

PHYSICS-ALIGNED DECODING (PAD) FOR DISCRETE PROTEIN STRUCTURE REPRESENTATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Discrete representations learned by deep autoencoders are increasingly reused as intermediate state spaces in generative, conditional, and autoregressive models. In this work, we empirically identify an objective-level failure mode in discrete protein structure tokenizers trained with reconstruction-aligned losses: despite low global reconstruction error, learned tokens encode locally unphysical geometry, including covalent distortions and steric clashes. We show that these violations are deterministic and persistent under reuse. We test the hypothesis that this behavior arises from objective misspecification rather than architectural limitations, and introduce Physics-Aligned Decoding (PAD), a minimal intervention that augments reconstruction objectives with differentiable physical priors. Without changing architecture or regenerating the codebook, PAD reshapes token semantics and restores physical validity while preserving reconstruction accuracy. Our results highlight how loss geometry determines representation semantics, and demonstrate the importance of objective alignment when discrete representations are reused beyond static reconstruction.

1 INTRODUCTION

Deep representation learning has traditionally emphasized reconstruction fidelity as a proxy for representation quality. In vector-quantized autoencoders (VQ-VAEs) (van den Oord et al., 2017; Razavi et al., 2019), this paradigm enables compression of high-dimensional inputs into discrete latent codes that can be reused for downstream modeling. Recently, this approach has been extended to protein structure, where discrete structural tokens serve as interfaces to generative models, conditional editors, and large language models (Gao et al., 2024; Lin et al., 2025; Yuan et al., 2025).

In these settings, discrete tokens are no longer treated as passive compression artifacts. Instead, they define a latent *state space* over which downstream models operate. This shift raises a critical question: what semantic content do these discrete states encode? When learned representations are reused as intermediate states, errors that are benign under static reconstruction may become amplified or irreparable.

Protein structures provide a stringent test case. Valid conformations lie on a highly constrained physical manifold defined by covalent geometry, steric exclusion, and torsional regularities (Engh & Huber, 1991; Ramachandran et al., 1963; Dunbrack, 2002). Yet most protein tokenizers are trained using reconstruction-aligned objectives—including coordinate-level L_2 losses or frame-aligned variants such as FAPE (Jumper et al., 2021)—that are insensitive to localized physical violations.

In this work, we show that this mismatch leads to a systematic failure mode: discrete tokens encode unphysical micro-geometry despite achieving low global reconstruction error. We argue that this behavior reflects objective misspecification rather than architectural deficiency, and we test this hypothesis through controlled intervention.

2 RECONSTRUCTION-ALIGNED DECODING AND ITS BLIND SPOTS

Reconstruction-aligned objectives minimize additive losses over geometric features such as Cartesian coordinates, distances, or frames (Ingraham et al., 2019; Gao et al., 2024). These losses induce

