MSAFLOW: A UNIFIED APPROACH FOR MSA REPRESENTATION, AUGMENTATION, AND FAMILY-BASED PROTEIN DESIGN

Anonymous authors Paper under double-blind review

000

002

003

006

008

009010011

012 013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Multiple Sequence Alignments (MSAs) provide fundamental information about protein evolution and play crucial roles in downstream applications, such as structure prediction and family-based design. However, constructing high-quality MSAs requires significant computational resources to query natural protein databases, and traditional techniques fail to retrieve sufficient data for proteins with limited homology. While recent generative models have been proposed for MSA augmentation, they often struggle to capture complex, high-order dependencies in sequence distributions while maintaining permutation invariance. To address these challenges, we introduce MSAFlow, a framework built on two key innovations. First, its core is a novel generative autoencoder that pairs a compressed AlphaFold3 (AF3) MSA representation with a conditional Statistical Flow Matching (SFM) decoder to faithfully model a family's sequence distribution that preserves permutation invariance. Second, we introduce a latent flow-matching model that performs zero-shot generation of MSA embeddings from a single sequence, enabling powerful augmentation for orphan proteins. By integrating these components, MSAFlow operates as a unified framework for MSA representation, augmentation, and family-based design. Our experiments demonstrate that MSAFlow significantly outperforms existing models on family-based protein design and MSA augmentation tasks, especially for low-homology proteins. MSAFlow is lightweight, fast, and memory-efficient, offering a single, versatile solution for diverse protein engineering tasks.

1 Introduction

Multiple Sequence Alignments (MSAs) provide fundamental information about protein evolution and play crucial roles in downstream tasks such as structure prediction and family-based sequence design (Gong et al., 2025; Truong Jr & Bepler, 2023; Chen et al., 2024; Zhang et al., 2024a; Cao et al., 2025). MSAs represent collections of homologous proteins that delineate the evolutionary history of a single query sequence, enabling identification of conserved regions (e.g., key active site residues for enzymes) and evolutionary couplings that inform three-dimensional structure.

Conventional homology search tools such as HHBlits (Remmert et al., 2012), MMSeqs (Steinegger & Söding, 2017), and JackHMMER (Johnson et al., 2010) often incur high computational costs for obtaining high-quality MSAs. More importantly, despite recent acceleration of MMSeqs2 with GPUs (Kallenborn et al., 2025), these methods fail to retrieve sufficient sequences for low-homology and orphan proteins. Therefore, tools that can generate MSAs and augment scarce evolution data are essential for expanding our capability to predict protein structure and functions. The challenge of MSA augmentation with generative tools has been partially addressed by Dense Homology Retriever (DHR) (Hong et al., 2024), which leverages pretrained embeddings from protein language models to identify homologous sequences more efficiently and with greater sensitivity. Several other models, including MSAGenerator (Zhang et al., 2024b), MSAGPT (Chen et al., 2024), and EvoDiff (Alamdari et al.), have subsequently emerged, employing autoregressive and discrete diffusion frameworks, respectively to generate MSAs. While they have shown promise in MSA generation and augmentation, they fail to capture rich distributional information while preserving permutation invariance. Furthermore, these methods typically utilize 2D positional encodings to represent row-

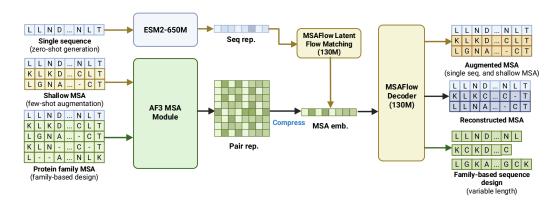


Figure 1: **General framework of MSAFlow.** Our approach supports three complementary pathways: (1) zero-shot generation from a single sequence using ESM2 embeddings, (2) few-shot augmentation of shallow MSAs, and (3) family-based design given MSAs embedded through the AF3 MSA Module and reconstructed through MSAFlow Decoder. All pathways leverage the latent flow-matching and decoder architecture to generate augmented or compressed MSAs, enabling both the enhancement of limited evolutionary information and the efficient representation of deep alignments.

wise and column-wise information present in MSAs. These approaches fail in critical aspects: they are substantially memory-intensive due to the $O(N^2)$ space complexity of self-attention operations, further exacerbated by the 2D nature of MSAs, and methods like MSAGPT lack true permutation invariance due to the left-to-right autoregressive decoding process.

On the other hand, better generative models for MSAs that capture higher-order evolutionary patterns can also serve as powerful tool for guiding functional protein design. Potts models (Seemayer et al., 2014) employ a pre-defined graphical model that is restricted to pairwise couplings. ProfileBFN (Gong et al., 2025) collapses sequence information into position-wise profiles that obscure higher-order dependencies, and methods such as MSA Transformer Rao et al. (2021) and EvoDiff (Alamdari et al.) flatten MSAs into 2D grids rather than explicitly modeling distributions over sequence space. This motivated us to develop a generative framework that can faithfully approximate the true sequence distribution within an MSA without imposing strong assumptions.

To address the limitations of the existing work, we introduce MSAFlow, a light-weight and effective framework that utilizes compressed latent MSA representations from AlphaFold3 (Abramson et al., 2024) (AF3) and conditional Statistical Flow Matching (Cheng et al., 2025) (SFM) as generative decoder to model the sequence distribution in an input MSA. Particularly, MSAFlow employs AF3's MSAModule as an encoder to produce pair presentations of MSAs, which is then mean-pooled and used as conditioning for the SFM decoder which is trained to reconstruct the original set of sequences in MSA (Figure 1). Unlike EVE (Frazer et al., 2021) that requires training separate VAE for each MSA, MSAFlow directly learns a generalizable generative auto-encoder over the space of MSAs (i.e. sets of sequences) with guaranteed permutation invariance. We further introduce a latent flow-matching model that generate MSA embeddings in a zero-shot manner from a single query sequence's ESM embedding. By learning from homology-rich MSA representations, our latent FM model can effectively augment proteins with shallow or no MSAs. By integrating these components, we introduce a unified end-to-end framework that is capable of MSA representation, augmentation and family-based protein design.

We summarize our contributions as follows:

- Novel architecture for modeling MSAs. We propose MSAFlow, an innovative generative autoencoding framework that operates on the space of sequence distributions. MSAFlow leverages
 compressed AF3 MSA embeddings to encode crucial evolutionary information, paired with a
 conditional Statistical Flow-matching decoder that reconstruct MSA sequences while maintaining
 permutation invariance.
- We enabled zero-shot generation of synthetic MSA with a two-stage approach combining a latent flow-matching over MSA embedding space and our MSAFlow decoder.

109

110

111

112

113

114

115

116

117 118 119

120 121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139 140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155 156 157

158

159

161

- We offer a unified framework for MSA representation, augmentation, and family-based sequence design. MSAFlow scales efficiently to large families, supports variable sequence lengths, and flexibly adapts to downstream design and analysis tasks—capabilities that prior models could not jointly achieve.
- Empirical significance. MSAFlow outperforms state-of-the-art models across multiple difficult protein structure prediction and family-based protein design tasks, including zero-shot and few-shot MSA generation for orphan and low-homology protein, and family-based enzyme design on EC classes with limited data. Notably, MSAFlow achieved better results despite being lightweight (130M parameters) and trained on smaller datasets, which offers better efficiency in terms of inference time and memory consumption (Table 5).

