

# GCP-VQVAE: A GEOMETRY-COMPLETE LANGUAGE FOR PROTEIN 3D STRUCTURE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Converting protein tertiary structure into discrete tokens via vector-quantized variational autoencoders (VQ-VAEs) creates a language of 3D geometry and provides a natural interface between sequence and structure models. While pose invariance is commonly enforced, retaining chirality and directional cues without sacrificing reconstruction accuracy remains challenging. In this paper, we introduce GCP-VQVAE, a geometry-complete tokenizer built around a strictly SE(3)-equivariant GCPNet encoder that preserves orientation and chirality of protein backbones. We vector-quantize pose-invariant readouts into a 4 096-token vocabulary, and a transformer decoder maps tokens back to backbone coordinates via a 6D rotation head trained with SE(3)-invariant objectives.

Building on these properties, we train GCP-VQVAE on a corpus of 24 million monomer protein backbone structures gathered from the AlphaFold Protein Structure Database. On the CAMEO2024, CASP15, and CASP16 evaluation datasets, the model achieves backbone RMSDs of 0.4377 Å, 0.5293 Å, and 0.7567 Å, respectively, and achieves 100% codebook utilization on a held-out validation set, substantially outperforming prior VQ-VAE-based tokenizers and achieving state-of-the-art performance. Beyond these benchmarks, on a zero-shot set of 2 261 completely new experimental structures, GCP-VQVAE attains a backbone RMSD of 0.8033 Å and a TM-score of 0.9747, demonstrating robust generalization to unseen proteins. Lastly, we elaborate on the various applications of this foundation-like model, such as protein structure compression and the integration of generative protein language models. We make the GCP-VQVAE source code, zero-shot dataset, and its pretrained weights fully open for the research community.

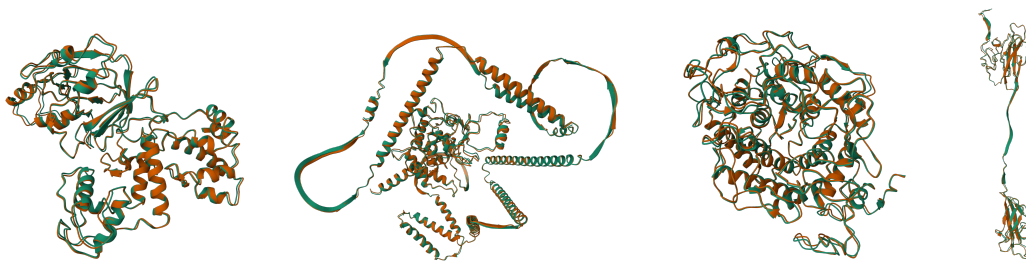


Figure 1: Superposition of GCP-VQVAE reconstructions (orange) with native structures (green) for four long, held-out backbone coordinates of proteins. Left→right: T1079 (CASP14; 482 residues; TM-score 0.9909; RMSD 0.7533 Å), T2272S8 (CASP16; 818 residues; TM-score 0.9967; RMSD 0.5639 Å), T1157s1-D1 (CASP15; 661 residues; TM-score 0.9920; RMSD 0.8172 Å), and 8BQU, chain A (CAMEO2024; 298 residues; TM-score 0.9764; RMSD 0.9981 Å).

## 1 INTRODUCTION

Proteins are the molecular machines of life, and their function is intricately tied to their three-dimensional structures (Bertoline et al., 2023; Tripathi et al., 2025). Understanding and predicting

054 these structures remains one of the central challenges in computational biology (Jänes & Beltrao,  
055 2024). Just as natural language is governed by grammatical and contextual rules, protein 3D struc-  
056 tures exhibit spatial patterns and constraints that suggest an underlying "grammar" of folds and  
057 interactions (Ruff et al., 2022; Kilgore et al., 2025; Weissenow & Rost, 2025).

058 Despite advances in protein structure prediction, effectively representing the 3D geometry of pro-  
059 teins in a form suitable for generative modeling remains an open problem (Draizen et al., 2024; Lu  
060 et al., 2025). While recent methods have begun to leverage generative AI, such as diffusion models  
061 and autoregressive frameworks, to produce full-atom structures or backbone coordinates, they still  
062 lag behind other domains such as language or vision in terms of reconstruction precision, scalability  
063 to large and diverse datasets, and openness for the broader research community, like natural lan-  
064 guage or image and video generation (Yuan et al., 2025b). These gaps continue to constrain our  
065 ability to build powerful, general-purpose protein models using modern AI techniques.

066 Beyond generative modeling, learning a discrete language for protein 3D structure opens up a wide  
067 range of downstream applications. First, compressing 3D coordinates into compact sequences of in-  
068 teger codes—while preserving accurate reconstruction—can substantially reduce storage and trans-  
069 mission costs for structural data (Kim et al., 2023). Second, discrete structural representations en-  
070 able fast, alignment-style comparison of protein shapes, analogous to multiple sequence alignment  
071 in sequence space (Van Kempen et al., 2024). Third, in learnable quantization frameworks such  
072 as vector-quantized autoencoders, these codes can be decoded into semantically rich continuous  
073 embeddings (Yuan et al., 2025b), facilitating structure-aware feature extraction for classification,  
074 clustering, structure-based comparison/search, and other predictive tasks. Finally, unifying protein  
075 sequence and structure through a shared discrete representation may pave the way for multimodal  
076 generative models that bridge amino acid sequences and 3D folds within a common language mod-  
077 eling framework (Hsieh et al., 2025).

078 However, most existing protein-structure VQVAEs are either closed-source or only partially re-  
079 leased (e.g., code without full evaluation scripts or strongest checkpoints), which impedes repro-  
080 ducible comparison (Gao et al., 2024c; Hayes et al., 2025). Furthermore, the publicly available  
081 baselines often generalize weakly to *unseen* proteins. Consequently, the field lacks a fully open-  
082 source, high-accuracy tokenizer with transparent training and evaluation that demonstrably transfers  
083 to new proteins.

084 **Contributions.** (1) We introduce GCP-VQVAE, a geometry-complete tokenizer that preserves ori-  
085 entation and chirality while producing pose-invariant codes that support missing coordinates. (2) We  
086 scale training to 24M monomers and report exhaustive evaluations (CAMEO2024, CASP14/15/16)  
087 and a zero-shot suite of newly deposited experimental structures. (3) Our model achieves state-of-  
088 the-art reconstruction on a diverse range of unseen 3D structures, and stays ahead of other open-  
089 source methods. On the zero-shot set, GCP-VQVAE attains a backbone RMSD of 0.8033 Å and a  
090 TM-score of 0.9747. (4) We release code, checkpoints, and unified evaluation scripts for ourselves  
091 and baselines, enabling reproducible comparison.

## 094 2 RELATED WORK

095  
096 An early and influential approach to casting protein 3D structure as a discrete language is FoldSeek  
097 van Kempen et al. (2022), which learns a 20-state 3Di alphabet with a VQ-VAE trained for evolu-  
098 tionary conservation and encodes structures as token sequences for ultra-fast k-mer-based local/  
099 global alignment. Building on the notion of a discrete structural language—but targeting generative  
100 reconstruction rather than search—the FoldToken series (Gao et al., 2025; 2024a;b;c) develops a  
101 VQ-VAE-style tokenizer and decoder: FoldToken introduces a SoftCVQ fold language with joint  
102 sequence–structure generation; FoldToken2 stabilizes quantization and extends to multi-chain set-  
103 tings; FoldToken3 mitigates gradient/class-space issues to reach 256-token compression with mini-  
104 mal loss; and FoldToken4 unifies cross-scale consistency and hierarchies in a single model, reducing  
105 redundant multi-scale training and code storage.

106 Yuan et al. (2025b) introduces AminoAseed codebook reparameterization plus Pareto-optimal  $K \times D$   
107 sizing, proposes structtokenbench—a fine-grained evaluation suite—and diagnoses codebook under-  
utilization in VQVAE PSTs. Concurrently, Gaujac et al. (2024) tokenizes protein backbones with

a VQ autoencoder (codebooks 4k–64k); its main open-source limitation is that it does not support sequences shorter than 50 or longer than 512 amino acids.

ESM-3 (Hayes et al., 2025) couples its multimodal transformer with a VQVAE structure tokenizer that discretizes local 3D geometry into structure tokens; structure, sequence, and function are jointly trained under a masked-token objective. The tokenizer uses an SE(3)-aware module within the encoder, and ESM-3 employs a 4 096-code structure codebook (plus special tokens) for downstream generation and masked reconstruction.

Across prior structure tokenizers, neither openness nor accuracy are yet satisfactory. The FoldToken line spans four variants, but to our knowledge, only FoldToken-4 offers a partial open release, without the strongest checkpoints, limiting reproducibility (Gao et al., 2024c). ESM-3 exposes an internal VQ-VAE tokenizer, yet public artifacts lack fully documented training/evaluation procedures and best weights, making directly comparable reconstruction benchmarking difficult (Hayes et al., 2025). The open VQ autoencoder of Gaujac et al. (2024) supports a restricted length window (e.g.,  $\sim 50$ –512 residues), precluding fair assessment on long chains where reconstruction accuracy degradation is most evident. Finally, FoldSeek’s learned 3Di alphabet targets ultra-fast structure search and does not provide a generative decoder from discrete codes back to coordinates, so reconstruction fidelity cannot be evaluated (van Kempen et al., 2022). Consequently, the community still lacks a fully open, end-to-end tokenizer with released best weights and source codes that attains high reconstruction accuracy on *unseen* proteins.

### 3 DATASET

We began with the latest release of UniRef50 (Consortium, 2019)<sup>1</sup>, which clusters protein sequences at 50% identity and thus offers a natural, non-redundant scaffold for large-scale structure modeling. For every UniRef50 entry with a corresponding model in the AlphaFold Database (AFDB; Varadi et al. (2022)), we downloaded the per-protein structure (PDB format at collection time). Limiting to UniRef50 reduces near-duplicate leakage by ensuring that homologs within clusters do not exceed 50% identity. After parsing and splitting multi-chain records into individual chains, this procedure yielded approximately 42M single-chain PDB samples.

From this  $\text{AFDB} \cap \text{UniRef50}$  pool, we drew a uniform random sample of 24M single-chain structures to form the training set. This down-sampling keeps training throughput tractable while preserving the global distribution of lengths, folds, and taxa present in the full pool.

