15):2642–2648, 03 2018. ISSN 1367-4803. doi:
780 10.1093/bioinformatics/bty178. URL [https://doi.org/10.1093/bioinformatics/
781 bty178](https://doi.org/10.1093/bioinformatics/bty178).
- 782 Cheng Zhang, Adam Leach, Thomas Makkink, Miguel Arbesú, Ibtissem Kadri, Daniel Luo, Liron
783 Mizrahi, Sabrine Krichen, Maren Lang, Andrey Tovchigrechko, Nicolas Lopez Carranza, Uğur
784 Şahin, Karim Beguir, Michael Rooney, and Yunguan Fu. FrameDiPT: SE(3) Diffusion Model for
785 Protein Structure Inpainting. *bioRxiv*, 2023. doi: 10.1101/2023.11.21.568057. URL <https://www.biorxiv.org/content/10.1101/2023.11.21.568057v2>.
- 786 Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure
787 template quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4):702–710, 2004. doi:
788 <https://doi.org/10.1002/prot.20264>. URL [https://onlinelibrary.wiley.com/doi/
789 abs/10.1002/prot.20264](https://onlinelibrary.wiley.com/doi/abs/10.1002/prot.20264).
- 790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A REPRESENTATION AND CONDITIONING ARCHITECTURE

RFdiffusion, like RoseTTAFold, maintains three representation “tracks”: an MSA track v^{MSA} , a pair track v^{pair} , and a per-residue state track v^{state} (Baek et al., 2021; Watson et al., 2023). For a protein of length L , we write these as

$$v^{\text{MSA}} \in \mathbb{R}^{L \times d_{\text{MSA}}}, \quad v^{\text{pair}} \in \mathbb{R}^{L \times L \times d_{\text{pair}}}, \quad v^{\text{state}} \in \mathbb{R}^{L \times d_{\text{state}}}. \quad (6)$$

In all experiments we instantiate PGEL’s denoiser on the released RFdiffusion “Base” model (`Base_ckpt.pt`), whose public configuration specifies $d_{\text{state}} = 16$, $d_{\text{MSA}} = 256$ and $d_{\text{pair}} = 128$; we keep all model weights and these dimensions fixed.

We partition the sequence positions into motif residues R_* with $|R_*| = L_*$ and scaffold residues R_c with $|R_c| = L_c$, so that $L = L_c + L_*$. Running the frozen RFdiffusion encoder on the scaffold with sequence s_c and structure x_c yields scaffold embeddings

$$v_c^{\text{MSA}} \in \mathbb{R}^{L_c \times d_{\text{MSA}}}, \quad v_c^{\text{pair}} \in \mathbb{R}^{L_c \times L_c \times d_{\text{pair}}}, \quad v_c^{\text{state}} \in \mathbb{R}^{L_c \times d_{\text{state}}}. \quad (7)$$

PGEL introduces learnable motif embeddings

$$v_* = (v_*^{\text{MSA}}, v_*^{\text{pair}}, v_*^{\text{state}}), \quad (8)$$

where $v_*^{\text{MSA}} \in \mathbb{R}^{L_* \times d_{\text{MSA}}}$ and $v_*^{\text{state}} \in \mathbb{R}^{L_* \times d_{\text{state}}}$. $v_*^{\text{pair}} \in \mathbb{R}^{L_* \times L_* \times d_{\text{pair}}}$ collects all pairwise features involving at least one motif residue (motif-motif and motif-scaffold blocks). The only trainable parameters in PGEL are the entries of v_* ; all RFdiffusion parameters are frozen.