2 Additional Related Work

Generative models for protein sequences Protein sequence generative modeling can be approached from both discrete and continuous perspectives. Discrete protein language models—such as autoregressive transformers and masked language models—treat amino acid sequences as token sequences, learning residue distributions through maximum likelihood estimation or masked denoising objectives. The ProGen series (Madani et al., 2020; Bhatnagar et al., 2025) and ESM (Lin et al., 2023) are notable examples that employ Transformer architectures (Vaswani et al., 2017) to model residue-residue dependencies across vast protein families. Recent research has also explored discrete diffusion frameworks, such as EvoDiff (Alamdari et al., 2023), which learns denoising processes in amino acid token space, generating novel sequences with desired structural or functional properties through sequential unmasking. Complementarily, flow-matching methods like MultiFlow (Campbell et al., 2024) and FlowSeq (Ma et al., 2019) approach protein generation from continuous spaces. These continuous methods typically offer greater flexibility in conditional generation and interpolation but require decoding mechanisms to map continuous representations back to valid sequences. Furthermore, these language models can be applied to numerous downstream tasks, including protein-binding peptide design (D-Flow (Wu et al., 2024), PepFlow (Li et al., 2024)), structure-based sequence design (LM-Design (Zheng et al., 2023), InstructPLM (Qiu et al., 2024), DRAKES (Wang et al., 2024)), and antibody engineering (Frey et al.). These existing approaches largely focus on single-sequence modeling and often overlook evolutionary information contained in MSAs, limiting their capacity to capture residue co-variation and functional diversity that are essential for robust protein design.

Latent diffusion for protein design Latent diffusion models were initially applied to protein structure generation tasks, as demonstrated in several works (Fu et al., 2024; Zhang et al., 2025; Xu et al., 2023; Yim et al.), leveraging the inherent advantages of continuous representations. Recently, latent diffusion models have gained attention for modeling complex sequence-structure relationships through continuous embeddings. The construction of latent spaces enables greater consistency across multiple protein modalities while achieving concise and efficient representations. CHEAP (Lu et al., 2024) introduces a compressed hourglass representation of protein embeddings through VAE or VQ techniques, creating an efficient protein latent space. Building upon CHEAP's latent space, PLAID (Lu et al.) introduces a latent diffusion model over folding model internal embeddings, simultaneously decoding generated latents into sequence and structure to enable joint sequencestructure generation. Similarly, ProteinGenerator (Lisanza et al., 2024) conducts diffusion directly in sequence space while using a folding model (RoseTTAFold (Baek et al., 2021)) to guide generation, enabling structural constraints during the design process. However, despite these advances, prior approaches remain limited in their ability to operate over MSAs, which are critical for capturing evolutionary variation and residue—residue dependencies. This gap leaves open the opportunity for methods such as MSAFlow that extend latent diffusion principles to the MSA domain.

3 METHOD

3.1 MSAFLOW: AN AUTO-ENCODING FRAMEWORK FOR MSAS

MSAs are mathematically represented as $S = \{s_1, s_2, ..., s_M\}$ where each sequence $s_i \in A^L$ consists of amino acids and gaps from alphabet A, aligned to a reference sequence s_{ref} of length L. Despite

containing hundreds to thousands of sequences, we hypothesize that the functional and evolutionary information within an MSA can be **compressed into a continuous latent representation** that captures the essential characteristics of the sequence distribution within that protein family.

This compression necessitates a permutation-invariant encoding method to avoid bias from sequence ordering. Formally, we seek an encoder $h_\phi:\mathcal{S}\to\mathbb{R}^d$ such that $h_\phi(\mathcal{S})=h_\phi(\pi(\mathcal{S}))$ for any permutation π of the sequences in $\mathcal{S}.$ We leverage the AF3 MSAModule architecture, which provides a computationally efficient framework for embedding evolutionary information (Abramson et al., 2024). The AF3 MSAModule processes an MSA by



Figure 2: MSAFlow lifts autocoder to the space of sequence distributions within MSAs and families.

computing a position-wise outer product for each sequence s_i with the reference sequence, resulting in pairwise representations $P_i \in \mathbb{R}^{L \times L \times h_{\text{pair}}}$. These representations are averaged across all sequences as $P_{\text{avg}} = \frac{1}{M} \sum_{i=1}^{M} P_i$. The averaged representation is then processed through multiple triangle self-attention blocks to produce a refined pair representation $P_{\text{refined}} \in \mathbb{R}^{L \times L \times H}$. We utilize Protenix (Team et al., 2025), a pretrained variant of AF3, to generate these embeddings for MSAs from the OpenFold dataset (Ahdritz et al., 2024). The resulting pair representation serves as our compressed MSA embedding $m = h_{\phi}(\mathcal{S}) \in \mathbb{R}^{L \times L \times H}$.

Viewed this way, MSAFlow realizes an autoencoding framework over sets: the encoder maps the finite set of sequences in an MSA to a latent embedding, while the decoder reconstructs the underlying family-level distribution of sequences conditioned on this latent representation. This perspective emphasizes that MSAFlow does not simply compress individual sequences, but rather learns a compact representation of the set as a distribution, enabling permutation-invariant and family-aware generative modeling.

3.1.1 STATISTICAL FLOW MATCHING FOR MSA SEQUENCE DECODING

We formulate MSA decoding as a conditional generation task over the sequences within a protein family. Given an MSA $\mathcal S$ and its embedding $m=h_\phi(\mathcal S)$, the decoder reconstructs sequence distribution. Let $\tilde{\mathcal S}=\{s_1,\ldots,s_n\}$ be n sequences drawn uniformly without replacement from $\mathcal S$. We model $p_\theta(\tilde{\mathcal S}\mid m)=\prod_{i=1}^n p_\theta(s_i\mid m)$, which is permutation-invariant by construction. The decoder $p_\theta(s\mid m)$ represents the probability of sampling a sequence s compatible with m.

To instantiate $p_{\theta}(s \mid m)$ for discrete (categorical) sequences, we adopt Statistical Flow Matching (SFM) (Cheng et al., 2024), which learns a continuous Riemannian flow over the statistical manifold of categorical distributions equipped with Fisher-Rao metric. Concretely, each sequence in the MSA is treated as a sample of the target distribution. We operate in the probability simplex $\Delta^{|\mathcal{A}| \times L}$, where each position in the sequence is represented by a one-hot categorical distribution μ over amino acids.