Table 1: Dataset and benchmark statistics. Training data is deduplicated at 100% sequence identity against all validation/test splits and external benchmarks (CAMEO2024, CASP14–16). We also include a zero-shot set of 2 261 newly deposited experimental monomer chains (PDB, 2024–2025).

Split	Source	Selection criterion	# Samples
Training	AFDB $\cap$ UniRef50	Uniform random sample from $\sim 42$ M pool	24 000 000
Validation	From Tests A/B/C	2k random per test (A,B,C) merged	6 000
Test set A	PDB DB	cross-referenced AF pLDDT $\leq 70$ ; len $\geq 25$ ; $< 25$ consecutive missing residues; dedup	31 112
Test set B	AFDB	97 species; $\sim 2.5$ k/species; len $\geq 25$ ; dedup	17 353
Test set C	PDB DB	under-represented taxonomies in SwissProt; pLDDT $\geq 95$ ; $< 25$ consecutive missing residues; dedup	31 867
Zero-shot (experimental)	PDB DB	recent experimental monomers (2024–2025); len 25–2048; dedup; $< 50\%$ seq id vs train	2 261
CAMEO2024	CAMEO	-	574
CASP14	CASP	-	33
CASP15	CASP	-	45
CASP16	AF3 predictions	-	104

We evaluate on three held-out test suites capturing complementary shifts (see Table 1): *A* low-confidence AF2 SwissProt models, *B* species-diverse (97-species) taxonomic shift, and *C* high-confidence structures from under-represented taxa. All suites are per-chain, length-filtered, deduplicated, and NaN-aware; a balanced validation set of 6,000 chains (2k per suite) is sampled disjointly. Full curation criteria, thresholds, and counts are provided in Appendix A.2 .

To probe generalization to genuinely new structures, we assembled a zero-shot suite from *recently deposited experimental* monomeric chains in the PDB (calendar years 2024–2025). Multi-chain entries were split into per-chain backbones, chains shorter than 25 residues or longer than 2 048 residues were discarded, and we retained entries even when they contained missing coordinates

<sup>1</sup>March 2024 release.

(NaN) to reflect typical experimental gaps, yielding 76 503 pdb's in total. We deduplicated using structural clustering with FoldSeek easy-cluster (alignment type 2, tmscore threshold 0.95), retaining one representative per cluster. This resulted in 19 080 clusters (van Kempen et al., 2022). For leakage control, each representative was compared with foldseek easy-search against (i) the AFDB  $\cap$  UniRef50 set used for our GCP-VQVAE training and (ii) the afdB\_rep\_v4 (Barrio-Hernandez et al., 2023) representative set used to pre-train the GCPNet encoder in ProteinWorkshop (Jamasp et al., 2024). We excluded any candidate with  $\geq 50\%$  structure similarity to any chain in those sources; the resulting zero-shot set comprises 2 261 monomer chains.

**Independent benchmarks.** In addition to the internal splits, we evaluate on community benchmarks to facilitate comparison with prior work: CAMEO2024 (Robin et al., 2021; Leemann et al., 2023), CASP14 (Kryshtafovych et al., 2021), CASP15 (Kryshtafovych et al., 2023), and AlphaFold3 (Abramson et al., 2024) predictions for CASP16 (Yuan et al., 2025a) targets. Each benchmark is processed with the same per-chain extraction and deduplication pipeline and is used strictly for out-of-distribution evaluation.

Table 1 summarizes the composition and selection criteria of all splits. Counts after converting multi-chain inputs into per-chain samples. Also, all samples are truncated to a maximum of 2048 amino acids in length.

## 4 GCP-VQVAE ARCHITECTURE

The proposed architecture leverages two main parts: (1) a GCPNet encoder to encode backbone coordinates into embeddings, and (2) a transformer-based VQVAE, which discretizes backbone embeddings and then converts them back into 3D coordinates.

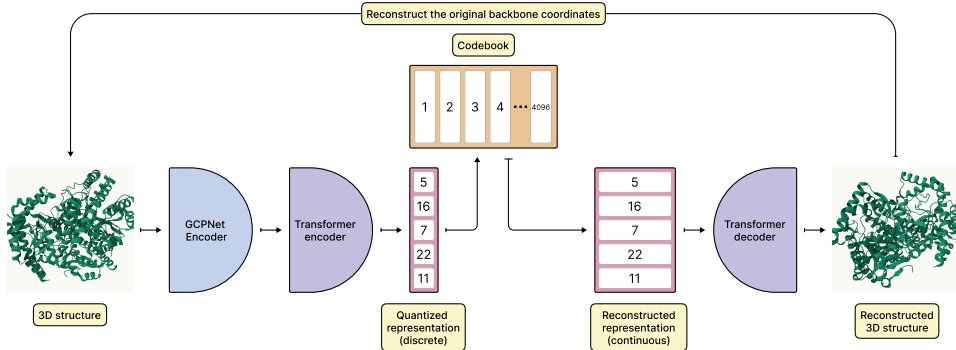


Figure 2: GCP-VQVAE overview. A protein backbone (N-C $\alpha$ -C) is first encoded by a SE(3)-equivariant GCPNet that preserves orientation and chirality. A Transformer encoder produces latents that are vector-quantized into a sequence of code indices from a 4 096-entry codebook (e.g., 5, 16, 7, 22, 11), yielding pose-invariant discrete tokens. For reconstruction, the indices are de-quantized to continuous embeddings and passed to a Transformer decoder equipped with a 6-D rotation head, which predicts rigid updates to recover the original backbone coordinates.

### 4.1 GCPNET ENCODER

GCPNet (Morehead & Cheng, 2024b) extends the scalar-vector message-passing philosophy of GVP-GNN to a geometry-complete, SE(3)-equivariant encoder. Every atom  $i$  in a molecular graph  $G = (V, E)$  carries scalars  $s_i \in \mathbb{R}^{d_s}$  and row-wise vectors  $\mathbf{v}_i \in \mathbb{R}^{d_v \times 3}$  that rotate as  $\mathbf{v}_i \mapsto R\mathbf{v}_i$  under  $g = (R, \mathbf{t}) \in \text{SE}(3)$ , while each edge  $(i, j)$  stores analogous features  $(s_{ij}, \mathbf{v}_{ij})$  and the relative displacement  $\mathbf{r}_{ij} = \mathbf{x}_j - \mathbf{x}_i$ . Before any update, the model attaches to every edge a right-handed orthonormal frame  $F_{ij} = [\mathbf{a}_{ij}, \mathbf{b}_{ij}, \mathbf{c}_{ij}] \in \text{SO}(3)$  with  $\mathbf{a}_{ij} = \mathbf{r}_{ij}/\|\mathbf{r}_{ij}\|$ ; this frame supplies a reference for chirality and orientation that GVP-GNN lacks.

The core computation is a Geometry-Complete Perceptron (GCP) micro-step that first down-scales the vectors and then projects them into the local frame to extract nine orientation-aware features.

Denoting  $\mathbf{z}_{ij} = W_d \mathbf{v}_{ij}$  and  $\text{vec}(\cdot)$  the row-wise vectorization, the joint update mixes the old scalars with their orientation signatures (the frame-projected vectors and their norms) and gates the vectors through a learnable row-wise gate to preserve equivariance;  $\phi_s$  is an MLP and  $\sigma$  denotes a learnable gating function that need not be a fixed sigmoid (Equation 1).

$$(s_{ij}, \mathbf{v}_{ij}) \mapsto \left( \phi_s[s_{ij}, \|\mathbf{z}_{ij}\|_2, \text{vec}(\mathbf{z}_{ij} F_{ij}^\top)], \sigma(W_g s_{ij}) \odot W_v \mathbf{z}_{ij} \right) \quad (1)$$

In the node update, the orientation features are averaged over neighbors. A sequence of such micro-steps, wrapped by a residual shortcut, forms a GCPCConv edge block. After aggregating messages  $m_i = \sum_{j \in \mathcal{N}(i)} \text{GCPCConv}(i, j)$ , a gated scalar–vector MLP updates the node features, and stacking  $L$  layers yields an invariant backbone. When tasks require coordinates, each layer appends an equivariant displacement head whose output is a learned 3-dimensional vector; this vector is added residually to  $\mathbf{x}_i$  and re-centered to ensure translation invariance, enabling force or trajectory prediction without breaking equivariance.

$$\mathbf{x}_i \leftarrow \mathbf{x}_i + \mathbf{f}_i, \quad \mathbf{f}_i = \text{MLP}_{\text{disp}}(m_i) \quad (2)$$

Equations 1 and 2 commute with every rigid motion, making the encoder strictly SE(3)-equivariant, while projection through  $F_{ij}$  preserves the complete set of edge orientations so that the latent remains *geometry-complete* at any depth. Eliminating the frame ( $F_{ij} = I$ ), replacing the orientation features by  $\|\mathbf{v}_{ij}\|_2$ , and dropping the coordinate head restores a frame-free, E(3)-equivariant network akin to GVP-GNN—thereby isolating the contributions responsible for our empirical gains. Because GCPNet keeps directional and chiral cues that GVP-GNN discards, it supports optional equivariant coordinate updates for force or dynamics prediction and attains improved performance across invariant (e.g., binding affinity; Morehead & Cheng (2024b)), equivariant (e.g., force regression; Morehead & Cheng (2024b)), and coordinate-generative (e.g., molecular diffusion; Morehead & Cheng (2024a)) tasks with only modest extra computation.

## 4.2 VQVAE

The transformer-based VQVAE employed in this work is organized into the classical three-stage pipeline of *encoder*, *vector-quantization*, and *decoder*. The encoder processes the given embeddings into a sequence of latent vectors, the quantizer discretizes these latents, and the decoder reconstructs the original signal from the resulting code indices.