At each optimization step, we construct composite tracks that align with the RFdiffusion denoiser inputs by combining scaffold and motif embeddings. For the MSA and state tracks we concatenate along the sequence dimension:

$$v^{\text{MSA}} = \text{concat}(v_c^{\text{MSA}}, v_*^{\text{MSA}}) \in \mathbb{R}^{L \times d_{\text{MSA}}}, \quad v^{\text{state}} = \text{concat}(v_c^{\text{state}}, v_*^{\text{state}}) \in \mathbb{R}^{L \times d_{\text{state}}}. \quad (9)$$

The full pair tensor $v^{\text{pair}} \in \mathbb{R}^{L \times L \times d_{\text{pair}}}$ is obtained by keeping the scaffold-scaffold block equal to v_c^{pair} and inserting v_*^{pair} into all rows/columns corresponding to motif positions. The frozen RFdiffusion denoiser then takes the noisy coordinates $x(t)$ together with $(v^{\text{MSA}}, v^{\text{pair}}, v^{\text{state}})$ and produces a prediction $\hat{x}^{(0)}(v_*)$. We use the motif slice $\hat{x}_*^{(0)}(v_*)$ in the PGEL loss, while the scaffold coordinates are clamped to x_c when applying the reverse diffusion update from $x(t)$ to $x(t-1)$.

Information exchange between the learned motif embedding and the fixed scaffold embedding occurs entirely through the existing RFdiffusion blocks: global self-attention in the MSA track, MSA-pair and pair-state couplings, and the $SE(3)$ -equivariant refinement network. Gradients from the structural and torsion losses propagate through these layers to update v_* , while v_c and all network weights remain fixed.

B MASKING

Multiple-sequence-alignment (MSA) based structure predictors, such as AlphaFold, RoseTTAFold and RFdiffusion, derive a set of MSA embeddings that encode evolutionary covariation between residue positions and therefore impose strong geometric constraints on the predicted structure. In RFdiffusion these embeddings can be represented as a matrix $v_c \in \mathbb{R}^{d_{\text{MSA}} \times L_c}$, where d_{MSA} is the MSA embedding depth and L_c is the number of scaffold residues. Motivated by recent work that uses MSA column masking in AlphaFold-based pipelines to weaken coevolutionary constraints and increase conformational diversity (Kalakoti & Wallner, 2025; Pionponi et al., 2025), we introduce a simple masking operator $\mathcal{M}_{\omega, \alpha}(\cdot)$ that perturbs the scaffold MSA embedding during PGEL generation.

Specifically, given a masking type $\omega \sim \text{Ber}(\frac{1}{2})$ and a masking rate $\alpha \sim \mathcal{U}[0, 1]$ we construct a binary mask on v_c :

- Row masking ($\omega = 0$) selects a fraction α of rows of v_c (MSA embedding features) and sets them to zero for all scaffold residues. This reduces the effective dimensionality of the MSA feature space while preserving which residues are constrained.
- Column masking ($\omega = 1$) selects a fraction α of columns of v_c (scaffold residue positions) and sets all their MSA features to zero. This locally removes the evolutionary context of specific scaffold residues while leaving the remaining scaffold positions fully constrained.

In both cases, masking weakens the coevolutionary signals that couple the scaffold to the motif and thereby relaxes the induced geometric constraints. In our implementation, masking is applied only to the scaffold MSA embedding: the scaffold state and pair embeddings, as well as the learned motif embedding v_* , are left unmodified so any changes in the motif under different masks are mediated by the way EvoFormer attention blocks mix MSA, state, and pair embeddings. The frozen RFdiffusion denoiser therefore receives two conditioning signals at each reverse step: the perturbed scaffold MSA embedding $\mathcal{M}_{\omega, \alpha}(v_c)$ and the learned motif embedding v_* .

Empirically, we observe that column masking ($\omega = 1$) is substantially more effective than row masking in generating diverse yet designable motifs (around 80% of successful generations in Table 4 employ column masking). This is consistent with the interpretation above: while row masking removes global features that are shared across all residues and therefore tends to preserve the relative pattern of constraints, zeroing entire columns selectively weakens constraints arising from a subset of scaffold residues that are strongly coevolving with the motif and which enforce the native motif geometry. (This local relaxation allows the motif to move relative to the scaffold frame, while the rest of the scaffold remains tightly constrained by the unmasked columns). An ablation where we disable masking altogether ($\mathcal{M}_{\omega, \alpha} \equiv \text{Id}$) shows that, while PGEL still produces designable structures, the diversity collapses to a single TM-score cluster per system at threshold $t_m = 0.6$, confirming that masking is necessary to move beyond the native motif basin while maintaining designability (Appendix D).