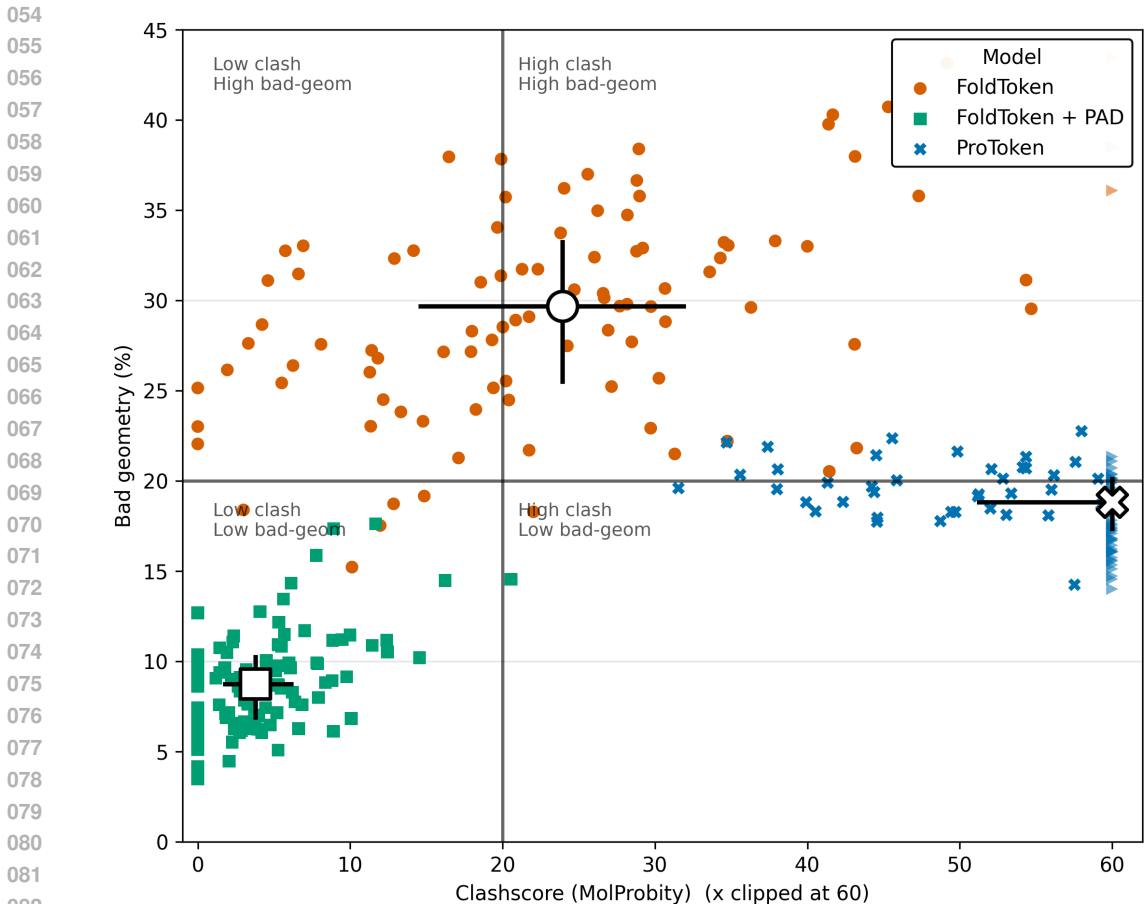


Figure 1: **Deterministic physical violations under reconstruction-aligned decoding.** Clashescore versus bad-geometry fraction for decoded structures on held-out conformations. Despite low global error, reconstruction-aligned tokenizers exhibit severe local violations. Physics-Aligned Decoding (PAD) restores physical validity without architectural changes.

smooth, symmetric penalties that treat all deviations as comparable, regardless of physical interpretation.

In physically constrained systems, this symmetry is problematic. A small displacement into empty space and an equally small displacement into another atom’s excluded volume incur similar penalties under an L_2 loss, even though the latter configuration is physically forbidden (Bondi, 1964; Chen et al., 2010). As a result, the optimizer may reduce average reconstruction error by introducing localized covalent distortions or steric overlaps.

Figure 1 illustrates this failure mode across representative tokenizers. Notably, violations are reproducible under a single deterministic encode–decode pass, indicating that they reflect stable optima of the training objective rather than stochastic decoding noise.

3 SCIENTIFIC-METHOD FRAMING

We cast this observation as a hypothesis-driven empirical study:

Observation. Reconstruction-aligned decoding produces discrete tokens whose decoded structures exhibit local physical violations.

Hypothesis. These violations arise from objective misspecification: symmetric reconstruction losses fail to encode asymmetric physical constraints such as steric exclusion.