Both stacks adopt a lightweight pre-layer-normalized Transformer that integrates several recent efficiency upgrades: (i) Pre-LayerNorm places the LayerNorm before each sub-block (Xiong et al., 2020), which keeps activations in a well-behaved range throughout the network, reduces gradient-scale drift, and therefore allows training with larger learning rates and much milder warm-up schedules; (ii) separately normalizes query and key vectors before computing attention logits (Henry et al., 2020). This prevents overly large dot-products, stabilizes attention distributions, and mitigates softmax saturation, especially in low-resource or small-batch training scenarios; (iii) Grouped-Query Attention shares key–value projections across groups of query heads (Ainslie et al., 2023), reducing both memory and compute without harming quality; (iv) Rotary Positional Embeddings (RoPE) inject relative-position information by applying position-dependent planar rotations to each query–key pair (Su et al., 2024), letting the model generalize to much longer sequences with virtually no extra computational cost; (v) We remove bias terms from projection and feed-forward layers, an established simplification that has negligible effect on accuracy while trimming a fraction of parameters and FLOPs.

In the architecture, the vector quantization layer provides the discrete bottleneck between the aforementioned encoder and decoder and follows the learnable formulation of VQVAE, augmented with several improvements described in the following to accelerate convergence, increase codebook usage, and stabilize large codebooks training.

Before training starts, we run  $k$ -means on the encoder outputs of the first mini-batch to seed the codebook as displayed in Equation 3, where  $Z^{(0)} \subset \mathbb{R}^d$  are the features,  $K$  the codebook size, and  $T$  the number of Lloyd iterations. Empirically, this step mitigates early code collapse and improves utilization when  $K$  is large.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

$$E^{(0)} = \text{KMEANS}(Z^{(0)}, K, T), \quad (3)$$

Given an encoder vector  $\mathbf{z} \in \mathbb{R}^d$ , quantization proceeds by nearest-neighbor lookup, Equation 4, which is identical to the vanilla VQ rule.

$$k = \arg \min_j \|\mathbf{z} - \mathbf{e}_j\|_2^2, \quad \mathbf{z}_q = \mathbf{e}_k, \quad (4)$$

To transmit gradients through this non-differentiable operation, instead of using the straight-through estimate (STE) introduced in Van Den Oord et al. (2017), we adopt the rotation trick (Fifty et al., 2024). During back-propagation, the Jacobian is replaced by Equation 5, where  $R$  is the shortest-arc rotation aligning the unit vectors  $\hat{\mathbf{z}}$  and  $\hat{\mathbf{e}}_k$ . This modification embeds both angular and magnitude mismatch into the back-propagated signal, yielding faster convergence and richer code usage in practice.

$$\mathbf{J}_k = \frac{\|\mathbf{e}_k\|}{\|\mathbf{z}\|} R(\hat{\mathbf{z}} \rightarrow \hat{\mathbf{e}}_k), \quad (5)$$

The codebook updates and the encoder commitment follow the standard VQ losses (Equation 6), with  $\beta$  balancing the commitment pressure.

$$\mathcal{L}_{\text{code}} = \|\text{sg}[\mathbf{z}] - \mathbf{e}_k\|_2^2, \quad \mathcal{L}_{\text{commit}} = \beta \|\mathbf{z} - \text{sg}[\mathbf{e}_k]\|_2^2, \quad (6)$$

To encourage diverse and well-separated embeddings, we regularize the codebook with the Frobenius norm (Equation 7) weighted by  $\lambda_{\text{orth}}$ . This orthogonality constraint spreads codes over the hypersphere and curbs under-utilization, which is a common failure mode in large books.

$$\mathcal{L}_{\text{orth}} = \|E^\top E - I_K\|_F^2, \quad (7)$$

To translate the decoder’s abstract embeddings into a physically meaningful 3D structure, we utilize a 6D rotation head on top of the decoder. This module’s primary purpose is to provide a stable and continuous parameterization of 3D rotations and translations, which is crucial for effective training of deep neural networks. This approach, notably used in AlphaFold2 and supported by the findings of Zhou et al. (2019) on rotation representations, avoids the well-known issues of Gimbal lock in Euler angles and the double-cover ambiguity of quaternions.

Intuitively, the head operates by predicting an intermediate representation for each residue, comprising two 3D direction vectors and a translation. The direction vectors are deterministically converted into a stable rotation matrix via the Gram-Schmidt process, while the translation is scaled by a hyper-parameter,  $\alpha$ , to an arbitrary range (e.g., Å). This resulting rigid transformation acts as an update, which is composed with the residue’s running pose. The final backbone coordinates are then generated by applying this new, refined pose to a fixed local atomic template ( $\mathbf{X}^{\text{local}}$ ). This iterative process ensures the structure is built in a geometrically consistent and equivariant manner, with the full operational details provided in Algorithm 1.

The decoder is supervised with a weighted sum of three geometric objectives. Building on the loss terms defined in Algorithm 2 (distance, direction and aligned MSE) and the Kabsch alignment returned by Algorithm 3, we supervise the decoder with the weighted sum of Equation 8 as the reconstruction part ( $\mathcal{L}_{\text{rec}}$ ) of the final loss.

$$\mathcal{L}_{\text{rec}} = \lambda_{\text{mse}} \mathcal{L}_{\text{MSE}} + \lambda_{\text{dist}} \mathcal{L}_{\text{dist}} + \lambda_{\text{dir}} \mathcal{L}_{\text{dir}}, \quad (8)$$

Here,  $\mathcal{L}_{\text{MSE}}$  is the mean-squared error between predicted and native backbones after Kabsch alignment,  $\mathcal{L}_{\text{dist}}$  penalises deviations in the  $3L \times 3L$  backbone distance matrix (clamped at  $5 \text{ \AA}^2$ ), and  $\mathcal{L}_{\text{dir}}$  measures squared differences of pairwise dot-product tensors over the six orientation vectors per residue (clamped at 20). Finally, the overall objective optimized for each iteration is demonstrated in Equation 9 where  $\mathcal{L}_{\text{rec}}$  is the decoder reconstruction.

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{code}} + \lambda_{\text{commit}} \mathcal{L}_{\text{commit}} + \lambda_{\text{orth}} \mathcal{L}_{\text{orth}}, \quad (9)$$

## 5 EXPERIMENTS

In this section, optimization is used with AdamW (Loshchilov, 2017) and a cosine-annealed learning-rate schedule with warm-up steps (Loshchilov & Hutter, 2016). Experiments were implemented in PyTorch 2.7 (Ansel et al., 2024)<sup>2</sup> with mixed-precision (BF16) training (Kalamkar et al., 2019; Micikevicius et al., 2017) on four nodes of NVIDIA 8 × A100 GPUs. For the GCP-Net encoder, we initialized from the ProteinWorkshop checkpoint (Jamasb et al., 2024). The full GCP-VQVAE configuration and exact hyperparameter values are listed in Appendix A.3; Table 5 and 6.

Table 2: Comparison with the available open-source structure tokenizer methods. Except for FoldToken-4, which uses 256 vocab size, all methods use 4096 vocab size. The Structure Tokenizer of Gaujac et al. (2024) only supports sequences of length 50–512, metrics for that baseline exclude out-of-range samples (other methods are evaluated on the full sets). Also, zero-shot results for it are omitted due to limited coverage.

Dataset	Metric	GCP-VQVAE (Ours)	FoldToken 4 (Gao et al., 2024c)	ESM-3 VQVAE (Hayes et al., 2025)	(Gaujac et al., 2024)
CASP14	TM-score ↑	<b>0.9890</b>	0.5410	0.5042	0.3624
	RMSD ↓	<b>0.5431</b>	8.9838	10.4611	10.5344
CASP15	TM-score ↑	<b>0.9884</b>	0.3289	0.3206	0.2329
	RMSD ↓	<b>0.5293</b>	14.6702	13.1877	14.8956
CASP16	TM-score ↑	<b>0.9857</b>	0.8055	0.7685	0.6058
	RMSD ↓	<b>0.7567</b>	5.5094	8.2640	8.7106
CAMEO2024	TM-score ↑	<b>0.9918</b>	0.4784	0.4633	0.3575
	RMSD ↓	<b>0.4377</b>	12.1089	12.1138	13.5360
Zero-Shot	TM-score ↑	<b>0.9747</b>	0.3324	0.3131	-
	RMSD ↓	<b>0.8033</b>	17.4449	18.9335	-

Training proceeds in two stages on the same training split. In Stage 1, sequences are truncated to 512 residues. In Stage 2, the maximum length is increased to 1280 residues. In Stage 2, we also introduce a simple NaN-masking augmentation to handle missing coordinates; implementation details are provided in Appendix A.3 (Figure 5). At the end of Stage 2 training, on the validation set, the model attains MAE 0.2239, RMSD 0.4281, GDT-TS 0.9856, and TM-score 0.9889 with 100% codebook utilization; per-test-set results appear in Table 3.

Table 3: Evaluate our GCP-VQVAE method on the predefined test sets.

Method	Test set A	Test set B	Test set C
RMSD ↓	0.5124	0.4027	0.4787
TM-score ↑	0.9823	0.9951	0.9885

We compare against open-source structure tokenizers: FoldToken-4, ESM-3 VQ-VAE, and the Structure Tokenizer of Gaujac et al. (2024). In practice, ESM-3 does not provide documented usage or official evaluation scripts for its VQ-VAE; we therefore reproduced its evaluation using the publicly released weights and their GitHub codebase. The original FoldToken-4 evaluation repository was prohibitively slow, so we re-implemented it, yielding a  $\sim 20\times$  speed-up with negligible loss in reconstruction rate. For Gaujac et al. (2024) we used the authors’ released scripts with the 4096-entry codebook checkpoint. See Table 2 for aggregate metrics; Figure 4 visualizes the full RMSD error across test sets.

We examined how reconstruction error varies with sequence length on the validation set (Appendix A.3; Figure 6). Errors increase only mildly with length, indicating stable accuracy up to 1280 residues with only a slight variance rise for very long chains. On the zero-shot set we obtain mean RMSD 0.8033 Å and TM-score 0.9747 (Table 2), with group-wise trends summarized in Table 4; extended zero-shot analyses appear in Appendix A.4.

<sup>2</sup>We built the VQVAE components extensively by using the x-transformers and vector-quantize-pytorch libraries.