C REVERSESTEP ALGORITHM

Let $x^{(t)} = \{(r_l^{(t)}, u_l^{(t)})\}_{l=1}^L$ denote the noisy protein backbone structure at diffusion step t , where each residue l is represented by a rotation $r_l^{(t)} \in SO(3)$, with $SO(3)$ the special orthogonal group in three dimensions, and a translation $u_l^{(t)} \in \mathbb{R}^3$. Let $\hat{x}^{(0)} = \{(\hat{r}_l^{(0)}, \hat{u}_l^{(0)})\}_{l=1}^L$ denote the predicted denoised structure. Let $\{\beta^{(t)}\}_{t=1}^T$ be a variance schedule with $\gamma^{(t)} = 1 - \beta^{(t)}$ and $\bar{\gamma}^{(t)} = \prod_{s=1}^t \gamma^{(s)}$. For translations, let $u_l^{(t-1)}$ be sampled from a Gaussian distribution with covariance $\beta^{(t)} I_3$. For rotations, let s_l denote the score approximation presented in Watson et al. (2023), $\epsilon_{l,d}$ isotropic Gaussian perturbations and $\{f_d\}_{d=1}^3$ a basis of the Lie algebra $SO(3)$.

Algorithm 3 REVERSESTEP function (Watson et al., 2023)

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Input: noisy structure $x^{(t)}$, denoised prediction $\hat{x}^{(0)}$.
Output: updated structure $x^{(t-1)}$.
for $l = 1, \dots, L$ **do**
 $(r_l^{(t)}, u_l^{(t)}) = x_l^{(t)}$
 $(\hat{r}_l^{(0)}, \hat{u}_l^{(0)}) = \hat{x}_l^{(0)}$
 $u_l^{(t-1)} \sim \mathcal{N}\left(\frac{\sqrt{\bar{\gamma}^{(t-1)}\beta^{(t)}}}{1-\bar{\gamma}^{(t)}}\hat{u}_l^{(0)} + \frac{\sqrt{\bar{\gamma}^{(t)}(1-\bar{\gamma}^{(t-1)})}}{1-\bar{\gamma}^{(t)}}u_l^{(t)}, \beta^{(t)}I_3\right)$
// Updating rotations below
 $s_l = \text{ROTATIONSCOREAPPROXIMATION}(r_l^{(t)}, \hat{r}_l^{(0)}, \sigma_t^2)$
 $\epsilon_{l,1}, \epsilon_{l,2}, \epsilon_{l,3} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$
 $r_l^{(t-1)} = r_l^{(t)} \exp_{I_3} \left\{ (\sigma_t^2 - \sigma_{t-1}^2) r_l^{(t)\top} s_l + \sqrt{\sigma_t^2 - \sigma_{t-1}^2} \sum_{d=1}^3 \epsilon_{l,d} f_d \right\}$
 $x_l^{(t-1)} = (r_l^{(t-1)}, u_l^{(t-1)})$
end for
Return $x^{(t-1)}$

D BASELINES

To demonstrate the necessity of both learned embeddings and zero masking, we consider simpler baselines in which one of these steps is omitted.

We first show that the presence of a fix scaffold, which informs the frozen denoiser via computed scaffold embeddings v_c (see Algorithm 1 and Appendix A), is insufficient to produce constraint-satisfying structures if the motif embeddings v_* are not learned. To test this, we skip Algorithm 1 and proceed directly to the generation with embedding masking (Algorithm 2). At each timestep, we now obtain the predicted motif structure $\hat{x}_*^{(0)}$ employing a masked version of the scaffold embeddings $\mathcal{M}_{\omega,\alpha}(v_c)$ and with v_* either set to zero or to random values drawn from a standard normal distribution. In both scenarios, we observe a total of 0 successful generations across all experiments. For zero embeddings, we obtain a mean mRMSD of 3.13Å and a mean pLDDT of 47 across the five examples, while these values are 3.48Å and 42, respectively, for random embeddings.