108
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
160
161

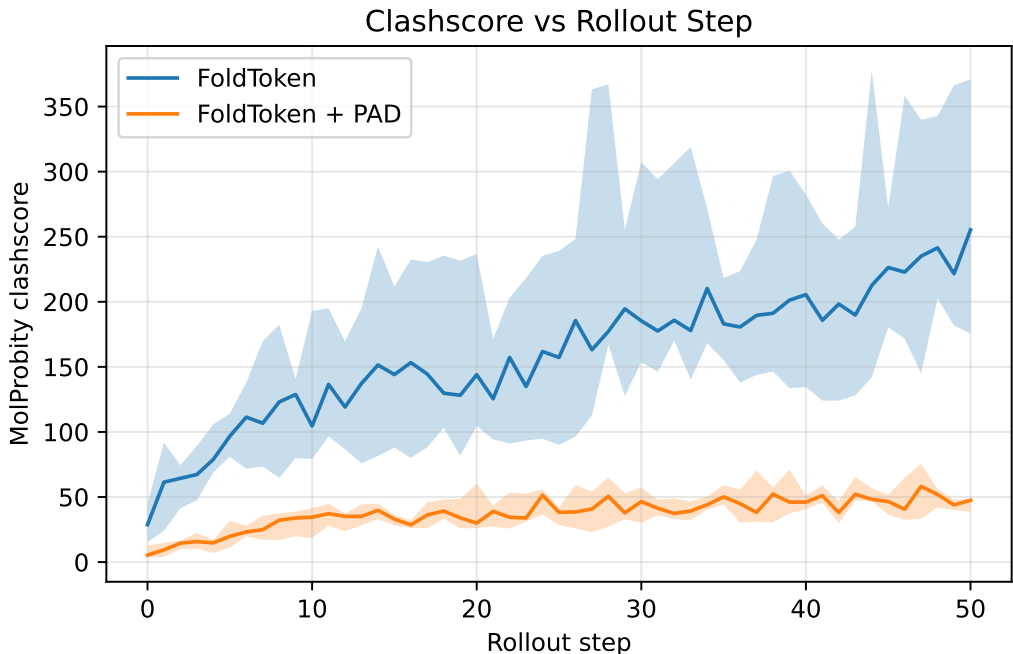


Figure 2: **Token reuse reveals semantic stability.** Median clashscore over multi-step autoregressive rollouts. Reconstruction-aligned tokens rapidly accumulate physical violations, whereas PAD tokens remain stable under reuse.

Prediction. If the hypothesis is correct, modifying only the decoding objective—while keeping architecture and codebook fixed—should reshape token semantics and eliminate violations under both single-step decoding and downstream reuse.

Test. We implement Physics-Aligned Decoding (PAD) and evaluate physical validity under controlled conditions.

4 PHYSICS-ALIGNED DECODING

Physics-Aligned Decoding augments the reconstruction objective with differentiable physical priors that encode local feasibility. The decoding loss takes the form

$$\mathcal{L}_{\text{PAD}} = \mathcal{L}_{\text{rec}} + \lambda_{\text{geom}}\mathcal{L}_{\text{geom}} + \lambda_{\text{vdw}}\mathcal{L}_{\text{vdw}},$$

where $\mathcal{L}_{\text{geom}}$ enforces empirical covalent and torsional constraints (Engh & Huber, 1991; Shapovalov & Dunbrack, 2011) and \mathcal{L}_{vdw} penalizes steric overlap using smooth van der Waals-inspired terms (Bondi, 1964). Full definition can be found in the appendix.

Crucially, PAD does not alter model architecture or regenerate the discrete codebook. We implement it via parameter-efficient fine-tuning using LoRA (Hu et al., 2021), ensuring that observed changes reflect objective-induced semantic reorganization rather than increased model capacity.

5 RESULTS

Applying PAD shifts decoded structures into physically admissible regimes while preserving reconstruction fidelity (Fig. 1). To test whether this improvement reflects genuine semantic change rather than post hoc smoothing, we evaluate token reuse in downstream generative settings.

Figure 2 shows that under autoregressive reuse, reconstruction-aligned tokens accumulate steric clashes rapidly, whereas PAD tokens maintain bounded physical quality. Importantly, improvements

162 are visible at early rollout steps, indicating that unphysical micro-geometry is encoded directly into
163 baseline token semantics.

164 165 6 DISCUSSION 166

167 Our results demonstrate that reconstruction fidelity alone is an insufficient training signal when dis-
168 crete representations are reused as state spaces. The geometry of the loss function determines which
169 regions of representation space are accessible, and symmetric objectives may admit semantically
170 invalid optima in constrained domains.

171 While our experiments focus on protein structure, the underlying principle is general. Any discrete
172 representation learned on a constrained manifold—and reused for planning, rollouts, or conditional
173 editing—may require objective terms that encode feasibility constraints rather than average recon-
174 struction error.

175 176 7 LIMITATIONS 177

178 Physics-Aligned Decoding enforces local feasibility but does not model full molecular energetics
179 or long-range interactions. More broadly, this work does not propose a universal objective, but
180 illustrates how objective misspecification can be diagnosed and corrected using controlled empirical
181 analysis.