Table 4: Metrics by structure group on the zero-shot dataset. Groups are defined by secondary-structure composition thresholds and are not mutually exclusive. Thresholds of each groups are described in Appendix A.2

Group	Count	RMSD ↓	TM-score ↑
All rows	2261	0.8033	0.9747
A (mainly $\alpha$ )	615	0.7671	0.9751
B (mainly $\beta$ )	582	0.6713	0.9737
AB (mixed)	97	0.6283	0.9883
D (coil-rich/disordered)	691	0.9004	0.9681
No-NaN residues	1314	0.5978	0.9796
Has-NaN residues ( $\geq 1$ )	947	1.0885	0.9680

## 6 DISCUSSION AND FUTURE WORK

GCP-VQVAE delivers high-fidelity reconstruction across diverse suites of benchmarks, with TM-scores typically  $\geq 0.98$  and mean RMSDs in the 0.40–0.80 Å range on CAMEO2024, CASP14/15/16, and comparable performance on our internal Test A/B/C splits (e.g., 0.40–0.51 Å RMSD).

On our diverse validation set, although overall reconstruction is strong, we observe a mild RMSD drift with sequence length (Figure 6). We attribute this to length imbalance in training; short chains dominate, leaving the model relatively underexposed to long-range constraints and rare structural motifs. Mitigations include targeted fine-tuning on a length-balanced subset, reweighting/upsampling long sequences, and length-aware objectives, which should reduce the slope without architectural changes.

On the zero-shot set of 2 261 newly deposited experimental monomers, the model maintains strong generalization with mean RMSD 0.8033 Å and TM-score 0.9747 (Table 4). Our diagnostics indicate that degradation is driven primarily by the ratio and count of missing coordinates (Appendix Figure 10, 9), and secondarily by sequence length (Appendix Figure 8), Similar to our interpretation of the zero-shot results. Consistent with this, sequences without missing residues exhibit substantially lower errors than those with gaps (Table 4). These trends suggest further gains from strengthening the NaN-masking augmentation (e.g., longer and multi-gap patterns) and increasing exposure to longer chains.

Relative to available open-source tokenizers, our method shows a large margin (often an order of magnitude lower RMSD) while maintaining near-ceiling TM-scores. We believe part of this disparity stems from incomplete public releases: several baselines either restrict sequence lengths or withhold their strongest checkpoints (as is evident for FoldToken-4 repository) and end-to-end evaluation scripts. To enable rigorous verification, we release our checkpoints and evaluation pipelines for both GCP-VQVAE and reproduced baselines, and we invite the community to run, audit, and extend these comparisons.

For the proposed GCP-VQVAE architecture, we see the following applications:

**Structure compression.** Using a 4 096-entry codebook yields about  $\sim 24\times$  compression of backbone coordinates for a 512-residue monomer (Appendix A.5). Given our low RMSDs of  $\sim 0.4\text{--}0.8$  Å (Table 2), this is a practical lossy backbone codec. Further gains are possible by (i) allowing a single code to represent short residue spans (e.g., 2–4 amino acids), reducing the token rate while the decoder expands span-tokens to per-residue poses, and/or (ii) increasing the codebook size (e.g., 8k–64k) to lower per-token quantization error so that fewer tokens are needed at a target distortion. These options trade bitrate against utilization and latency, and remain compatible with our VQ training approach.

**Downstream 3D learning representation.** De-quantized code embeddings yield continuous, structure-aware features that potentially are more separable and clusterable than base GCPNet outputs, providing a strong encoder for downstream tasks, as observed in Yuan et al. (2025b).

**Generative modeling.** Because our structure representation is discrete and pose-invariant, it slots directly into autoregressive protein language model (PLM) pre-training: tokens can be interleaved with amino-acid symbols, letting us exploit next-token scaling laws observed in autoregressive Large

432 Language Models (LLM) (Kaplan et al., 2020). This unifies sequence and structure in a single  
433 model, enabling controllable generation via simple conditioning; e.g., sequence-given-structure and  
434 structure-given-sequence modes. Structure tokens serve as an explicit geometric prior for PLMs  
435 (Hayes et al., 2025; Su et al., 2023), potentially enabling more control over protein design. More-  
436 over, the same machinery supports higher-throughput structure prediction from sequence by first  
437 predicting structural tokens and then mapping them to backbones with GCP-VQVAE’s decoder  
438 (Pourmirzaei et al., 2025; Lu et al., 2024; Chen et al., 2024).

439 Several avenues look promising to explore: scaling GCP-VQVAE tokenization to multi-chain com-  
440 plexes, providing a lightweight variant for low-latency use, enriching tokens with side-chain infor-  
441 mation, and evaluating GenAI-friendliness of this new language to unify sequence and structure  
442 generation via autoregressive PLMs.

## 443 7 CONCLUSION

444 In this work, we have introduced GCP-VQVAE, a state-of-the-art open-source protein structure tok-  
445 enizer. GCP-VQVAE offers improved structure reconstruction quality and generalization compared  
446 to prior methods, which could enable new (sequence and structure-based) predictive and generative  
447 modeling applications in future work.

## 448 REFERENCES

- 449 Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf  
450 Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure  
451 prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, 2024.
- 452 Joshua Ainslie, James Lee-Thorp, Michiel De Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit  
453 Sanghai. Gqa: Training generalized multi-query transformer models from multi-head check-  
454 points. *arXiv preprint arXiv:2305.13245*, 2023.
- 455 Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky,  
456 Bin Bao, Peter Bell, David Berard, Evgeni Burovski, et al. Pytorch 2: Faster machine learning  
457 through dynamic python bytecode transformation and graph compilation. In *Proceedings of the*  
458 *29th ACM International Conference on Architectural Support for Programming Languages and*  
459 *Operating Systems, Volume 2*, pp. 929–947, 2024.
- 460 Inigo Barrio-Hernandez, Jinger Yeo, Jürgen Jänes, Milot Mirdita, Cameron LM Gilchrist, Tanita  
461 Wein, Mihaly Varadi, Sameer Velankar, Pedro Beltrao, and Martin Steinegger. Clustering pre-  
462 dicted structures at the scale of the known protein universe. *Nature*, 622(7983):637–645, 2023.
- 463 Letícia MF Bertoline, Angélica N Lima, Jose E Krieger, and Samantha K Teixeira. Before and after  
464 alphafold2: An overview of protein structure prediction. *Frontiers in bioinformatics*, 3:1120370,  
465 2023.
- 466 Zhiyuan Chen, Tianhao Chen, Chenggang Xie, Yang Xue, Xiaonan Zhang, Jingbo Zhou, and Xi-  
467 aomin Fang. Unifying sequences, structures, and descriptions for any-to-any protein generation  
468 with the large multimodal model helixprotx. *arXiv preprint arXiv:2407.09274*, 2024.
- 469 UniProt Consortium. Uniprot: a worldwide hub of protein knowledge. *Nucleic acids research*, 47  
470 (D1):D506–D515, 2019.
- 471 Eli J Draizen, Stella Veretnik, Cameron Mura, and Philip E Bourne. Deep generative models of  
472 protein structure uncover distant relationships across a continuous fold space. *Nature Communi-*  
473 *cations*, 15(1):8094, 2024.
- 474 Christopher Fifty, Ronald G Junkins, Dennis Duan, Aniketh Iyengar, Jerry W Liu, Ehsan Amid,  
475 Sebastian Thrun, and Christopher Ré. Restructuring vector quantization with the rotation trick.  
476 *arXiv preprint arXiv:2410.06424*, 2024.
- 477 Zhiyuan Gao, Cheng Tan, and Stan Z Li. Foldtoken2: Learning compact, invariant and generative  
478 protein structure language. *arXiv preprint arXiv:2407.00050*, 2024a.

- 486 Zhangyang Gao, Cheng Tan, and Stan Z Li. Foldtoken3: Fold structures worth 256 words or less.  
487 *bioRxiv*, pp. 2024–07, 2024b.  
488
- 489 Zhangyang Gao, Cheng Tan, and Stan Z Li. Foldtoken4: Consistent & hierarchical fold language.  
490 *bioRxiv*, pp. 2024–08, 2024c.  
491
- 492 Zhangyang Gao, Cheng Tan, Jue Wang, Yufei Huang, Lirong Wu, and Stan Z Li. Foldtoken: Learn-  
493 ing protein language via vector quantization and beyond. In *Proceedings of the AAAI Conference*  
494 *on Artificial Intelligence*, volume 39, pp. 219–227, 2025.
- 495 Benoit Gaujac, Jérémie Donà, Liviu Copoiu, Timothy Atkinson, Thomas Pierrot, and Thomas D  
496 Barrett. Learning the language of protein structure. *arXiv preprint arXiv:2405.15840*, 2024.  
497
- 498 Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert  
499 Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years  
500 of evolution with a language model. *Science*, 387(6736):850–858, 2025.
- 501 Alex Henry, Prudhvi Raj Dachapally, Shubham Pawar, and Yuxuan Chen. Query-key normalization  
502 for transformers. *arXiv preprint arXiv:2010.04245*, 2020.  
503
- 504 Cheng-Yen Hsieh, Xinyou Wang, Daiheng Zhang, Dongyu Xue, Fei Ye, Shujian Huang, Zaixiang  
505 Zheng, and Quanquan Gu. Elucidating the design space of multimodal protein language models.  
506 *arXiv preprint arXiv:2504.11454*, 2025.
- 507 Arian Rokkum Jamasb, Alex Morehead, Chaitanya K Joshi, Zuobai Zhang, Kieran Didi, Simon V  
508 Mathis, Charles Harris, Jian Tang, Jianlin Cheng, Pietro Lio, et al. Evaluating representation  
509 learning on the protein structure universe. In *The Twelfth International Conference on Learning*  
510 *Representations*, 2024.
- 511 Jürgen Jänes and Pedro Beltrao. Deep learning for protein structure prediction and design—progress  
512 and applications. *Molecular Systems Biology*, 20(3):162–169, 2024.  
513
- 514 Dhiraj Kalamkar, Dheevatsa Mudigere, Naveen Mellempudi, Dipankar Das, Kunal Banerjee,  
515 Sasikanth Avancha, Dharm Teja Vooturi, Nataraj Jammalamadaka, Jianyu Huang, Hector Yuen,  
516 et al. A study of bfloat16 for deep learning training. *arXiv preprint arXiv:1905.12322*, 2019.
- 517 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
518 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
519 models. *arXiv preprint arXiv:2001.08361*, 2020.  
520
- 521 Henry R Kilgore, Itamar Chinn, Peter G Mikhael, Ilan Mitnikov, Catherine Van Dongen, Guy Zyl-  
522 berberg, Lena Afeyan, Salman F Banani, Susana Wilson-Hawken, Tong Ihn Lee, et al. Protein  
523 codes promote selective subcellular compartmentalization. *Science*, 387(6738):1095–1101, 2025.
- 524 Hyunbin Kim, Milot Mirdita, and Martin Steinegger. Foldcomp: a library and format for compress-  
525 ing and indexing large protein structure sets. *Bioinformatics*, 39(4):btad153, 2023.  
526
- 527 Andriy Kryshchak, Torsten Schwede, Maya Topf, Krzysztof Fidelis, and John Moult. Critical  
528 assessment of methods of protein structure prediction (casp)—round xiv. *Proteins: Structure,*  
529 *Function, and Bioinformatics*, 89(12):1607–1617, 2021.
- 530 Andriy Kryshchak, Torsten Schwede, Maya Topf, Krzysztof Fidelis, and John Moult. Critical  
531 assessment of methods of protein structure prediction (casp)—round xv. *Proteins: Structure,*  
532 *Function, and Bioinformatics*, 91(12):1539–1549, 2023.  
533
- 534 Michèle Leemann, Ander Sagasta, Jerome Eberhardt, Torsten Schwede, Xavier Robin, and Janani  
535 Durairaj. Automated benchmarking of combined protein structure and ligand conformation pre-  
536 diction. *Proteins: Structure, Function, and Bioinformatics*, 91(12):1912–1924, 2023.
- 537 I Loshchilov. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.  
538
- 539 Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv*  
*preprint arXiv:1608.03983*, 2016.