It could be argued that the stochasticity of RFdiffusion’s denoiser, even with frozen parameters, is sufficient to generate diversified motifs after embedding learning. Nevertheless, when computing at each timestep the predicted motif structure as $\hat{x}_*^{(0)} = \text{DENOISER}(v_c, v_*)$, only a single self-consistent cluster is obtained at $t_m = 0.6$ per experiment.

E MASKING RATE AND ITS IMPACT ON STRUCTURAL DIVERSITY

We studied the effect of the masking rate $\alpha \in [0, 1]$ on the resulting structural diversity, quantified by the TM-score between each generated motif and the native structure. We plotted $1 - \text{TM-score}$ as a function of α for all the experiments, obtaining moderate correlations as measured by the Pearson correlation score r with statistical significance (Figure 6).

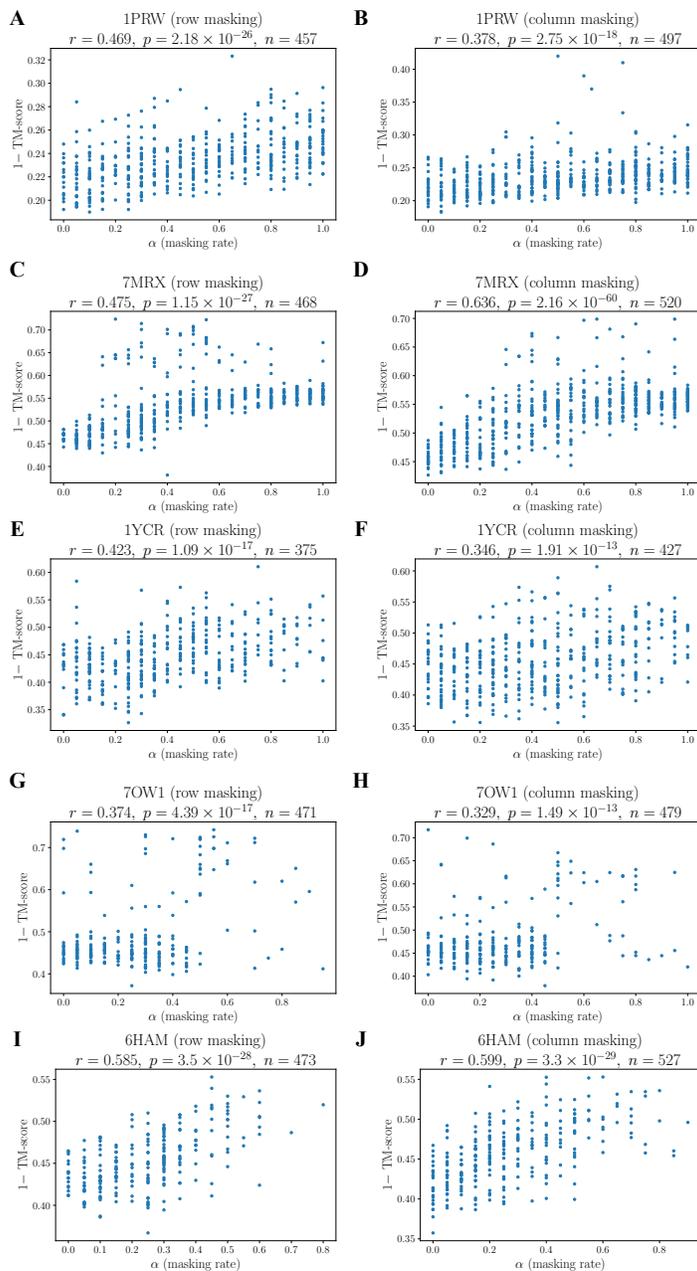


Figure 6: $1 - \text{TM-score}$ of each generated structure deemed designable (with respect to the native motif) as a function of the masking rate α , for row and column masking and five examples. For each case, we report the Pearson correlation coefficient r , its associated p -value, and the sample size.

F COMPUTATIONAL PERFORMANCE: STANDALONE vs. LINKED MOTIFS

We observed that the learning process described in Algorithm 1 requires more time when the designed motif is a standalone protein, *i.e.*, a full chain. To enable a fairer comparison, we expanded the original set of five cases by adding four additional examples. The results indicate that the computational performance appears largely independent of the motif length L_* (Table 2).