182 183 REFERENCES 184

- 185 A. Bondi. van der waals volumes and radii. *Journal of Physical Chemistry*, 68(3):441–451, 1964.
- 186 Vincent B. Chen, William B. Arendall, Jeffrey J. Headd, et al. Molprobity: all-atom structure
187 validation for macromolecular crystallography. *Acta Crystallographica Section D: Biological*
188 *Crystallography*, 66(1):12–21, 2010. doi: 10.1107/S0907444909042073.
- 189 Roland L Dunbrack. Rotamer libraries in the 21st century. *Current Opinion in Structural Biology*,
190 12(4):431–440, 2002. ISSN 0959-440X. doi: [https://doi.org/10.1016/S0959-440X\(02\)00344-5](https://doi.org/10.1016/S0959-440X(02)00344-5).
- 191 R. A. Engh and R. Huber. Accurate bond and angle parameters for x-ray protein structure refinement.
192 *Acta Crystallographica Section A*, 47:392–400, 1991.
- 193 Zhangyang Gao, Cheng Tan, Jue Wang, Yufei Huang, Lirong Wu, and Stan Z. Li. Foldtoken:
194 Learning protein language via vector quantization and beyond. 2024. URL <https://arxiv.org/abs/2403.09673>.
- 195 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
196 and Weizhu Chen. Lora: Low-rank adaptation of large language models. 2021. URL <https://arxiv.org/abs/2106.09685>.
- 197 John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph-
198 based protein design. In *Advances in Neural Information Processing Systems*, volume 32. Curran
199 Associates, Inc., 2019.
- 200 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
201 Kathryn Tunyasuvunakool, Russ Bates, et al. Highly accurate protein structure prediction with
202 AlphaFold. *Nature*, 596(7873):583–589, 2021. doi: 10.1038/s41586-021-03819-2.
- 203 Xiaohan Lin, Zhenyu Chen, Yanheng Li, et al. Unifying sequence-structure coding for advanced
204 protein engineering via a multimodal diffusion transformer. *Chem. Sci.*, 16:11087–11102, 2025.
205 doi: 10.1039/D5SC02055G. URL <http://dx.doi.org/10.1039/D5SC02055G>.
- 206 Christine A. Orengo, Andrew D. Michie, et al. CATH: a hierarchic classification of protein domain
207 structures. *Structure*, 5(8):1093–1108, 1997. doi: 10.1016/S0969-2126(97)00260-8.
- 208 G. N. Ramachandran, C. Ramakrishnan, and V. Sasisekharan. Stereochemistry of polypeptide chain
209 configurations. *Journal of Molecular Biology*, 7:95–99, 1963.

216 Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with
217 vq-vae-2. 2019. URL <https://arxiv.org/abs/1906.00446>.
218

219 Maxim V. Shapovalov and Roland L. Dunbrack. A smoothed backbone-dependent rotamer library
220 for proteins derived from adaptive kernel density estimates and regressions. *Structure*, 19(6):
221 844–858, 2011.

222 Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learn-
223 ing. In *NeurIPS*, 2017.
224

225 Yann Vander Meersche, Gabriel Cretin, Aria Gheeraert, Jean-Christophe Gelly, and Tatiana Ga-
226 lochkina. Atlas: protein flexibility description from atomistic molecular dynamics simulations.
227 *Nucleic Acids Research*, 52(D1):D384–D392, 11 2023. ISSN 0305-1048. doi: 10.1093/nar/
228 gkad1084. URL <https://doi.org/10.1093/nar/gkad1084>.

229 Xinyu Yuan, Zichen Wang, Marcus Collins, and Huzefa Rangwala. Protein structure tokenization:
230 Benchmarking and new recipe. 2025. URL <https://arxiv.org/abs/2503.00089>.
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

A DETAILED VAN DER WAALS AND STERIC POTENTIALS

This appendix provides the explicit functional forms of the differentiable steric-exclusion and van der Waals terms used in Physics-Aligned Decoding (PAD). The forms are inspired by Rosetta-style Lennard–Jones potentials, but are modified to ensure smooth, numerically stable gradients for end-to-end training (soft-floored distances, a linear continuation at short range, and a smooth cutoff for attraction).

A.1 SOFT-FLOORED INTERATOMIC DISTANCES

For atoms indexed by p, q with decoded coordinates \hat{x}_p, \hat{x}_q , define

$$d_{pq} = \|\hat{x}_p - \hat{x}_q\|.$$

To avoid unbounded gradients as $d_{pq} \rightarrow 0$, we use a soft-floored distance

$$\tilde{d}_{pq} = \sqrt{d_{pq}^2 + d_0^2}, \tag{1}$$

with small constant $d_0 > 0$. All Lennard–Jones terms below are evaluated using \tilde{d}_{pq} .