- 540 Jiarui Lu, Xiaoyin Chen, Stephen Zhewen Lu, Chence Shi, Hongyu Guo, Yoshua Bengio, and  
541 Jian Tang. Structure language models for protein conformation generation. *arXiv preprint*  
542 *arXiv:2410.18403*, 2024.
- 543 Tianyu Lu, Melissa Liu, Yilin Chen, Jinho Kim, and Po-Ssu Huang. Assessing generative model  
544 coverage of protein structures with shapes. *bioRxiv*, 2025.
- 545 Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia,  
546 Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision  
547 training. *arXiv preprint arXiv:1710.03740*, 2017.
- 548 Alex Morehead and Jianlin Cheng. Geometry-complete diffusion for 3d molecule generation and  
549 optimization. *Communications Chemistry*, 7(1):150, 2024a.
- 550 Alex Morehead and Jianlin Cheng. Geometry-complete perceptron networks for 3d molecular  
551 graphs. *Bioinformatics*, 40(2):btac087, 2024b.
- 552 Mahdi Pourmirzaei, Farzaneh Esmaili, Salhuldin Alqarghuli, Mohammadreza Pourmirzaei, Ye Han,  
553 Kai Chen, Mohsen Rezaei, Duolin Wang, and Dong Xu. Prot2token: A unified framework for  
554 protein modeling via next-token prediction. *arXiv preprint arXiv:2505.20589*, 2025.
- 555 Xavier Robin, Juergen Haas, Rafal Gumienny, Anna Smolinski, Gerardo Tauriello, and Torsten  
556 Schwede. Continuous automated model evaluation (cameo)—perspectives on the future of fully  
557 automated evaluation of structure prediction methods. *Proteins: Structure, Function, and Bioin-*  
558 *formatics*, 89(12):1977–1986, 2021.
- 559 Kiersten M Ruff, Yoon Hee Choi, Dezerae Cox, Angelique R Ormsby, Yoochan Myung, David B  
560 Ascher, Sheena E Radford, Rohit V Pappu, and Danny M Hatters. Sequence grammar underlying  
561 the unfolding and phase separation of globular proteins. *Molecular cell*, 82(17):3193–3208, 2022.
- 562 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-  
563 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 564 Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein  
565 language modeling with structure-aware vocabulary. *BioRxiv*, pp. 2023–10, 2023.
- 566 Timir Tripathi, Vladimir N Uversky, and Alessandro Giuliani. ‘intelligent’ proteins. *Cellular and*  
567 *Molecular Life Sciences*, 82(1):239, 2025.
- 568 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in*  
569 *neural information processing systems*, 30, 2017.
- 570 Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Cameron LM Gilchrist,  
571 Johannes Söding, and Martin Steinegger. Foldseek: fast and accurate protein structure search.  
572 *Biorxiv*, pp. 2022–02, 2022.
- 573 Michel Van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee,  
574 Cameron LM Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein struc-  
575 ture search with foldseek. *Nature biotechnology*, 42(2):243–246, 2024.
- 576 Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina  
577 Yordanova, David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. Alphafold protein  
578 structure database: massively expanding the structural coverage of protein-sequence space with  
579 high-accuracy models. *Nucleic acids research*, 50(D1):D439–D444, 2022.
- 580 Konstantin Weissenow and Burkhard Rost. Are protein language models the new universal key?  
581 *Current Opinion in Structural Biology*, 91:102997, 2025.
- 582 Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang,  
583 Yanyan Lan, Liwei Wang, and Tiejian Liu. On layer normalization in the transformer architecture.  
584 In *International conference on machine learning*, pp. 10524–10533. PMLR, 2020.
- 585 Rongqing Yuan, Jing Zhang, Andriy Kryshchak, R Dustin Schaeffer, Jian Zhou, Qian Cong, and  
586 Nick V Grishin. Casp16 protein monomer structure prediction assessment. *Proteins: Structure,*  
587 *Function, and Bioinformatics*, 2025a.

594 Xinyu Yuan, Zichen Wang, Marcus Collins, and Huzefa Rangwala. Protein structure tokenization:  
595 Benchmarking and new recipe. *arXiv preprint arXiv:2503.00089*, 2025b.

597  
598 Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation  
599 representations in neural networks. In *Proceedings of the IEEE/CVF conference on computer  
600 vision and pattern recognition*, pp. 5745–5753, 2019.

## 603 A APPENDIX

### 606 A.1 METHODS

608 Algorithms 1–3 define the decoder head, geometric losses, and alignment used in training. Given  
609 per-residue embeddings, the 6D head projects to two direction vectors and a translation, constructs  
610 a proper rotation via Gram–Schmidt, composes this rigid update with the running pose, and applies  
611 it to a fixed (N–C $_{\alpha}$ –C) template to produce backbone coordinates. Reconstruction is supervised by  
612 the weighted sum in Equation 8: (i) aligned MSE after optimal Kabsch alignment, (ii) a clamped  
613 backbone distance-matrix loss, and (iii) a clamped pairwise orientation (direction) loss built from  
614 six backbone vectors per residue.

---

#### 616 **Algorithm 1** Pseudocode for 6D rotation-based structure prediction

---

```
618 def dim6rot_structure_head(h: Tensor[N, d],
619                          g_initial: Tuple[Tensor[N,3,3], Tensor[N,3,1]],
620                          alpha: float = 1.0):
621     """
622     Predicts backbone coordinates from embeddings using a 6D rotation representation.
623     This process is vectorized over N residues.
624
625     h: [N, d] float32          -- input embeddings for N residues
626     g_initial: (R_0, t_0)      -- initial SE(3) pose (running transformation)
627     alpha: float              -- translation scaling hyper-parameter
628     """
629     # Define a fixed local reference frame for backbone atoms (N, C_alpha, C)
630     X_local = torch.tensor(...) # Shape: [3, 3]
631
632     # 1) Project embeddings to unconstrained translations and direction vectors
633     # This represents the internal FFN and projection layers.
634     h_ffn = LayerNorm(GELU(Linear(h)))
635     # W_rigid projects h_ffn to 9 dimensions for t_tilde, a, and b
636     rigid_proj = Linear(h_ffn) # Output shape: [N, 9]
637     t_tilde, a, b = torch.split(rigid_proj, 3, dim=-1)
638
639     # 2) Create rotation R_i via Gram-Schmidt orthonormalization
640     a_hat = a / torch.norm(a, dim=-1, keepdim=True)
641     c = torch.cross(a_hat, b, dim=-1)
642     c_hat = c / torch.norm(c, dim=-1, keepdim=True)
643     b_hat = torch.cross(c_hat, a_hat, dim=-1)
644
645     # Stack column vectors to form the rotation matrix R
646     R = torch.stack([a_hat, b_hat, c_hat], dim=-1) # Shape: [N, 3, 3]
647
648     # 3) Scale the translation vector
649     t = t_tilde * alpha # Shape: [N, 3]
650
651     # 4) Form the local rigid body update g = (R, t) and compose it
652     # with the running pose g_initial = (R_0, t_0)
653     R_new, t_new = g_initial
654     R_new = R_0 @ R
655     t_new = (R_0 @ t.unsqueeze(-1)) + t_0
656
657     # 5) Apply the new pose g_new to the local template to get final coordinates
658     # Unsqueeze t_new for broadcasting over the 3 atoms in X_local
659     X_pred = (R_new @ X_local.T).transpose(-1, -2) + t_new.transpose(-1, -2)
660
661     return X_pred, (R_new, t_new) # Return coords and the updated pose
```