Table 2: Learning time comparison for standalone vs. linked motif designs.

PDB	Standalone	L_*	Learning time (min.)
1PRW	×	42	3.20
7MRX	×	22	2.06
7OW1	×	8	4.18
6HAM	×	22	2.48
1YCR	✓	13	118.78
7WT5	✓	8	34.87
3MZW	✓	58	23.41
1ATP	✓	20	97.55

G STUDY OF REGULARIZATION HYPERPARAMETERS

We studied which combinations of λ_{DM} and λ_{torsion} minimize the mRMSD at convergence of Algorithm 1 (see equation 2). Further, we examined the contribution of individual loss terms \mathcal{L}_{DM} and $\mathcal{L}_{\text{torsion}}$ by performing ablation experiments in which either the distance matrix term (equation 4) or the backbone torsion angle term (equation 5) were removed.

Although all combinations in Table 3 satisfy the mRMSD $< 2\text{\AA}$ constraint, the configurations with $\lambda_{\text{torsion}} = 0.05$ and $\lambda_{\text{DM}} \in \{0.01, 0.05, 0.1\}$, as well as $\lambda_{\text{torsion}} = \lambda_{\text{DM}} = 0.1$ appear to provide the best performance. While these results offer some guidance on the appropriate orders of magnitude of the regularization strengths, the full validation pipeline should be run for each case to assess the number of successes.

Regarding the ablation study, in the three examples the pLDDT values are no longer constraint-satisfying at convergence, with a mean value of 65 across the three experiments, when $\mathcal{L}_{\text{torsion}}$ is set to zero. When $\mathcal{L}_{\text{DM}} = 0$, mRMSD values are worse at convergence (as also reflected in the first column of Table 3, and the center of mass suffers a progressive deviation from that of the native structure as the iterations progress.

Table 3: Grid search over regularization strengths λ_{DM} and λ_{torsion} . Each cell reports the motif RMSD at convergence, averaged across three test cases (PDB IDs: 1PRW, 7MRX and 6HAM).

$\lambda_{\text{torsion}} \backslash \lambda_{\text{DM}}$	0.005	0.01	0.05	0.1	0.5
0.005	0.95 ± 0.32	0.99 ± 0.29	0.93 ± 0.19	0.94 ± 0.18	1.02 ± 0.21
0.01	1.82 ± 0.40	1.02 ± 0.29	0.77 ± 0.21	1.11 ± 0.30	1.04 ± 0.22
0.05	1.58 ± 0.37	0.98 ± 0.33	0.81 ± 0.24	0.96 ± 0.27	0.99 ± 0.24
0.1	1.03 ± 0.36	1.00 ± 0.31	0.78 ± 0.28	0.80 ± 0.29	0.83 ± 0.25
0.5	1.19 ± 0.42	1.15 ± 0.34	0.90 ± 0.17	0.87 ± 0.25	0.92 ± 0.22

H PGEL EVALUATION RESULTS

Table 4: Detailed results of PGEL successes for examples 1 to 5.