Following SFM, we construct flow paths along geodesics on the positive orthant of the unit sphere by applying the mapping: $\pi: x=\pi(\mu)=\sqrt{\mu}$. SFM demonstrated that such a mapping to the unit sphere preserves the metric, which coincides with the canonical spherical geometry. Therefore, we can operate on the unit sphere with the standard spherical geometry. Mathematically, given a sequence s_i from the MSA and its corresponding categorical representation $x_1=\pi(\mu_1)$ (e.g., one-hot encoding) and the noise representation $x_0=\pi(\mu_0)$, the time-dependent interpolation follows:

$$x_t = \exp_{x_0}(t \cdot \log_{x_0}(x_1)) \tag{1}$$

where exp and log are the spherical exponential and logarithm maps on the manifold, respectively, and can be calculated in closed form as

$$\exp_x(u) = x \cos \|u\|_2 + \frac{u}{\|u\|_2} \sin \|u\|_2, \quad \log_x(y) = \frac{\arccos(\langle x, y \rangle)}{\sqrt{1 - \langle x, y \rangle^2}} (y - x - \langle x, y - x \rangle x), \quad (2)$$

After transforming back to the simplex with $\mu_t = \pi^{-1}(x_t)$, the interpolation in Equation 1 traces the geodesic between μ_0 and μ_1 with respect to the Fisher information metric, ensuring we follow the

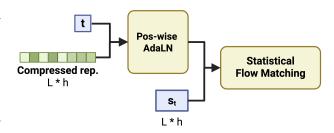
shortest path on the statistical manifold. The corresponding vector field for this mapped geodesic flow is given by $u_t(x_t|x_0,x_1)=\log_{x_t}(x_1)/(1-t)$. Instead of an unconditional model, our MSAFlow decoder employs a conditional parameterization where $v_{\theta}(x_t|m,t)$ is trained to approximate the vector field conditioning on the MSA embedding $m=h_{\phi}(\mathcal{S})$:

$$\mathcal{L}_{SFM}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1], s_i \sim \mathcal{S}, \mu_0 \sim \pi_* p_0, \mu_1 \sim \pi_* \delta(s_i)} \left[\| v_{\theta}(x_t | m, t) - u_t(x_t | x_0, x_1) \|^2 \right]$$
(3)

where π_* denotes the pushforward operation of applying the mapping π , x_t is obtained via the geodesic interpolation, and $\delta(s_i)$ represents the categorical distribution corresponding to sequence s_i (typically a one-hot encoding) in an MSA. During sampling, we first follow the learned marginal vector field on the sphere to obtain x_1 , then discrete generations of MSAs can be sampled from the categorical distribution $\mu_1 = \pi^{-1}(x_1)$.

3.1.2 Model Architecture and Implementation

We implement the vector field model v_{θ} using a modified conditional Diffusion Transformer (DiT) (Peebles & Xie, 2023) architecture. Since the output of the AF3 MSAModule is the pair representation of dimension $L \times L \times H$, we first compress it along the second dimension through mean pooling to obtain a sequence-level representation of dimension $L \times H$:



$$m_{\text{seq}} = \frac{1}{L} \sum_{j=1}^{L} m_{:,j,:} \in \mathbb{R}^{L \times H}$$
 (4)

Figure 3: DiT architecture for MSAFlow decoder.

This compressed representation serves as conditional information for the DiT model, which consists of 12 transformer blocks with a hidden dimension of 768, totaling approximately 130M parameters. The architecture incorporates sinusoidal time embeddings for the diffusion timestep t, token embeddings for each amino acid position, conditional embeddings from the compressed MSA representation, and multi-headed self-attention blocks with adaptive layer normalization. Notably, the MSA embedding conditioning is applied per-residue through a **position-wise AdaLN**, which introduces a novel mechanism for residue-level control. Unlike global conditioning schemes that broadcast the same modulation across all tokens, this design injects fine-grained, position-specific information into each layer normalization step, allowing for more precise alignment between evolutionary context and sequence generation. This innovation enhances the expressivity of the conditioning pathway and represents a new approach for leveraging MSAs in diffusion-based protein design. At inference time, we sample sequences by starting with random noise $x_1 \sim \text{Uniform}(\mathcal{A})$ and iteratively applying:

$$x_{t-\Delta t} = x_t - v_{\theta}(x_t|m, t) \cdot \Delta t \tag{5}$$

for timesteps $t = 1, 1 - \Delta t, 1 - 2\Delta t, ..., 0$, where Δt is a small step size (typically 0.01). At t = 0, we obtain the final sequence by taking the argmax over the amino acid probabilities at each position.

3.2 CONDITIONAL LATENT FLOW MATCHING FOR ZERO-SHOT MSA EMBEDDING GENERATION

While our decoder model generates sequences from MSA embeddings, we also develop a complementary approach to generate synthetic MSA embeddings themselves. This enables us to create artificial MSAs for proteins with limited evolutionary data. Let $z_1 = h_\phi(\mathcal{S}) \in \mathbb{R}^{L \times H}$ be the compressed MSA embedding for a reference sequence s_{ref} , and let $e = g_\psi(s_{\text{ref}}) \in \mathbb{R}^{d_e}$ be its ESM embedding. We aim to learn a conditional generative model $p_\theta(z_1|e)$ that can produce plausible MSA embeddings given only the reference sequence embedding.

Latent Flow Matching: We train a *conditional rectified flow* that maps a standard Gaussian $z_0 \sim \mathcal{N}(0,I)$ on the distribution of MSA embeddings $p(z\mid e)$ conditioned on the ESM embedding e (Lin et al., 2023). We use a straight-line path $z_t=(1-t)\,z_1+t\,z_0$ from target z_1 (the ground-truth MSA embedding) to noise z_0 , whose reference velocity is the constant field $u_t^\star(z_t;z_0,z_1)=z_0-z_1$. A time-dependent, conditional velocity $v_\theta(z_t,e,t)$ is learned by least-squares flow matching:

$$\mathcal{L}_{RFM} = \mathbb{E}_{t \sim \mathcal{U}[0,1], z_0 \sim \mathcal{N}(0,I), z_1} \| v_{\theta}(z_t, e, t) - (z_0 - z_1) \|_2^2,$$
(6)

which provides a simple, stable objective without explicit score estimation.

Generative Sampling Process: At inference, we draw $z_0 \sim \mathcal{N}(0, I)$ and integrate the learned conditional velocity backward from t=1 to t=0 with an explicit Euler solver. By default we use the deterministic probability-flow ODE (T=0); optionally, we add isotropic noise with temperature $T \in [0, 1]$ to trade fidelity for diversity:

$$z_{t-\Delta t} = z_t - v_\theta(z_t, e, t) \, \Delta t + T \sqrt{\Delta t} \, \varepsilon, \qquad \varepsilon \sim \mathcal{N}(0, I).$$
 (7)

Empirically, smaller T (e.g., T=0.5) improves alignment to e, while larger T increases sample diversity. Full SDE variants and discretization details are provided in Appendix 6.7.

3.3 END-TO-END UNIFIED PIPELINE FOR MSA REPRESENTATION, AUGMENTATION AND FAMILY-BASED SEQUENCE DESIGN

Our complete framework enables three complementary paths for MSA generation (as shown in Figure 1), each tailored to specific protein scenarios:

MSA Compression and Reconstruction: For deep MSAs with abundant evolutionary information, we first compress the multidimensional sequence information through the AF3 MSAModule into a compact latent representation. This compressed embedding effectively captures the evolutionary and functional signals present in the original MSA. We then use our SFM decoder to selectively reconstruct sequences, maintaining the key evolutionary characteristics while reducing redundancy.

Zero-shot MSA Augmentation: For orphan or de novo proteins with limited evolutionary data, we first generate the ESM embedding of the single available sequence. Our latent diffusion model then transforms this single-sequence representation into a synthetic MSA embedding that emulates the evolutionary diversity typically found in natural protein families. Finally, we decode multiple diverse sequences from this embedding using our SFM decoder, effectively bootstrapping evolutionary information where none previously existed.