---

**Algorithm 2** Pseudocode for backbone *distance*, *direction* and *MSE* losses

```

648
649
650
651 def compute_backbone_vectors(coords: Tensor[L, 3, 3]) -> Tensor[L, 6, 3]:
652     """
653     coords[... , 0/1/2] = (N, CA, C) atom positions for each residue.
654     Returns six normalised vectors per residue:
655         1) N → CA          4) -n_CA   (binormal of NCA plane)
656         2) CA → C          5) n_N     (binormal of CN_prev-N plane)
657         3) C → N_{i+1}    6) n_C     (binormal of CA C N_{i+1} plane)
658     """
659     # bond vectors -----
660     v1 = coords[:, 1] - coords[:, 0]          # N → CA
661     v2 = coords[:, 2] - coords[:, 1]          # CA → C
662     v3 = pad_end(coords[1:, 0] - coords[:-1, 2]) # C → N_{i+1}
663     v4 = -torch.cross(v1, v2)                 # -n_CA
664     v5 = torch.cross(v3, v1)                  # n_N
665     v6 = torch.cross(v2, v3)                  # n_C
666
667     V = torch.stack([v1, v2, v3, v4, v5, v6], dim=1)
668     return F.normalize(V, dim=-1)             # shape = [L, 6, 3]
669
670
671 def backbone_distance_loss(x_pred, x_true, mask):
672     """
673     x_pred/x_true : [L, 3, 3] predicted & native backbone coords
674     mask          : [L]      boolean, valid residues
675     """
676     P = x_pred[mask, :3].reshape(-1, 9)      # flatten 3 atoms → 9-D point
677     T = x_true[mask, :3].reshape(-1, 9)
678
679     D_pred = pairwise_l2(P, P)                # (N, N) distance matrix
680     D_true = pairwise_l2(T, T)
681
682     diff = (D_pred - D_true).pow(2).clamp(max=25.)
683     return diff.mean()
684
685
686 def backbone_direction_loss(x_pred, x_true, mask):
687     """
688     Uses six orientation vectors per residue as returned by
689     'compute_backbone_vectors'.
690     """
691     V_pred = compute_backbone_vectors(x_pred[mask, :3])
692     V_true = compute_backbone_vectors(x_true[mask, :3])
693
694     # Pairwise dot-product tensors: S_{ij}^{kl} = v_i^{(k)} · v_j^{(l)}
695     S_pred = torch.einsum('ikd,jld->ijkl', V_pred, V_pred)
696     S_true = torch.einsum('ikd,jld->ijkl', V_true, V_true)
697
698     diff = (S_pred - S_true).pow(2).clamp(max=20.)
699     return diff.mean()
700
701
702 def aligned_mse_loss(x_pred, x_true, mask):
703     """
704     Mean-squared error after optimal rigid alignment
705     (Kabsch) between predicted and native backbone coords.
706
707     Args
708     ----
709     x_pred : Tensor[L, 3, 3] # predicted (N, CA, C)
710     x_true : Tensor[L, 3, 3] # ground truth
711     mask   : BoolTensor[L]   # valid residues
712     """
713     # 1) select masked residues and flatten atoms ----- #
714     P = x_pred[mask].reshape(-1, 3) # (3N, 3)
715     T = x_true[mask].reshape(-1, 3) # (3N, 3)
716
717     # 2) best-fit rotation / translation via Kabsch ----- #
718     R, t = kabsch(P, T) # R ∈ SO(3), t ∈ ℝ³
719     T_aln = (T @ R.T) + t # align native coords
720
721     # 3) mean-squared error ----- #
722     mse = ((P - T_aln).pow(2)).mean()
723     return mse

```

**Algorithm 3** Pseudocode for rigid alignment via Kabsch

```

def kabsch_align(A: Tensor[N, 3],
                B: Tensor[N, 3],
                allow_reflections : bool = False):
    """
    Optimal rigid alignment (rotation R, translation t) of
    point cloud A onto reference cloud B.

    A, B : [N, 3] float32 -- corresponding coordinates
    """
    # 1) centre the clouds ----- #
    centroid_A = A.mean(0); centroid_B = B.mean(0)
    AA = A - centroid_A      # centred source
    BB = B - centroid_B      # centred target

    # 2) SVD of covariance ----- #
    H = AA.T @ BB           # 3x3
    U, _, Vh = torch.linalg.svd(H)      # H = U S Vh

    # 3) ensure proper rotation (det = +1) ----- #
    if not allow_reflections and torch.det(Vh.T @ U.T) < 0:
        Vh[-1, :] *= -1

    R = Vh.T @ U.T          # rotation
    t = centroid_B - R @ centroid_A      # translation

    return (R @ A.T).T + t      # aligned A

```

## A.2 DATASETS

*Test A (low-confidence)* evaluates robustness under structural uncertainty using a curated subset of AlphaFold2 (AF2) models for SwissProt proteins. We enumerated all AF2-predicted structures available for SwissProt entries in AFDB and retained only those with average pLDDT  $\leq 70$  (AFDB’s 0–100 scale), discarding higher-confidence models. The surviving records were *cross-referenced* against the PDB to annotate the presence of experimental structures for the same proteins. Multimer entries were decomposed into per-chain monomer samples, after which we removed chains shorter than 25 amino acids and retained only chains with fewer than 25 consecutive missing residues. A final deduplication step produced 33 112 single-chain samples for Test A.

*Test B (species-diverse, taxonomic shift)* assesses generalization under distributional shift in species rather than structure confidence. We assembled a broad species roster (expanded to 100 and finalized at 97 species) and targeted roughly 2 500 proteins per species from AFDB without pLDDT filtering to reflect natural variability, dropping species with fewer than 1 000 available structures. After aggregation, per-chain conversion, and deduplication, this suite comprised 19 353 single-chain samples.

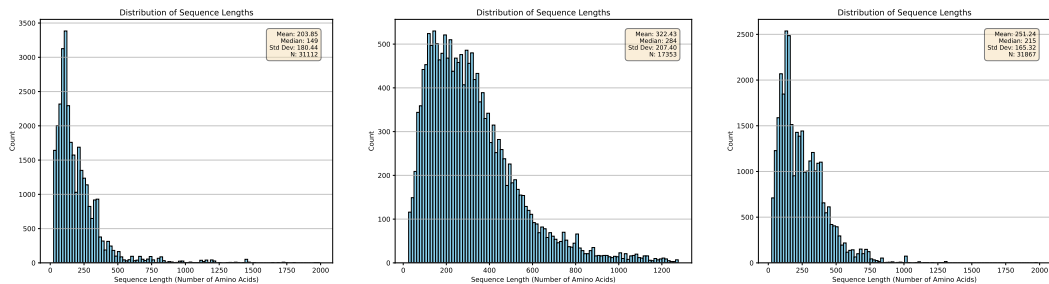


Figure 3: Distribution of sequence lengths in the three test sets: (left) Test set A, (middle) Test set B, and (right) Test set C. The majority of samples have sequence length  $< 1280$  amino acids.

*Test C (high-confidence, under-represented taxa)* measures performance on high-confidence structures from taxa under-represented in SwissProt. We selected species with low SwissProt coverage and required high average pLDDT ( $> 95$ ). The selected taxonomies were cross-referenced against

the PDB to find the presence of experimental structures for the proteins from those species. Applying the same preprocessing and deduplication as Test A above yielded 33 867 single-chain samples.

To obtain a balanced validation set independent of training, we first constructed three targeted hold-out suites and then set aside a fixed 2 000 randomly sampled chains from each to form a combined 6 000-chain validation set; the remaining examples constitute Test sets A, B, and C, respectively. The distribution of length sequences in the test sets is displayed in Figure 3.

From the zero-shot set, we defined interpretable structural subsets based on each protein’s secondary-structure composition. Coil-rich (D) includes sequences with at least 60% coil/disordered content. Mainly  $\alpha$  (A) requires an  $\alpha$ -helical fraction of at least 45% and an  $\alpha - \beta$  difference of at least 15 percentage points; mainly  $\beta$  (B) requires a  $\beta$ -strand fraction of at least 25% and a  $\beta - \alpha$  difference of at least 10 points. Balanced (AB) captures mixed folds with  $\alpha \geq 25\%$ ,  $\beta \geq 15\%$ , and  $|\alpha - \beta| \leq 10$  points. We also report two data-quality subsets: sequences with no missing residues and sequences with at least one missing residue. Groups are defined independently (i.e., they may overlap).

### A.3 EXPERIMENTS

Table 5: Training hyperparameters for Stages 1 and 2.

Hyperparameter	Stage 1	Stage 2
Optimizer type	AdamW	
$\beta_1$	0.9	
$\beta_2$	0.98	
$\epsilon$	$10^{-7}$	
Weight decay	$10^{-3}$	
Gradient clipping (max $L_2$ norm)	1.0	
Gradient accumulation	1	
8-bit parameters	✗	✓
Batch size (per GPU)	16	4
Learning-rate strategy	Cosine annealing	
Base learning-rate	$1e - 4$	$5e - 5$
Min learning-rate	$1e - 6$	
Warm-up steps	16 000	
Total steps	375 000	1 500 000
MSE coefficient ( $\lambda_{\text{MSE}}$ )	$1 \times 10^{-3}$	$5 \times 10^{-3}$
Backbone distance coefficient ( $\lambda_{\text{dist}}$ )	$10^{-2}$	
Backbone direction coefficient ( $\lambda_{\text{dir}}$ )	$5 \times 10^{-2}$	

This subsection compiles all experimental assets: Table 5 lists the full optimization schedule and training hyperparameters for Stages 1–2; Table 6 specifies the exact model configuration (GCPNet, VQ, and Transformer stacks). Figure 4 reports RMSD error distributions on CASP14/15/16 and CAMEO2024 across methods; Figure 5 shows qualitative robustness to contiguous missing-residue segments from benchmark structures; Figure 6 plots sequence length versus reconstruction RMSD with a least-squares fit on the validation set; Figure 7 visualizes the highest-RMSD (worst-case) examples per benchmark suite. Metrics use backbone atoms ( $C_\alpha$ ) after Kabsch alignment, and statistics are over the full, unfiltered test sets.

We maintain the codebook with the exponential–moving–average (EMA) variant of VQ-VAE: for each code  $k$  we keep decayed counts  $N_k$  and sums  $M_k$  of the encoder vectors assigned at the current step and update the entry by  $\mathbf{e}_k \leftarrow M_k / (N_k + \epsilon)$ , with decay  $\gamma$  (Stage 1: 0.99, Stage 2: 0.995; Table 6). Code entries are *not* updated by gradient; only the encoder receives the commitment penalty  $\mathcal{L}_{\text{commit}} = \beta \|\mathbf{z} - \text{sg}[\mathbf{e}_k]\|_2^2$ . We detect “dead” codes via an EMA count threshold (2 by default) and reinitialize them from recent encoder features using the same  $k$ -means seeding as at start-up. For the encoder pathway, we still use the rotation-trick Jacobian (Fifty et al., 2024) to pass informative gradients through the nearest-neighbour assignment; this affects only the encoder and does not modify the EMA update rule.

We set the coefficients  $\lambda$  (Equation 8) by monitoring the per-term gradient norms and equalizing their contribution to shared parameters. This gradient-norm balancing prevents any single term (e.g., direction or distance) from dominating updates and yields stable, faster convergence.

Throughout all experiments, sequences longer than the model’s input window are truncated to the first  $L$  residues (Stage 1:  $L = 512$ , Stage 2:  $L = 1,280$ ), and all metrics are computed on the truncated region.

Table 6: Configurations of GCP-VQVAE architecture.