A.2 DIFFERENTIABLE STERIC CLASH PENALTY

A coarse differentiable clash barrier is applied over backbone atoms and pseudo- C_β atoms with radii r_p . For atom pairs $p < q$, we penalize overlaps using

$$\mathcal{L}_{\text{clash}}(\hat{X}) = \mathbb{E}_{p < q} \left[\left(\max\{0, (r_p + r_q + \delta) - \tilde{d}_{pq}\} \right)^2 \right], \tag{2}$$

where $\delta > 0$ is a small margin. Using \tilde{d}_{pq} (rather than d_{pq}) ensures bounded gradients even under severe overlap. This definition is used consistently throughout the paper; the squared term applies to distance penetration quadratically.

A.3 ROSETTA-INSPIRED VAN DER WAALS ENERGY

We include a smooth van der Waals loss \mathcal{L}_{vdw} comprising a repulsive (fa_rep) and an attractive (fa_atr) component. For atom pairs p, q , let σ_{pq} and ϵ_{pq} denote Lennard–Jones radius and well-depth parameters. Define

$$x_{pq} = \frac{\sigma_{pq}}{\tilde{d}_{pq}}.$$

Shifted Lennard–Jones form. We use the shifted potential

$$E_{pq}^{\text{LJ}^+}(\tilde{d}) = \epsilon_{pq} (x_{pq}^{12} - 2x_{pq}^6 + 1), \tag{3}$$

which is zero at $\tilde{d} = \sigma_{pq}$ and increases sharply for $\tilde{d} < \sigma_{pq}$.

Repulsive term (fa_rep). Let $\tilde{d}_t = 0.6\sigma_{pq}$. For $\tilde{d} \leq \tilde{d}_t$, we use a linear continuation

$$E_{pq}^{\text{lin}}(\tilde{d}) = a_{pq}\tilde{d} + b_{pq},$$

where (a_{pq}, b_{pq}) are chosen to match both the value and slope of $E_{pq}^{\text{LJ}^+}$ at \tilde{d}_t . The repulsive energy is

$$E_{pq}^{\text{rep}}(\tilde{d}) = \begin{cases} E_{pq}^{\text{lin}}(\tilde{d}), & \tilde{d} \leq \tilde{d}_t, \\ E_{pq}^{\text{LJ}^+}(\tilde{d}), & \tilde{d}_t < \tilde{d} \leq \sigma_{pq}, \\ 0, & \tilde{d} > \sigma_{pq}. \end{cases} \tag{4}$$

324 **Attractive term** (fa_atr). Define the standard Lennard–Jones form

$$325 E_{pq}^{\text{LJ}}(\tilde{d}) = \epsilon_{pq}(x_{pq}^{12} - 2x_{pq}^6).$$

326 We apply a smooth cutoff between an inner radius d_{in} and an outer cutoff d_{out} using a cubic Hermite

327 switch

$$328 s(t) = 2t^3 - 3t^2 + 1, \quad t = \frac{\tilde{d} - d_{\text{in}}}{d_{\text{out}} - d_{\text{in}}}. \quad (5)$$

329 The attractive energy is

$$330 E_{pq}^{\text{atr}}(\tilde{d}) = \begin{cases} -\epsilon_{pq}, & \tilde{d} \leq \sigma_{pq}, \\ E_{pq}^{\text{LJ}}(\tilde{d}), & \sigma_{pq} < \tilde{d} \leq d_{\text{in}}, \\ E_{pq}^{\text{LJ}}(\tilde{d})s(t), & d_{\text{in}} < \tilde{d} \leq d_{\text{out}}, \\ 0, & \tilde{d} > d_{\text{out}}. \end{cases} \quad (6)$$

331 In our implementation, $(d_{\text{in}}, d_{\text{out}})$ are chosen as fixed multiples of σ_{pq} , and the switch ensures

332 continuous energy and gradients at the cutoff.

333 **Total van der Waals loss.** The full van der Waals loss sums over valid atom pairs, excluding

334 bonded and near-bonded interactions via an exclusion mask $\chi_{pq} \in \{0, 1\}$:

$$335 \mathcal{L}_{\text{vdw}}(\hat{X}) = w_{\text{rep}} \sum_{p < q} \chi_{pq} E_{pq}^{\text{rep}}(\tilde{d}_{pq}) + w_{\text{atr}} \sum_{p < q} \chi_{pq} E_{pq}^{\text{atr}}(\tilde{d}_{pq}). \quad (7)$$

336 Together, $\mathcal{L}_{\text{clash}}$ and \mathcal{L}_{vdw} provide complementary shaping: the clash barrier suppresses forbidden

337 overlaps, while the van der Waals term provides a smooth repulsive/attractive landscape that guides

338 decoded structures toward physically plausible local minima under differentiable optimization.