Family-based Sequence Design: To perform family-based protein design, we first align all sequences belonging to the family (e.g., enzyme class) for a given query. These sequences are compressed into a latent representation using our AF3-based MSA encoder. Our SFM decoder then generates new sequences conditioned on this latent embedding, effectively producing new sequence designs that share a similar distribution to the given family. Because the generated sequences may include gaps, we can support both variable-length and fixed-length designs: gaps can be ignored when constructing the final sequence, enabling flexible design strategies.

This approach combines both MSA compression and generation capabilities in a unified framework. For data-rich scenarios, our method enables efficient information extraction from deep MSAs while preserving their evolutionary signals. For data-limited proteins, it allows the creation of synthetic alignments that capture potential evolutionary diversity. The integration of these complementary pathways addresses a fundamental limitation in protein analysis by extending evolutionary context to proteins that previously lacked sufficient homologous sequences, potentially improving downstream structure prediction, functional annotation tasks, and family-based design ability.

4 EXPERIMENTS

4.1 BENCHMARKING MSA AUTOENCODING

We evaluate the reconstruction ability of our model on 50 proteins released by CAMEO on May 10, 2025, where the ground truth MSA is generated using the same procedure as described in (Team et al., 2025). We took rigorous measures to avoid data leakage (maximum sequence identity from training set of 0.72, average 0.55) an ensured clear temporal separation between training and evaluation sets as described in Appendix 6.1. We then compute the embedding for each MSA using the AF3 MSAModule and generate 32 sequences for each latent MSA representation. We find that the relatively shallow MSAs generated by our model come close to matching the deep, ground-truth MSAs in terms of pLDDT (89.0 vs. 91.6) and TM-scores (0.86 vs. 0.89) while only consuming 6.5% of the overall bits required to represent a deep MSA (this is for an average sequence length of 365 and number of alignments being more than 7,000 from the CAMEO dataset). We perform conditional

generation given an embedding of 16-bit floats with an average size of 365×128 from the CAMEO dataset.

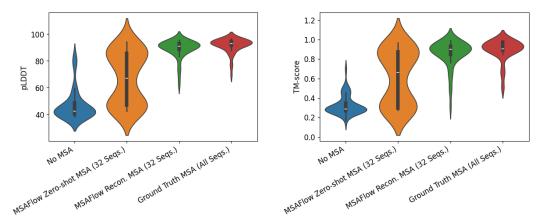


Figure 4: pLDDT and TM-scores for AF3 predictions of proteins from CAMEO with no MSA, MSAs generated through the MSAFlow-based zero-shot augmentation method, the MSAFlow-based reconstructed MSA (32 sequences), and the ground truth deep MSA (approximately 7k sequences).

Furthermore, when attempting to build synthetic MSA embedding (i.e. MSAs generated via our latent diffusion model), we find that our decoder is able to reconstruct some signal from the generated MSA latents, achieving much higher quality than without using an MSA altogether, although the structure prediction accuracy remains worse than using the ground truth embedding itself. Moreover, our model effectively compresses the heavy signal of full-depth evolutionary information encoded in thousands of aligned sequences into a single, fixed-size latent tensor that can be dynamically decoded into a range of sequences that remain evolutionarily related to the query, further evidenced in Table 7. As a result, we keep almost all of the functional signal that matters for folding accuracy.

We further evaluate the intrinsic quality of generated MSAs by comparing their residue-level entropy statistics to ground truth alignments. Following the evaluation setup in Zhang et al. (2024b), we generate 1000 sequences per MSA for our CAMEO test set and compute per-position entropies. MSAFlow's alignments closely mirror ground-truth entropy profiles, with generated sequences centered almost exactly on the true distribution (average entropy difference of 0.076 vs. 0.136 for ProfileBFN). Moreover, residue-level conservation patterns are preserved with high fidelity, as reflected by a markedly lower variance (0.294 vs. 0.724), demonstrating that MSAFlow achieves unsupervised alignment quality substantially closer to ground truth statistics.

Table 1: Comparison of entropy statistics between generated MSAs and ground truth (GT). MSAFlow more accurately recapitulates GT distributions, with lower entropy differences and variance.

	MSAFlow	ProfileBFN	GT
Average entropy	2.755 ± 0.294	2.838 ± 0.724	2.68 ± 0.589
Average entropy difference from GT	0.076	0.136	

4.2 Augmenting shallow and single-sequence MSAs

We further evaluate our model on a dataset of sequences with limited evolutionary information derived from MSAGPT (Chen et al., 2024), which includes 200 proteins from CAMEO (Haas et al., 2018), CASP14, CASP15, and PDB (Berman et al., 2000) with either few or no sequences in their MSA (few-shot and zero-shot cases, respectively). For the zero-shot case, we embed the query sequence with ESM and use it as conditioning for our latent diffusion model, which generates a synthetic MSA embedding for the reference sequence. We generate embeddings using 10 different seeds and employ low-temperature sampling during the SDE forward pass for higher-fidelity reconstructions, as detailed in (Geffner et al., 2025b). We then decode 32 sequences from each of the 10 synthetic MSA embeddings and report the best pLDDT and TM-scores. We find that our model significantly

outperforms prior state-of-the-art MSA augmentation tools, which also yield poorer results when evaluated using AF3.

Table 2: The accuracy of MSAFlow-generated multiple sequence alignments compared to other state-of-the-art methods, as evaluated by AlphaFold3 protein structure prediction performance on a naturally scarce MSA dataset curated from CAMEO, PDB, and CASP.

	AF3 pLDDT		TM-score	
	Zero-shot	Few-shot	Zero-shot	Few-shot
No/Shallow MSA	73.1	70.8	0.55	0.58
EvoDiff (650M)	67.7	67.5	0.49	0.55
MSAGPT (3B)	71.6	70.3	0.53	0.58
ESMFold	_	_	0.58	-
MSAFlow (Ours,130M)	75.2	70.4	0.62	0.60

For the few-shot augmentation case, we use our latent flow matching model to generate synthetic embeddings for each sequence over 5 different seeds, and decode 32 sequences from each MSA embedding. We then decode 64 sequences from the ground-truth shallow MSA embedding and extract the 16 most diverse sequences across all generations, following Chen et al. (2024). We concatenate our generated sequences with the original shallow MSA and find that our model improves upon structure prediction accuracy for such cases. We detail ablations motivating this reconstruction and augmentation scheme in Appendix 6.5 and 6.6.

4.3 CASE STUDIES ON de novo AND INTRINSICALLY DISORDERED PROTEINS

We show that MSAFlow markedly improves structure prediction for notoriously difficult proteins by generating high-quality synthetic MSAs. We focus on three challenging cases from a sparse MSA dataset:

- **8B4K**: the N-terminal domain of Rfa1 complexed with a phosphorylated Ddc2 peptide—only 133 residues, with scarce evolutionary relatives.
- **8G8I**: a Rosetta-designed four-helix bundle with rigid backbone constraints, extraordinary thermal stability $(T_m > 90^{\circ}\text{C})$, and NMR-validated topology (backbone RMSD = 1.11 Å).
- **80KH**: the crystal structure of *Bdellovibrio bacteriovorus* Bd1399.

MSAFlow's synthetic MSAs significantly outperform both MSA-free predictions and those using MSAGPT, which lacks sufficiently precise coevolutionary signals. This highlights MSAFlow's strengths in addressing two key failure modes: (i) limited sequence homology and (ii) intrinsically flexible or disordered regions—by synthesizing information-rich, high-fidelity MSAs in latent space that modern folding models require. We provide further such examples of case studies in Appendix 6.2.