Hyperparameter	Stage 1	Stage 2
GCPNet Encoder		
Initialization	(Jamasp et al., 2024)	Stage 1
Input dimension (node scalars $h$ )	6	
Input dimension (node vectors $\chi$ )	3	
Input dimension (edge scalars $e$ )	8	
Input dimension (edge vectors $\xi$ )	1	
Hidden dimension (node scalars $h$ )	128	
Hidden dimension (node vectors $\chi$ )	16	
Hidden dimension (edge scalars $e$ )	32	
Hidden dimension (edge vectors $\xi$ )	4	
Number of layers	6	
SE(3) equivariance	✓	
Vector Quantization		
Initialization	Random	Stage 1
Hidden dimension	256	
Vocab size	4 096	
EMA Decay	0.99	0.995
Threshold EMA dead code	2	
Commitment weight ( $\lambda_{\text{commit}}$ )	0.5	0.25
Orthogonal reg weight ( $\lambda_{\text{orth}}$ )	10	
Orthogonal reg max codes	512	
Orthogonal reg active codes only	✓	
Rotation trick	✓	
KMeans initialization	✓	
KMeans iteration	10	
Transformer Encoder		
Initialization	Random	Stage 1
Hidden dimension	1024	
FF Multiplier	4×	
Blocks	12	
Query heads	12	
Key-value heads	3	
Positional encoding	RoPE	
Query-key norm	✓	
Pre-norm	✓	
Transformer Decoder		
Initialization	Random	Stage 1
Hidden dimension	1024	
FF Multiplier	4×	
Blocks	16	
Query heads	16	
Key-value heads	1	
Positional encoding	RoPE	
Query-key norm	✓	
Pre-norm	✓	

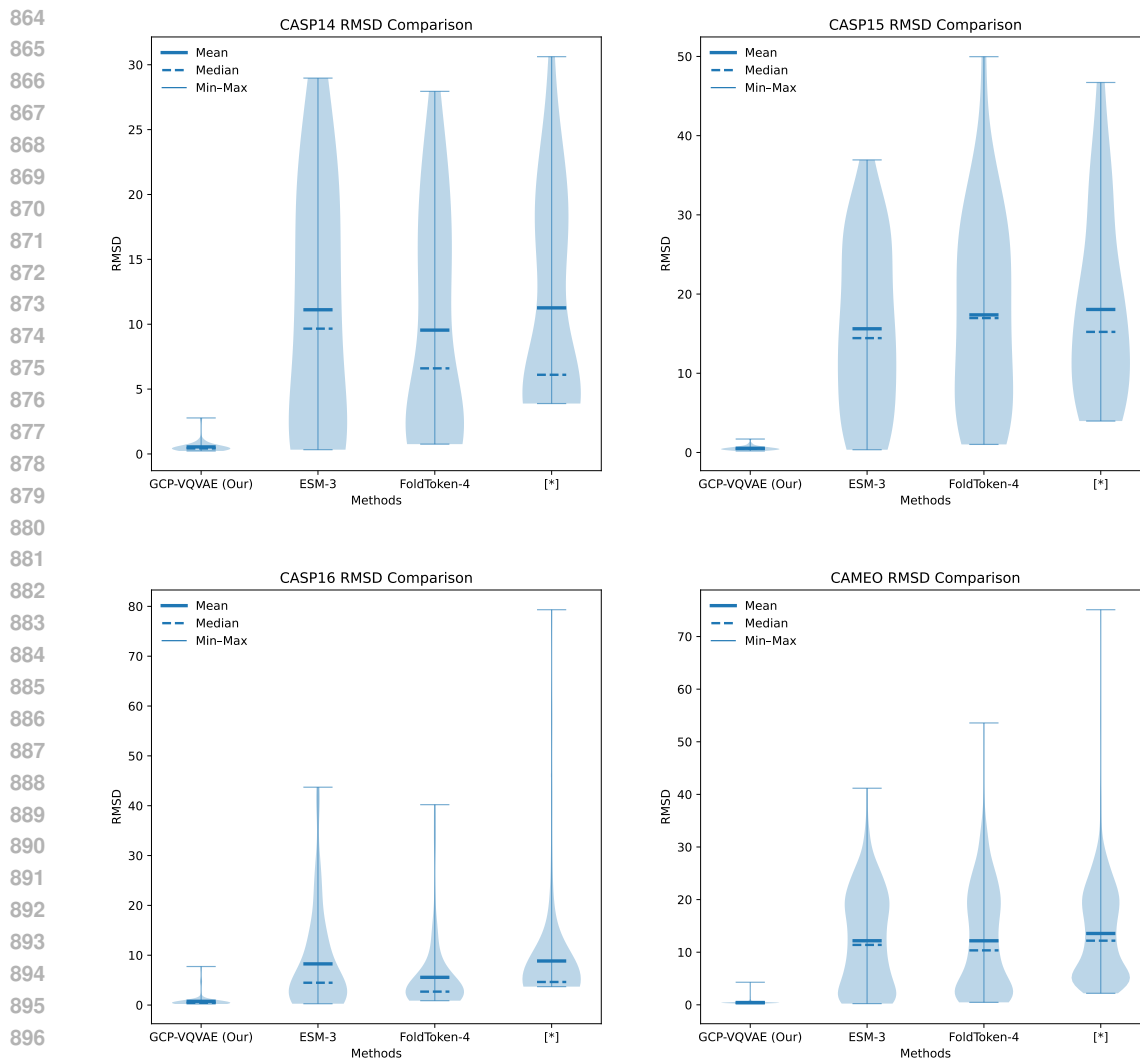
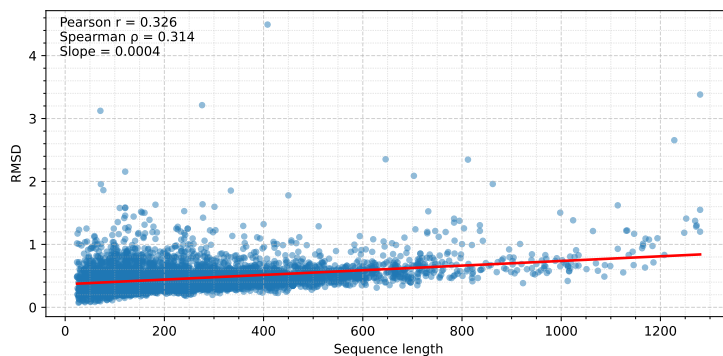


Figure 4: RMSD ( $\text{\AA}$ ) error distributions on CASP14, CASP15, CASP16, and CAMEO2024. Violin plots show the density of errors per method; horizontal markers denote the *mean* (solid), *median* (dashed), and *min-max* range (thin). Methods compared: GCP-VQVAE (ours), ESM-3 VQ-VAE, FoldToken-4, and the Structure Tokenizer of Gaujac et al. (2024) (shown as [\*]). For Gaujac et al. (2024), samples outside its supported length range (50–512 residues) are excluded; other methods use the full sets.



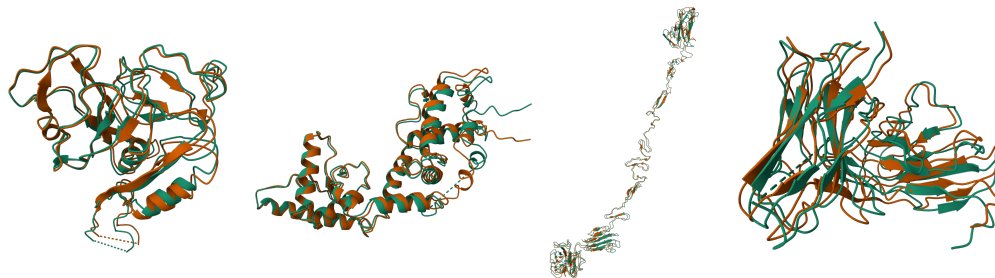
Figure 5: Superposition of GCP-VQVAE reconstructions (orange) with native backbones (green) for three proteins drawn from our external benchmark suites. Blue circles mark contiguous missing-residue segments (dashed guides span regions without deposited atoms); the model tracks the native structure outside the gaps.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929



930 Figure 6: Sequence length vs. reconstruction error (validation set). Each point is one protein; the  
931 red line is a least-squares fit. The dependence on length is modest (Pearson  $r = 0.326$ , Spearman  
932  $\rho = 0.314$ ) with a small slope of  $\approx 4 \times 10^{-4}$  Å/residue ( $\sim 0.4$  Å per 1k residues).  
933

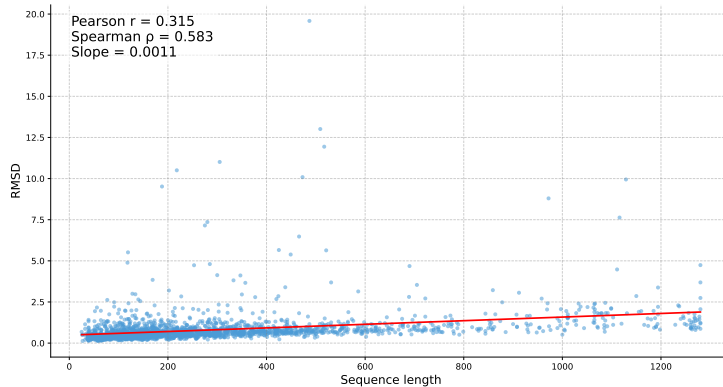
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951



952 Figure 7: Superposition of GCP-VQVAE (orange) and native backbones (green) for the  
953 highest-backbone-RMSD example in each suite, left→right: CASP14, CASP15, CASP16,  
954 CAMEO2024.  
955

#### 956 A.4 ZERO-SHOT PERFORMANCE

957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971



972 Figure 8: RMSD as a function of sequence length (number of amino acids) on the zero-shot set. Each  
973 dot is one protein. The red line is an ordinary-least-squares (OLS) fit; the panel reports Pearson  $r$ ,  
974 Spearman  $\rho$ , and the fitted slope (RMSD units per residue).  
975

972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025

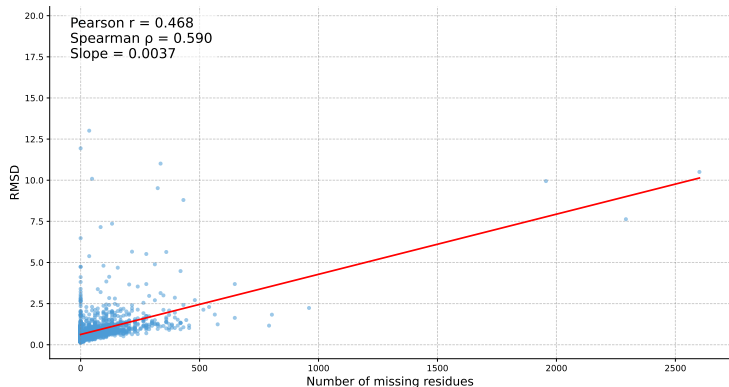


Figure 9: RMSD as a function of the number of missing residues on the zero-shot set. Points show individual proteins; the red line is an OLS fit. The panel reports Pearson  $r$ , Spearman  $\rho$ , and the fitted slope (RMSD units per missing residue).