339 B AUTOREGRESSIVE ROLLOUT MODEL.

340 For the downstream multi-step reuse stress test (Sec. ??), we train a lightweight Transformer encoder

341 that predicts the next-frame token at each residue given the current-frame token sequence. Con-

342 cretely, the model embeds per-residue discrete tokens with an embedding dimension of $d = 256$,

343 adds sinusoidal positional encodings over residue index, and applies $L = 4$ TransformerEncoder

344 layers with $H = 8$ attention heads and feedforward width $4d$, followed by a linear classifier to K

345 token logits per residue. The model is trained with cross-entropy on next-frame prediction pairs

346 (z_t, z_{t+1}) extracted from tokenized MD trajectories, where $z_t \in \{0, \dots, K - 1\}^L$ is the full-length

347 residue token vector for a frame. We use an 80/20 split over trajectory frames (first $0.8T$ frames for

348 training, remaining $0.2T$ frames for seeding held-out rollouts), AdamW optimization with learning

349 rate 3×10^{-4} and weight decay 10^{-2} , batch size 16, gradient clipping at 1.0, and train for 10 epochs

350 (dropout 0.1; GELU; pre-norm). The vocabulary size K is inferred from the maximum token id

351 observed in the training frames ($K = \max(z_{\leq 0.8T}) + 1$) unless provided explicitly. All rollout

352 experiments use identical model architecture, training procedure, and sampling hyperparameters

353 for FoldToken and FoldToken+PAD, so differences in physical validity arise from token semantics

354 rather than downstream model differences.

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

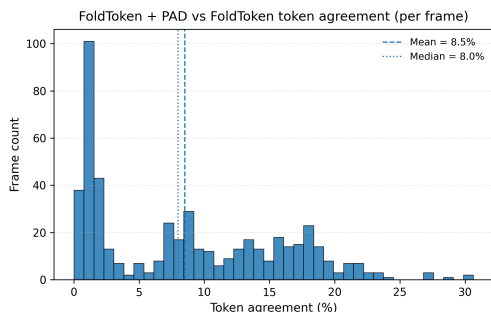
397

398

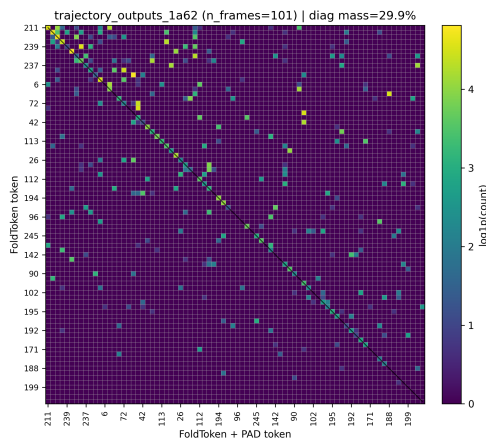
399

400

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392



(a) Per-frame token agreement between FoldToken and FoldToken+PAD, defined as the fraction of residues whose discrete token assignment is unchanged. Dashed and dotted lines indicate the mean (8.5%) and median (8.0%) agreement, respectively.



(b) Aggregated token transition matrix between FoldToken and FoldToken+PAD over 101 frames. Each entry counts how often a residue assigned token i under the baseline model is reassigned to token j under PAD. Color indicates $\log(1 + \text{count})$.

Figure 3: Aggregate and token-level views of token reassignment under Physics-Aligned Decoding. Left: histogram of per-frame token agreement between FoldToken and FoldToken+PAD, defined as the fraction of residues whose discrete token assignment is unchanged. Right: aggregated token transition matrix between FoldToken and FoldToken+PAD over 101 frames, restricted to the top- k most frequently used tokens ($k = 20$). Each entry (i, j) counts the number of residue positions for which token i under the baseline FoldToken model is reassigned to token j under PAD; color indicates $\log(1 + \text{count})$.