4.4 FAMILY-BASED PROTEIN DESIGN

To better demonstrate the strength of MSAFlow on few-shot generation and generalization to other downstream applications than AF3 prediction, we now provide new results on family-based enzyme design. Our experiments demonstrate clear and significant advantages of MSAFlow, particularly for EC classes with limited sequences. Following ProfileBFN (Gong et al., 2025), we generate sequences in a single shot using our model, for enzymes with less than 20 sequences in their corresponding EC class, using the sequences from the EC class as an MSA. We then use CLEAN (Yu et al., 2023) to determine their EC number, and compute the accuracy (i.e. how many generated designs match the ground truth EC number) and the uniqueness across all generated designs. We report the accuracy × uniqueness score as done by ProfileBFN, the current SOTA for this task. MSAFlow exhibits SOTA performance on family-based enzyme design in both fixed and variable length settings. Notably, ProfileBFN is confined to fixed-length generation, whereas MSAFlow learns a meaningful homology distribution that guides the placement of gaps, which effectively enables variable-length design with unprecedented success rate.

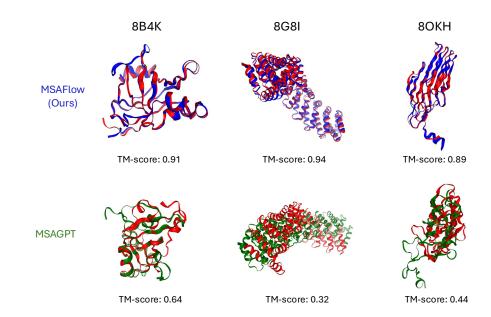


Figure 5: Visualization of improved structure prediction for zero-shot augmentation on de novo and disordered proteins with MSAFlow decoded synthetic MSAs, as compared to MSAs generated with MSAGPT. Blue represents predictions with an MSAFlow-generated MSA and green represents predictions with an MSAGPT-generated MSA. Red indicates the ground truth structure.

Table 3: Performance comparison of MSAFlow with baseline methods on family-based enzyme design task across different EC classes.

		Q15I65	Q15BH7	P13280	P57298
MSA Depth		15	12	13	15
# of Generated Sequences		1000	100	100	100
Acc. × Uniqueness (Fixed Length)	EvoDiff	1.39% (Gong et al., 2025)	0%	80%	5%
	ProfileBFN	42.67% (Gong et al., 2025)	89%	100%	82%
	MSAFlow	83.10%	84%	100%	95 %
Acc. × Uniqueness (Variable Length)	EvoDiff	-	0%	0%	0%
	MSAGPT	-	35.59%	37.5%	24.98%
	MSAFlow	-	92 %	92%	84%

5 CONCLUSION

MSAFlow integrates statistical flow matching with latent space optimization to enable bidirectional manipulation of multiple sequence alignments. By combining AlphaFold3-inspired permutation-equivariant embeddings with diffusion-based generation, it uniquely achieves both evolutionary signal compression and biologically plausible augmentation of sparse alignments. Comprehensive benchmarking across critical applications—latent space reconstruction fidelity, shallow MSA augmentation, synthetic alignment generation, and enzyme design—demonstrates MSAFlow's superiority, achieving state-of-the-art performance with only 130M parameters. MSAFlow's ability to generate evolutionarily coherent sequence ensembles creates new opportunities for designing orphan proteins and tackling de novo structure prediction challenges. Importantly, our framework also enables family-based design, where latent representations distilled from enzyme or protein families can guide the generation of sequences that remain faithful to family-level constraints while still exploring novel sequence diversity. Overall, MSAFlow advances both computational efficiency and conceptual modeling of protein sequence spaces through flow-based generation, paving the way for conditional protein engineering, resource-efficient applications, and family-level design of functional proteins.

ETHICS AND REPRODUCIBILITY STATEMENT

We have taken several steps to ensure reproducibility of our findings. The full model description, including encoder, decoder, and flow-matching components, is detailed in Section 3. Hyperparameters, training/test splits, and dataset sources are provided in Section 6.1. Ablation studies (Section 6.6, Table 9) clarify the contributions of different components, and additional case studies (Section 6.2, Table 4) demonstrate robustness across diverse proteins. Experimental comparisons with baselines are presented in Section 4 (Tables 1–3). To further facilitate reproducibility, we plan to release the anonymized source code and trained models in the supplementary material soon.

REFERENCES

- Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, Sebastian W Bodenstein, David A Evans, Chia-Chun Hung, Michael O'Neill, David Reiman, Kathryn Tunyasuvunakool, Zachary Wu, Akvilė Žemgulytė, Eirini Arvaniti, Charles Beattie, Ottavia Bertolli, Alex Bridgland, Alexey Cherepanov, Miles Congreve, Alexander I Cowen-Rivers, Andrew Cowie, Michael Figurnov, Fabian B Fuchs, Hannah Gladman, Rishub Jain, Yousuf A Khan, Caroline M R Low, Kuba Perlin, Anna Potapenko, Pascal Savy, Sukhdeep Singh, Adrian Stecula, Ashok Thillaisundaram, Catherine Tong, Sergei Yakneen, Ellen D Zhong, Michal Zielinski, Augustin Žídek, Victor Bapst, Pushmeet Kohli, Max Jaderberg, Demis Hassabis, and John M Jumper. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature*, 630(8016):493–500, June 2024.
- Gustaf Ahdritz, Nazim Bouatta, Christina Floristean, Sachin Kadyan, Qinghui Xia, William Gerecke, Timothy J. O'Donnell, Daniel Berenberg, Ian Fisk, Niccolò Zanichelli, Bo Zhang, Arkadiusz Nowaczynski, Bei Wang, Marta M. Stepniewska-Dziubinska, Shang Zhang, Adegoke Ojewole, Murat Efe Guney, Stella Biderman, Andrew M. Watkins, Stephen Ra, Pablo Ribalta Lorenzo, Lucas Nivon, Brian Weitzner, Yih-En Andrew Ban, Shiyang Chen, Minjia Zhang, Conglong Li, Shuaiwen Leon Song, Yuxiong He, Peter K. Sorger, Emad Mostaque, Zhao Zhang, Richard Bonneau, and Mohammed AlQuraishi. Openfold: retraining alphafold2 yields new insights into its learning mechanisms and capacity for generalization. *Nature Methods*, 21(8):1514–1524, May 2024. ISSN 1548-7105. doi: 10.1038/s41592-024-02272-z. URL http://dx.doi.org/10.1038/s41592-024-02272-z.
- Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex X. Lu, Nicolo Fusi, Ava P. Amini, and Kevin K. Yang. Protein generation with evolutionary diffusion: sequence is all you need. doi: 10.1101/2023.09.11.556673. URL https://www.biorxiv.org/content/early/2023/09/12/2023.09.11.556673. Publisher: Cold Spring Harbor Laboratory _eprint: https://www.biorxiv.org/content/early/2023/09/12/2023.09.11.556673.full.pdf.
- Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Neil Tenenholtz, Bob Strome, Alan Moses, Alex Xijie Lu, Nicolo Fusi, Ava Pardis Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *BioRxiv*, pp. 2023–09, 2023.
- Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, Sergey Ovchinnikov, Gyu Rie Lee, Jue Wang, Qian Cong, Lisa N Kinch, R Dustin Schaeffer, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science*, 373(6557):871–876, 2021.
- H M Berman, J Westbrook, Z Feng, G Gilliland, T N Bhat, H Weissig, I N Shindyalov, and P E Bourne. The protein data bank. *Nucleic Acids Res.*, 28(1):235–242, January 2000.
- Aadyot Bhatnagar, Sarthak Jain, Joel Beazer, Samuel C Curran, Alexander M Hoffnagle, Kyle Ching, Michael Martyn, Stephen Nayfach, Jeffrey A Ruffolo, and Ali Madani. Scaling unlocks broader generation and deeper functional understanding of proteins. *bioRxiv*, pp. 2025–04, 2025.
- Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design. arXiv preprint arXiv:2402.04997, 2024.
- Hanqun Cao, Xinyi Zhou, Zijun Gao, Chenyu Wang, Xin Gao, Zhi Zhang, Chunbin Gu, Ge Liu, and Pheng-Ann Heng. Plame: Leveraging pretrained language models to generate enhanced protein multiple sequence alignments. *arXiv preprint arXiv:2507.07032*, 2025.