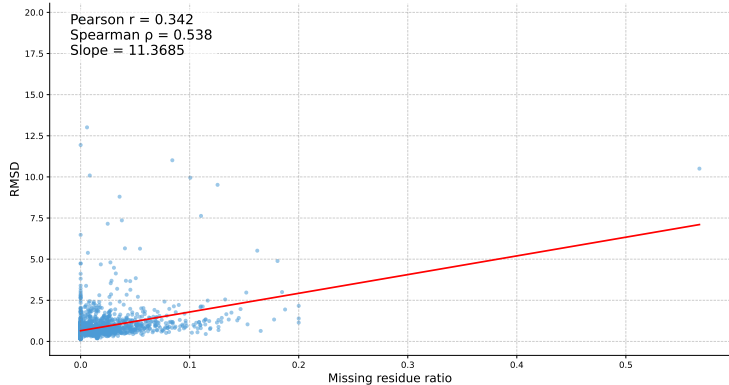


Figure 10: RMSD as a function of the missing-residue ratio on the zero-shot set. The red line is an OLS fit; the panel reports Pearson  $r$ , Spearman  $\rho$ , and the fitted slope (RMSD units per unit fraction; 1.0=100%).

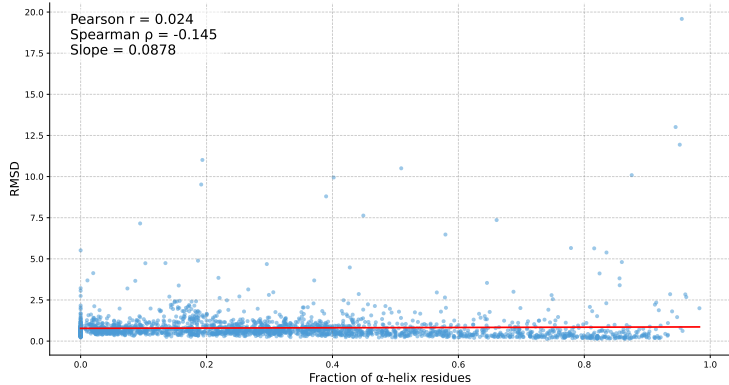


Figure 11: RMSD as a function of the fraction of  $\alpha$ -helix residues on the zero-shot set. The red line indicates an OLS fit. The panel reports Pearson  $r$ , Spearman  $\rho$ , and the fitted slope (RMSD units per unit fraction; 1.0=100%).

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

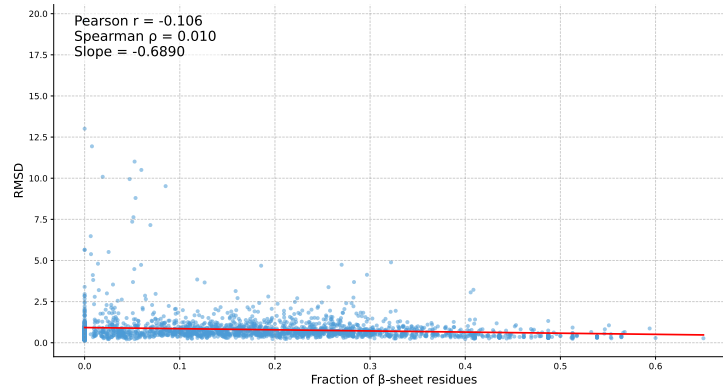


Figure 12: RMSD as a function of the fraction of  $\beta$ -sheet residues on the zero-shot set. The red line indicates an OLS fit. The panel reports Pearson  $r$ , Spearman  $\rho$ , and the fitted slope (RMSD units per unit fraction; 1.0=100%).

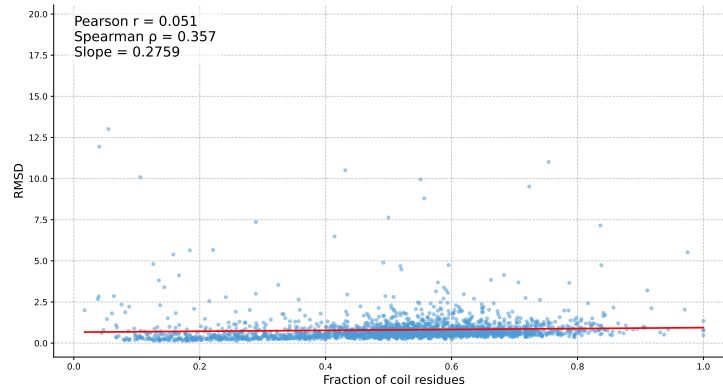


Figure 13: RMSD as a function of the fraction of coil residues on the zero-shot set. The red line indicates an OLS fit. The panel reports Pearson  $r$ , Spearman  $\rho$ , and the fitted slope (RMSD units per unit fraction; 1.0=100%).

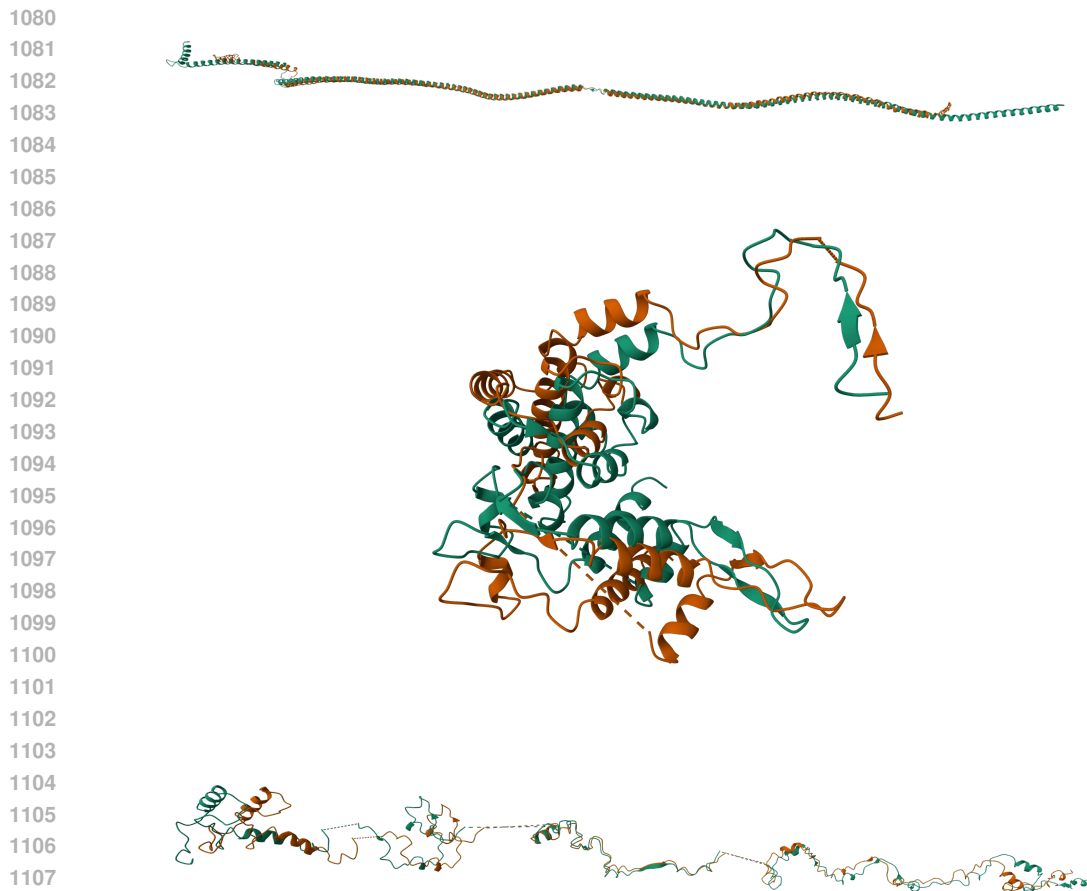


Figure 14: Worst zero-shot reconstructions (orange: GCP-VQVAE; green: native). Top: RMSD 19.5731, TM-score 0.6935, length 487, NaN residues 0. Middle: RMSD 12.276, TM-score 0.575, length 218, NaN residues 2601. Bottom: RMSD 11.0486, TM-score 0.4766, length 305, NaN residues 336. The middle and lower cases exhibit an extreme number of NaN residues introduced by our current NaN-handling/augmentation pipeline, which forces the model to contend with long masked segments and degrades reconstruction accuracy; improving the NaN strategy (e.g., capping gap ratio, better segmentation, and length-aware masking) should mitigate this failure mode.

#### A.5 COMPRESSION CALCULATION

Each residue is encoded by one code from a 4096-entry book, i.e.,  $\log_2(4096) = 12$  bits = 1.5 bytes/residue. For  $L = 512$  residues, the token stream is  $512 \times 1.5 = 768$  bytes. Raw backbone coordinates (N,  $C_\alpha$ , C) stored as 32-bit floats require  $3 \text{ atoms} \times 3 \text{ coords} \times 4 \text{ bytes} = 36$  bytes/residue, i.e.,  $36 \times 512 = 18,432$  bytes. If we include a tiny global pose header (e.g., rotation+translation) of  $\approx 36$  bytes, the coded footprint is  $768+36 = 804$  bytes. Thus the compression ratio is  $18\,432/804 \approx 22.9\times$  (approximately  $24\times$  if the pose header is omitted). This estimate concerns backbone only; metadata and format containerization add negligible overhead relative to the raw-float baseline.

#### A.6 LLM USAGE

In this manuscript, we used large language models only for copy-editing: improving grammar, clarity, and style of author-written text.