393
394
395
396
397
398
399
400
401
402
403
404
405
406
407

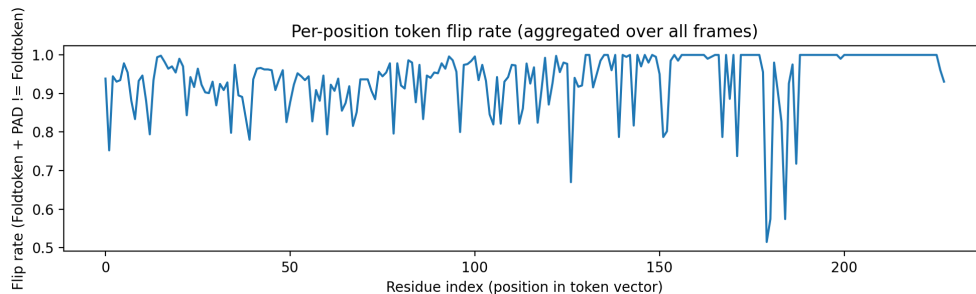


Figure 4: Per-position token flip rate aggregated across all frames, defined as the fraction of frames in which the FoldToken+PAD assignment differs from the baseline FoldToken assignment at a given residue index. While most positions exhibit high stability, several localized regions show consistently elevated flip rates, indicating position-dependent semantic reorganization under PAD.

417
418
419
420
421
422

423 D LORA FINE-TUNING CONFIGURATION

424
425
426
427
428
429

This appendix summarizes the LoRA configuration used for the geometry-push experiments reported in the main text. We list only the hyperparameters and architectural targets that materially affect the learned representation; standard training, logging, and data-loading details are omitted for brevity.

430
431

LoRA hyperparameters. We use LoRA with rank $r = 32$ and scaling factor $\alpha = 64$, with dropout 0.1. The base FoldToken model is frozen during fine-tuning, and no bias parameters are trained. LoRA adapters are trained in feature-extraction mode.

432 **Targeted modules.** LoRA adapters are applied broadly across the encoder and decoder, including
433 attention projections, feed-forward layers, and geometric feature pathways, while preserving the
434 original coordinate prediction heads. Specifically, adapters are attached to:

- 435 • the VQ projection layers and embedding MLPs,
- 436 • encoder attention (W_Q, W_K, W_V, W_O), feed-forward, and geometric feature modules,
- 437 • decoder GNN attention, feed-forward, and geometric feature modules,
- 438 • decoder coordinate and quaternion prediction projections.

439
440
441 **Optimization and geometry losses.** Fine-tuning uses a learning rate of 1×10^{-4} with cosine
442 scheduling and a warmup of 100 steps. Gradients are clipped to a norm of 2.0 and accumulated
443 over 12 steps. Geometry-aware losses are enabled throughout training, including bond length, bond
444 angle, Ramachandran, and ω -torsion penalties (each with weight 0.5), with an additional peptide
445 planarity term (weight 0.05). A van der Waals clash loss is enabled with a small weight (0.01),
446 ramped during training.

447 448 E TRAINING DATA, SPLITS, AND PROCEDURE

449
450 **Training data.** LoRA fine-tuning for Physics-Aligned Decoding is performed on molecular dy-
451 namics (MD) trajectory frames derived from the CATH domain dataset (Orengo et al., 1997). Train-
452 ing frames are extracted from equilibrium MD simulations of CATH domains, providing physically
453 realistic intermediate conformations beyond static experimental structures. Only backbone atoms
454 (N, C_α, C, O) and pseudo- C_β atoms are used during training, consistent with the representation
455 described in the main text.

456
457 **Data splits and coverage.** Trajectories are split at the protein level to avoid information leakage
458 across train and evaluation sets. Domains used for PAD fine-tuning are disjoint from those used in all
459 static reconstruction, trajectory-level, and autoregressive rollout evaluations reported in Section 4.
460 For downstream MD-based evaluations, additional out-of-distribution trajectories are drawn from
461 the ATLAS dataset (Vander Meersche et al., 2023), which is never used during fine-tuning.

462
463 **Frame sampling.** From each training trajectory, frames are subsampled uniformly in time to re-
464 duce temporal correlation and to expose the decoder to a diverse range of physically valid conforma-
465 tions. No frame-level augmentation or artificial perturbation is applied. Each training batch consists
466 of independently sampled frames rather than contiguous trajectory segments.

467
468 **Training procedure.** Fine-tuning is performed for a fixed number of epochs over the training
469 frames using the AdamW optimizer. The pretrained FoldToken encoder, decoder, and codebook
470 remain frozen except for the LoRA adapters described above. All reconstruction losses used in the
471 base tokenizer are retained unchanged, with additional physics-aligned terms applied only at the
472 decoding stage. Gradients are propagated through the straight-through estimator to the encoder,
473 allowing token assignment boundaries to adapt without regenerating the discrete codebook.