- Bo Chen, Zhilei Bei, Xingyi Cheng, Pan Li, Jie Tang, and Le Song. Msagpt: Neural prompting protein structure prediction via msa generative pre-training, 2024. URL https://arxiv.org/abs/2406.05347.
- Chaoran Cheng, Jiahan Li, Jian Peng, and Ge Liu. Categorical flow matching on statistical manifolds.
 arXiv preprint arXiv:2405.16441, 2024.
 - Chaoran Cheng, Jiahan Li, Jian Peng, and Ge Liu. Categorical flow matching on statistical manifolds, 2025. URL https://arxiv.org/abs/2405.16441.
 - Jonathan Frazer, Pascal Notin, Mafalda Dias, Aidan Gomez, Joseph K Min, Kelly Brock, Yarin Gal, and Debora S Marks. Disease variant prediction with deep generative models of evolutionary data. *Nature*, 599(7883):91–95, November 2021.
 - Nathan C Frey, Dan Berenberg, Karina Zadorozhny, Joseph Kleinhenz, Julien Lafrance-Vanasse, Isidro Hotzel, Yan Wu, Stephen Ra, Richard Bonneau, Kyunghyun Cho, et al. Protein discovery with discrete walk-jump sampling. In *The Twelfth International Conference on Learning Representations*.
 - Cong Fu, Keqiang Yan, Limei Wang, Wing Yee Au, Michael Curtis McThrow, Tao Komikado, Koji Maruhashi, Kanji Uchino, Xiaoning Qian, and Shuiwang Ji. A latent diffusion model for protein structure generation. In *Learning on Graphs Conference*, pp. 29–1. PMLR, 2024.
 - Tomas Geffner, Kieran Didi, Zuobai Zhang, Danny Reidenbach, Zhonglin Cao, Jason Yim, Mario Geiger, Christian Dallago, Emine Kucukbenli, Arash Vahdat, and Karsten Kreis. Proteina: Scaling flow-based protein structure generative models, 2025a. URL https://arxiv.org/abs/2503.00710.
 - Tomas Geffner, Kieran Didi, Zuobai Zhang, Danny Reidenbach, Zhonglin Cao, Jason Yim, Mario Geiger, Christian Dallago, Emine Kucukbenli, Arash Vahdat, and Karsten Kreis. Proteina: Scaling flow-based protein structure generative models, 2025b. URL https://arxiv.org/abs/2503.00710.
 - Jingjing Gong, Yu Pei, Siyu Long, Yuxuan Song, Zhe Zhang, Wenhao Huang, Ziyao Cao, Shuyi Zhang, Hao Zhou, and Wei-Ying Ma. Steering protein family design through profile bayesian flow. February 2025.
 - Jürgen Haas, Alessandro Barbato, Dario Behringer, Gabriel Studer, Steven Roth, Martino Bertoni, Khaled Mostaguir, Rafal Gumienny, and Torsten Schwede. Continuous automated model Evaluation (CAMEO) complementing the critical assessment of structure prediction in CASP12. *Proteins*, 86:387–398, March 2018.
 - Liang Hong, Zhihang Hu, Siqi Sun, Xiangru Tang, Jiuming Wang, Qingxiong Tan, Liangzhen Zheng, Sheng Wang, Sheng Xu, Irwin King, Mark Gerstein, and Yu Li. Fast, sensitive detection of protein homologs using deep dense retrieval. *Nat. Biotechnol.*, pp. 1–13, August 2024. ISSN 1546-1696. doi: 10.1038/s41587-024-02353-6.
 - L. Steven Johnson, Sean R. Eddy, and Elon Portugaly. Hidden Markov model speed heuristic and iterative HMM search procedure. *BMC Bioinformatics*, 11(1):431, August 2010. ISSN 1471-2105. doi: 10.1186/1471-2105-11-431. URL https://doi.org/10.1186/1471-2105-11-431.
 - Felix Kallenborn, Alejandro Chacon, Christian Hundt, Hassan Sirelkhatim, Kieran Didi, Sooyoung Cha, Christian Dallago, Milot Mirdita, Bertil Schmidt, and Martin Steinegger. Gpu-accelerated homology search with mmseqs2. *Nature Methods*, 2025. doi: 10.1038/s41592-025-02819-8. URL https://doi.org/10.1038/s41592-025-02819-8.
 - Jiahan Li, Chaoran Cheng, Zuofan Wu, Ruihan Guo, Shitong Luo, Zhizhou Ren, Jian Peng, and Jianzhu Ma. Full-atom peptide design based on multi-modal flow matching. *arXiv* preprint *arXiv*:2406.00735, 2024.

Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/science.ade2574. URL https://www.science.org/doi/abs/10.1126/science.ade2574.