PDB & design ID	α	ω	mRMSD (Å)	Motif pLDDT	Sequence	mRMSD AF3 (Å)	pAE
1PRW							
17	0.75	1	1.78	79	FRVIAGGEDGLVLTLEQLARYVRRVAGRGGRLISFEDFLAI	1.54	4.43
860	0.5	1	1.96	81	ARWLDKGGSGAVFGEQLGEEVAAALEGGKEARLEEWFLNY	1.25	4.84
7MRX							
0	0.55	0	0.74	79	GLPESVSGNNQALYDSIMYDVE	0.89	3.17
31	0.4	0	1.35	78	GLPDTVKGLYAIGREAGYYYGD	0.83	3.30
100	0.95	1	1.26	73	GLPSTVTGLAGIGEDIRKGLLE	1.78	3.73
114	0.75	1	0.83	75	GASEGDAPAIYLAEDYCRYDLD	1.22	3.79
145	0.4	0	0.78	78	EPPEWVNGTLDAIYDGLIYTE	0.69	4.93
308	0.5	1	1.57	74	GIPESFKATTEAIGDWIRSNAD	1.25	4.97
352	0.85	1	1.53	75	GLPEEFTGNPYAIGEEAKRRLD	1.98	3.81
730	0.8	1	1.85	72	GIPEYLMGPDDESLLDWLKSLS	1.42	4.90
744	0.8	1	0.88	75	EVPGLEITGLSSLESTIRGYGS	1.96	2.42
814	0.3	1	0.77	78	GIKNLELSGLDAIKAAIDDLG	1.13	3.79
1YCR							
14	0.3	1	1.36	75	GSLEELLAQLDGG	1.56	2.88
36	0.7	1	1.09	72	MGLMELLGVPGS	1.25	3.19
285	0.3	1	1.59	72	ASGLELLKKFKLP	1.60	3.53
334	0.65	1	1.03	72	SSLMELLEKIDIE	1.25	2.88
619	0.2	1	0.59	87	ESFEELWKKIPGS	1.70	2.45
695	0.25	1	1.04	74	SSMWELWQEIEGE	0.86	2.42
7OW1							
136	0.2	1	1.92	70	SIDFGDGA	1.93	4.92
287	0.05	1	0.53	91	GSGGSLDT	0.62	3.54
637	0.35	0	1.00	77	GSRSGELV	0.77	3.59
837	0.15	1	1.32	77	GDNSGEAV	0.99	4.26
6HAM							
40	0.45	1	1.29	76	IATRAGGLGGLIREWAEKVG	1.09	4.08
64	0.15	1	0.80	87	ILLSRAGYLSREAARWISEKYG	1.40	4.33
104	0.15	0	1.23	71	IALTNAGWLDGLIAEFMKEKTG	1.88	3.45
115	0.05	1	1.10	86	ILLARPGSLGAEARHLSLTG	0.90	4.64
131	0.55	0	1.94	78	IALSRLVLDPELAQLMKELCG	1.98	3.71
257	0.25	1	1.28	71	IALERAGYRDRIKELGKELLG	1.40	4.14
313	0.25	1	1.78	71	IALSTKGGLSGLIADFAKEVLG	1.43	4.51
578	0.05	1	1.88	81	VLLHRPGADELARWLAKEVGG	0.83	3.75
653	0.6	0	1.85	74	ICLSRAGVFSGLFREIAEEFGK	0.89	3.88
658	0.1	1	1.16	81	IVLSRPGANGSVAREYAKEKLG	0.79	4.67
669	0.3	1	1.75	72	VFLSRNGFLHGLAREIAEELLG	1.46	4.38
694	0.5	0	1.89	73	IFLTRGGLLHGLVRELARELLK	1.53	4.31
716	0.05	1	1.08	81	IVLSRAGARSSEAGEWIKERLG	1.22	4.19
839	0.15	1	1.28	81	IALS RAGARS GEMGQIFKEITG	1.19	4.30
937	0.4	1	1.80	71	VLLSRGPYGSSLARELADYFGL	1.49	4.92

I BINDING AFFINITY RESULTS

Table 5: PRODIGY-predicted binding affinities for examples 2 and 3.

PDB & design ID	Method	ΔG (kcal/mol)
7MRX		
Native	-	-11.4
0	PGEL	-9.3
31	PGEL	-11.4
100	PGEL	-10.8
114	PGEL	-9.5
145	PGEL	-10.1
308	PGEL	-10.4
352	PGEL	-12.8
730	PGEL	-11.6
744	PGEL	-9.5
814	PGEL	-9.7
1	Partial diff.	-9.5
18	Partial diff.	-10.9
307	Partial diff.	-8.5
327	Partial diff.	-8.1
513	Partial diff.	-9.7
780	Partial diff.	-9.4
1YCR		
Native	-	-7.7
14	PGEL	-7.1
36	PGEL	-5.9
285	PGEL	-6.8
334	PGEL	-6.8
619	PGEL	-6.6
695	PGEL	-7.6
101	Partial diff.	-6.4

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187