474
475 **Computational cost.** Because PAD is implemented via parameter-efficient fine-tuning, training
476 requires substantially less computation than full retraining. Inference-time decoding incurs only a
477 small constant-factor overhead from evaluating the additional geometry-aware loss terms, with no
478 change to model architecture or codebook size.

486 F SCIENCE OF DL IMPROVEMENT CHALLENGE SUBMISSION
487488 F.1 WHAT MODEL ARE YOU TARGETING?
489

490 *Provide a summary of the problem the deep net model is designed to solve. Good summaries should*
491 *outline the state of the literature, provide an overview that domain experts would consider reason-*
492 *able, and cite relevant sources.*
493

494 We target deep discrete representation models trained with reconstruction-aligned objectives, with a
495 particular focus on vector-quantized autoencoders (VQ-VAEs) used for protein structure tokeniza-
496 tion. Such models map high-dimensional geometric inputs to a finite set of discrete latent states,
497 enabling downstream generative, conditional, and autoregressive modeling.

498 In the protein domain, discrete structure tokenizers such as FoldToken and ProToken have been pro-
499 posed to serve as compact interfaces between three-dimensional molecular structure and sequence-
500 based models, including transformers and large language models. These models are typically trained
501 using reconstruction-aligned losses (e.g., coordinate-level L_2 loss or frame-aligned variants such as
502 FAPE) that emphasize global geometric fidelity.

503 While effective for static reconstruction benchmarks, these models are increasingly reused as latent
504 state spaces for downstream tasks, including conditional editing, autoregressive rollouts, and latent
505 dynamics modeling. In these regimes, discrete tokens no longer function merely as compression
506 artifacts, but as semantic states over which learning and inference are performed. Our work targets
507 this class of models and interrogates whether reconstruction-aligned training objectives produce
508 representations whose discrete states are semantically valid when reused beyond reconstruction.
509

510 F.2 HOW DO YOUR RESULTS CONTRIBUTE—OR COULD POTENTIALLY CONTRIBUTE—TO
511 UNDERSTANDING THESE MODELS?
512

513 *What aspects of the models become better understood thanks to your work?*
514

515 Our results clarify how the geometry of the training objective shapes the semantic content of discrete
516 representations. We show that reconstruction-aligned objectives can admit stable optima in which
517 discrete latent states encode locally invalid configurations, even when global reconstruction error
518 is low. These violations are deterministic, reproducible, and persist under reuse, indicating that
519 they arise from objective-level properties rather than stochastic decoding noise or insufficient model
520 capacity.

521 By isolating objective effects through a controlled intervention that modifies only the decoding loss
522 while keeping architecture and codebook fixed, we demonstrate that unphysical behavior is not an
523 inherent limitation of discrete autoencoders, but a consequence of symmetric reconstruction losses
524 that fail to encode asymmetric feasibility constraints. This provides empirical evidence that repre-
525 sentation failure modes can arise from loss misspecification even when models appear successful
526 under standard evaluation metrics.

527 More broadly, our work highlights that when learned representations are reused as state spaces,
528 reconstruction fidelity alone is an insufficient proxy for representational validity. Understanding
529 deep models in such settings requires analyzing not only architecture and capacity, but also how
530 objective geometry constrains which regions of representation space are accessible.
531

532 F.3 HOW DO YOU EXPECT YOUR SUBMISSION TO INFLUENCE FUTURE WORK?
533

534 *Propose ways in which your insights, findings, or methodologies could shape subsequent research*
535 *directions, model design choices, or scientific applications.*

536 We expect this work to encourage more systematic investigation of objective-induced failure modes
537 in deep representation learning, particularly in settings where representations are reused as inter-
538 mediate states. Our findings suggest that future model design should explicitly consider whether
539 training objectives encode the feasibility constraints of the underlying data manifold, rather than
relying solely on symmetric reconstruction losses.

540 Methodologically, we hope this work promotes the use of hypothesis-driven, controlled interven-
541 tions to disentangle the roles of architecture, capacity, and objective geometry in learned represen-
542 tations. Substantively, the Physics-Aligned Decoding framework illustrates that minimal objective-
543 level modifications can reshape representation semantics without architectural changes, offering a
544 practical template for diagnosing and correcting misaligned representations.

545 While our experiments focus on protein structure, the underlying insights apply to any domain in-
546 volving constrained state spaces, such as robotics, physical simulation, planning, or discrete latent
547 dynamics. We anticipate that future work will extend these ideas to other modalities and develop
548 principled objective designs that ensure learned representations remain semantically valid under
549 reuse.

550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593