- Sidney Lyayuga Lisanza, Jacob Merle Gershon, Samuel WK Tipps, Jeremiah Nelson Sims, Lucas Arnoldt, Samuel J Hendel, Miriam K Simma, Ge Liu, Muna Yase, Hongwei Wu, et al. Multistate and functional protein design using rosettafold sequence space diffusion. *Nature biotechnology*, pp. 1–11, 2024.
- Amy X Lu, Wilson Yan, Sarah A Robinson, Simon Kelow, Kevin K Yang, Vladimir Gligorijevic, Kyunghyun Cho, Richard Bonneau, Pieter Abbeel, and Nathan C Frey. All-atom protein generation with latent diffusion. In *ICLR 2025 Workshop on Generative and Experimental Perspectives for Biomolecular Design*.
- Amy X Lu, Wilson Yan, Kevin K Yang, Vladimir Gligorijevic, Kyunghyun Cho, Pieter Abbeel, Richard Bonneau, and Nathan Frey. Tokenized and continuous embedding compressions of protein sequence and structure. *bioRxiv*, pp. 2024–08, 2024.
- Xuezhe Ma, Chunting Zhou, Xian Li, Graham Neubig, and Eduard Hovy. Flowseq: Non-autoregressive conditional sequence generation with generative flow, 2019. URL https://arxiv.org/abs/1909.02480.
- Ali Madani, Bryan McCann, Nikhil Naik, Nitish Shirish Keskar, Namrata Anand, Raphael R. Eguchi, Po-Ssu Huang, and Richard Socher. Progen: Language modeling for protein generation, 2020. URL https://arxiv.org/abs/2004.03497.
- William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL https://arxiv.org/abs/2212.09748.
- Jiezhong Qiu, Junde Xu, Jie Hu, Hanqun Cao, Liya Hou, Zijun Gao, Xinyi Zhou, Anni Li, Xiujuan Li, Bin Cui, et al. Instructplm: Aligning protein language models to follow protein structure instructions. *bioRxiv*, pp. 2024–04, 2024.
- Roshan Rao, Jason Liu, Robert Verkuil, Joshua Meier, John F. Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. Msa transformer. bioRxiv, 2021. doi: 10.1101/2021.02.12.430858. URL https://www.biorxiv.org/content/early/2021/02/13/2021.02.12.430858.
- Michael Remmert, Andreas Biegert, Andreas Hauser, and Johannes Söding. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nature Methods*, 9(2):173–175, February 2012. ISSN 1548-7091, 1548-7105. doi: 10.1038/nmeth.1818. URL https://www.nature.com/articles/nmeth.1818.
- Stefan Seemayer, Markus Gruber, and Johannes Söding. Ccmpred—fast and precise prediction of protein residue–residue contacts from correlated mutations. *Bioinformatics*, 30(21):3128–3130, 07 2014. ISSN 1367-4803. doi: 10.1093/bioinformatics/btu500. URL https://doi.org/10.1093/bioinformatics/btu500.
- Martin Steinegger and Johannes Söding. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature Biotechnology*, 35(11):1026–1028, November 2017. ISSN 1087-0156, 1546-1696. doi: 10.1038/nbt.3988. URL https://www.nature.com/articles/nbt.3988.
- ByteDance AML AI4Science Team, Xinshi Chen, Yuxuan Zhang, Chan Lu, Wenzhi Ma, Jiaqi Guan, Chengyue Gong, Jincai Yang, Hanyu Zhang, Ke Zhang, Shenghao Wu, Kuangqi Zhou, Yanping Yang, Zhenyu Liu, Lan Wang, Bo Shi, Shaochen Shi, and Wenzhi Xiao. Protenix -advancing structure prediction through a comprehensive alphafold3 reproduction. *bioRxiv*, 2025. doi: 10.1101/2025.01.08.631967. URL https://www.biorxiv.org/content/early/2025/01/11/2025.01.08.631967.

- Timothy Truong Jr and Tristan Bepler. Poet: A generative model of protein families as sequences-of-sequences. *Advances in Neural Information Processing Systems*, 36:77379–77415, 2023.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
 - Chenyu Wang, Masatoshi Uehara, Yichun He, Amy Wang, Tommaso Biancalani, Avantika Lal, Tommi Jaakkola, Sergey Levine, Hanchen Wang, and Aviv Regev. Fine-tuning discrete diffusion models via reward optimization with applications to dna and protein design. *arXiv preprint arXiv:2410.13643*, 2024.
 - Fang Wu, Tinson Xu, Shuting Jin, Xiangru Tang, Zerui Xu, James Zou, and Brian Hie. D-flow: Multi-modality flow matching for d-peptide design. *arXiv preprint arXiv:2411.10618*, 2024.
 - Minkai Xu, Alexander S Powers, Ron O Dror, Stefano Ermon, and Jure Leskovec. Geometric latent diffusion models for 3d molecule generation. In *International Conference on Machine Learning*, pp. 38592–38610. PMLR, 2023.
 - Jason Yim, Marouane Jaakik, Ge Liu, Jacob Gershon, Karsten Kreis, David Baker, Regina Barzilay, and Tommi Jaakkola. Hierarchical protein backbone generation with latent and structure diffusion.
 - Tianhao Yu, Haiyang Cui, Jianan Canal Li, Yunan Luo, Guangde Jiang, and Huimin Zhao. Enzyme function prediction using contrastive learning. *Science*, 379(6639):1358–1363, March 2023.
 - Le Zhang, Jiayang Chen, Tao Shen, Yu Li, and Siqi Sun. Msa generation with seqs2seqs pretraining: Advancing protein structure predictions. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 57324–57348. Curran Associates, Inc., 2024a. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/694be3548697e9cc8999d45e8d16fele-Paper-Conference.pdf.
 - Le Zhang, Jiayang Chen, Tao Shen, Yu Li, and Siqi Sun. MSA generation with seqs2seqs pretraining: Advancing protein structure predictions, 2024b. URL https://openreview.net/forum?id=bM6LUC2lec.
 - Rongchao Zhang, Yu Huang, Yiwei Lou, Yi Xin, Haixu Chen, Yongzhi Cao, and Hanpin Wang. Exploit your latents: Coarse-grained protein backmapping with latent diffusion models. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 39, pp. 1111–1119, 2025.
 - Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei Ye, and Quanquan Gu. Structure-informed language models are protein designers. In *International conference on machine learning*, pp. 42317–42338. PMLR, 2023.

6 Additional results

6.1 DETAILS ON MSAFLOW TRAIN/TEST SPLIT

The maximum sequence identity for sequences in our CAMEO reconstruction dataset to our training set is 0.72, when run at 80% coverage against the consensus sequence for each MSA in the training set (with the average maximum identity across all sequences in the test set being 0.55). This is an even stricter threshold than MSAGPT (which uses 90% coverage instead). Furthermore, the MSAs we used for training come from the OpenProteinSet, which consists of sequences searched from Uniclust30 v2018-8. The cutoff for AlphaFold3 training data is September of 2021, and the cutoff for ESM2 training data is February of 2020. The CAMEO structures we used for reconstruction evaluation, however, were all deposited in May of 2025. This rigorous separation ensures the novelty of our test set. This is in line with ProfileBFN, which trains on the same corpus as ESM2, while evaluating their model on CAMEO structures deposited in May of 2024. For the zero-shot/few-shot augmentation task, we use the same test set as MSAGPT, which is also trained on the OpenProteinSet. The authors ensure minimal data leakage between the train and test set during their experiments, which implies the same for MSAFlow.

6.2 Additional Case Studies

To further validate the robustness of MSAFlow's zero-shot predictions, we provide more cases for comparison. From the table 4, we can observe that MSAFlow achieves improvement on cases with different structural patterns as well as different families.

PDB ID	Length	Description	GT	MSAGPT	MSAFlow
6NW8_A	27	Scorpion venom toxin	0.39	0.40	0.53
6WKK_X	280	Phage capsid	0.28	0.27	0.55
7EQB_B	80	Central spindle assembly	0.65	0.58	0.71
7QRR_L	153	Noumeavirus	0.31	0.61	0.83
7ZOL_A	151	Cas 7-11 regulator	0.33	0.34	0.67

Table 4: Performance comparison of MSAFlow with baseline methods on clinically relevant proteins showing TM-Score improvements across different structural patterns and protein families.

6.3 INFERENCE SPEED AND MEMORY COST

In order to demonstrate that MSAFlow exhibits notable improvements in sampling efficiency compared to other MSA-based generative models, We benchmark MSAFlow against existing tools, attempting to generate 100 sequences conditioned on an existing MSA with 6 sequences on an NVIDIA A40 GPU, and observe the following:

	Latency Per Sequence	Memory Consumption
MSAFlow	1.02s	5.8 GiB
ProfileBFN	8.49s	7.7 GiB
MSAGPT	62.46s	41.6 GiB
EvoDiff	478.24s	4.0 GiB

Table 5: Sampling efficiency comparison of MSAFlow with baseline methods showing latency per sequence and memory consumption on NVIDIA A40 GPU for generating 100 sequences conditioned on an MSA with 6 sequences.

We find that MSAFlow has better sampling efficiency, both in terms of speed and memory. We can attribute this to the fact that our model only has to deal with $L \times H$ embedding of the MSA, rather than carry the quadratic cost of representing an MSA in the ambient space. The result shows that MSAFlow has the potential to be a highly light-weight and accurate MSA designer.

Moreover, our pipeline utilizes outputs from tools like MMseqs and HMMER for Multiple Sequence Alignment (MSA) reconstruction. A key advantage of this approach is its ability to generate high-quality MSAs even when these standard homology search methods fail to find sufficient homologous

information. To provide a quantitative comparison of computational cost, we evaluated our MSAFlow model against HMMER and MMseqs2 for generating an MSA from a single query sequence (PDB 9BCZ_A from CAMEO, 644 amino acids). The empirical results are detailed below.

Method	Wall Clock Time (s)
MSAFlow (100 seqs)	153.93
HMMER	310.92
MMseqs2	497.73

Table 6: Computational cost comparison for generating MSA from query sequence alone (PDB 9BCZ_A from CAMEO, 644 AA) showing wall clock time in seconds.

These results show that MSAFlow achieves over $2 \times$ speedups compared to HMMER and MMseqs2, while still providing the ability to operate in settings where homology search fails. This confirms that MSAFlow not only addresses the coverage gap but also offers computational efficiency advantages over traditional methods.

6.4 ABLATION STUDY OF RECONSTRUCTION SEQUENCES

We address using the additional ablation study on the reconstruction task with 2, 4, 8, 16, and 32 decoded MSA sequences, as well as the comparison with natural-MSA depth on 3 samples from the CAMEO reconstruction test set.

When we keep 2-4 sequences, the MSAFlow reconstructions beat the random ground-truth subsample. As we generate more sequences, the designed MSAs generally match that of the ground-truth samples (AlphaFold3 searched MSA), indicating that MSAFlow accurately captures structure patterns of protein families.

	PDB ID	2	4	8	16	32
	9EJY	0.59	0.55	0.85	0.80	0.86
Ground Truth Random Sample	9BIX	0.19	0.32	0.35	0.32	0.49
	9CVV	0.35	0.31	0.93	0.97	0.98
	9EJY	0.61	0.61	0.84	0.83	0.84
MSAFlow Reconstruction	9BIX	0.28	0.22	0.20	0.30	0.26
	9CVV	0.43	0.62	0.87	0.87	0.97

Table 7: Ablation study comparing MSAFlow reconstruction performance against ground truth random samples across different sequence counts on CAMEO reconstruction test set. Values represent performance metrics for MSA reconstruction quality. Numbers in the first row denotes the amounts of decoding MSA sequences.

6.5 ABLATION STUDY ON SYNTHETIC AND RECONSTRUCTED MSAS

The reconstruction pathway preserves the authentic signal from a limited, shallow MSA, while the latentflow pathway generates evolutionary diversity generalized from other MSA-rich proteins. These two tracks provide complementary signals that make the few-shot augmentation stronger. To provide evidence for this, we detail the separate contributions of each track below:

As shown in the table, the reconstruction path focuses on preserving crucial motif information within the limited observed sequences, which is reflected in the lower entropy signals in the shallow MSA. In contrast, the latentflow path generates synthetic MSAs that provide evolution-consistent diversity, resulting in higher entropy.

The combination of both tracks leads to an improvement in TM score and an increase in entropy. This observation confirms that the two tracks offer complementary signals, which synergistically improve quality. Finally, by augmenting the shallow ground truth MSA with the combined generation output, we improve prediction accuracy and achieve a better TM score than the MSAGPT baseline, which is what we report in Table 1. As can be seen, MSAFlow is the only method to achieve a better TM score than the ground truth, with an entropy value closest to it.

Few-shot task	TM Score	Avg Per-position Entropy
Syn-16	0.54	2.23
Rec-16	0.52	1.33
Syn+Rec-32	0.57	2.69
Syn+Rec+GT	0.60	2.58
MSAGPT+GT	0.58	1.33
GT	0.58	2.16

Table 8: Ablation study showing the complementary contributions of synthetic and reconstructed MSA pathways in few-shot tasks, demonstrating improved TM scores and entropy characteristics. **Syn** represents Synthetic MSAs; **Rec** represents Reconstructed MSAs. The number denotes amount of MSA sequences.

6.6 ABLATION STUDY ON ESM EMBEDDINGS

To clarify the individual contributions of the ESM embeddings and our proposed Statistical Flow-matching decoding mechanism, we perform an ablation study on the zero-shot augmentation track of MSAFlow. Specifically, we compare:

- A simple feature regression task that learns MSA embeddings from ESM2 features
- Replacing ESM2 embeddings with one-hot encodings of the query sequence
- Full ESM2 embeddings with our latent statistical flow-matching decoder

Method	TM Score
MSAGPT (3B)	0.53
MSAFlow Latent w/ ESM2 regression (128M)	0.54
MSAFlow Latent w/ one-hot (130M)	0.55
MSAFlow Latent w/ ESM2 (130M)	0.62

Table 9: Ablation study comparing the contribution of ESM embeddings versus one-hot sequence encoding in MSAFlow's zero-shot MSA augmentation performance.

The results demonstrate that the efficiency of our method. Moreover, ESM2 encoding provides more useful signals to address the evolutionary information.

6.7 GENERATIVE SAMPLING PROCESS

To sample a synthetic MSA embedding, we convert the ODE flow into an SDE following Geffner et al. (2025a), and integrate the reverse-time stochastic differential equation:

$$dz_t = \left(v_t^{\theta} - \frac{1}{2}g_t^2 \cdot s_t^{\theta}(x_t)\right)dt + T \cdot g_t \cdot d\bar{W}_t, \quad t \in [0, 1]$$
(8)

where $f_t = -\frac{z_t}{1-t}$ is the drift term of the forward rectified flow, $g_t = \sqrt{\frac{2t}{1-t}}$ is the diffusion coefficient, $s_t^\theta = \nabla_{z_t} \log p_\theta(z_t|e,t)$ is the score function that can be converted from our predicted $v^\theta, T \in [0,1]$ is a temperature parameter, and $d\bar{W}_t$ is the standard Wiener process running backward in time.

We implement the sampling using the Euler-Maruyama discretization with steps of size Δt :

$$z_{t-\Delta t} = z_t - \left(v_t^{\theta} - \frac{1}{2}g_t^2 \cdot s^{\theta}(z_t, e, t)\right) \Delta t - T \cdot g_t \sqrt{\Delta t} \cdot \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I)$$
 (9)

where $v_{\theta}(z_t, e, t)$ is the time-dependent vector field predicted by the DiT. The temperature parameter T controls the stochasticity of the generation: T=1 reproduces the exact generative SDE used during training, while $T \to 0$ suppresses the noise, approaching the deterministic probability-flow ODE.

7 USAGE OF LANGUAGE MODELS

We use large language model (LLM) to aid in the preparation of this manuscript. Its use was limited to editorial tasks, including proofreading for typographical errors, correcting grammar, and improving the clarity and readability of the text.