

---

# Rethinking Self-Consistency in Protein Generative Models

---

Anonymous Authors<sup>1</sup>

## Abstract

Self-consistency is the standard metric for evaluating structural plausibility in protein generative models: a generated sequence–structure pair is accepted if an independent folding model refolds the sequence to the generated structure within a fixed RMSD threshold. However, self-consistency rejects 44% of native sequence–structure pairs, a gap previously overlooked as a folding model limitation. In this paper, we connect self-consistency to the biophysical ideal of structural plausibility—sufficiently low free energy in that protein’s equilibrium ensemble—which makes two failure modes visible: *flexibility failures*, where a broad ensemble cannot be summarized by a single refold, and *folder failures*. We use ATLAS molecular dynamics trajectories to empirically analyze native structure rejections under this lens, and find that both failure modes are enriched among flexible proteins. We propose ensemble self-consistency using conformational ensemble predictors, improving native protein recall across the flexibility spectrum and rescuing co-design generations previously rejected, without compromising specificity on a Rosetta decoy dataset. Since flexibility is often central to protein function, a metric that penalizes it misdirects generative model development; ensemble self-consistency offers a more faithful framework, and our formalization and diagnostics let it evolve as structure prediction models improve.

## 1. Introduction

Protein structure generative models (Watson et al., 2023; Geffner et al., 2025a; Lin et al., 2024; Zambaldi et al., 2024) are transforming protein design, generating nanomolar binders to viral proteins (Didi et al., 2026b) and small

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Submitted to the 2026 Workshop on Generative and Agentic AI for Biology (ICML 2026). Do not distribute.

molecules (Stark et al., 2025). Quality evaluation of generated protein structures is crucial to guide the development of generative models and to select the structures for costly experimental verification. Yet assessing the plausibility of protein tertiary structure is nontrivial. Standard physical validity metrics (Chen et al., 2010) are largely local; a generative model can easily satisfy the bond geometry and avoid steric clashes while still producing unrealistic folds.

To assess global plausibility, the field has converged on the self-consistency pipeline (Yim et al., 2023). The premise is that folding models have learned the manifold of plausible folds, so an implausible tertiary structure should not be reproduced from its sequence. In this pipeline, given an all-atom structure (which implies a sequence), an independent folding model—commonly ESMFold (Lin et al., 2023)—predicts a structure from the sequence, and the structure is accepted if the two agree within a fixed RMSD threshold, typically 2Å (Yim et al., 2023). Self-consistency is now a headline metric in protein co-design (Geffner et al., 2025a; Rector-Brooks et al., 2026) and is often used as an early filter before experimental validation (Didi et al., 2026a).

Despite its central role, the self-consistency pipeline has received little systematic scrutiny. In our PDB subset of 1,520 native protein structures, 44% fail the standard self-consistency test, even though such structures should be accepted by any reasonable plausibility metric. Similar native failure rates have also been noted in Lu et al. (2025), but the underlying failure modes still remain unclear. In this work, we ask a more nuanced, fundamental question: *when self-consistency rejects a real protein, what exactly has failed?*

Answering this requires first specifying what self-consistency is meant to approximate. We first formalize structural plausibility as membership in the high probability region of a protein’s equilibrium ensemble. Canonical self-consistency, then, can be viewed as a proxy that replaces ensemble membership with proximity to a single predicted structure. This formalization makes two failure modes of native rejection visible. A *flexibility failure* arises when the equilibrium ensemble is broad, so a plausible conformation need not match a single refolded conformation within 2Å. A *folder failure* arises when the folding model’s prediction is itself a poor surrogate for the sequence-conditioned

ensemble. Empirically, using the ATLAS database (Vander Meersche et al., 2024), we decompose native protein failures into these two modes and find that **both** are enriched among flexible proteins. The self-consistency false negatives are not random: they concentrate on flexibility, the core property of protein biology that underlies allostery, catalysis, and binding.

To address this bias, we propose *ensemble self-consistency*, which replaces comparison to a single deterministic re-fold with comparison to a sampled conformation ensemble. This better approximates the ideal acceptance criterion of whether a structure lies within the sequence-conditioned plausible ensemble. Empirically, flexibility failures are nearly eliminated by moderate sample sizes ( $N \geq 16$ ), and replacing deterministic ESMFold self-consistency with SimpleFold ensemble self-consistency significantly improves native-protein recall, while preserving high specificity on a hard Rosetta decoy benchmark (Park et al., 2016). Finally, we re-evaluate published co-design models and show that ensemble self-consistency rescues flexible generations that canonical self-consistency rejects.

**Our contributions are threefold:** (i) A formal account of self-consistency as an approximation to an ideal structural plausibility, and use this to decompose native failures into flexibility and folder modes, (ii) Empirically characterize these two failure modes on native proteins and show that both increase with conformational flexibility, (iii) *Ensemble self-consistency*, a drop-in replacement that nearly eliminates flexibility failures at moderate sample sizes ( $N \geq 16$ ).

## 2. Related works

**Local physical validity metrics.** Protein structure quality has long been assessed with local physical and stereochemical criteria, including bond geometry, Ramachandran outliers, steric clashes (Chen et al., 2010; Adams et al., 2010). These metrics are essential for refining and ranking experimental structure characterizations, where the global fold is established but local details may be uncertain; however, they do not assess whether the global fold is plausible.

**Self-consistency and (co-)designability.** The standard per-sample evaluation metric for protein generative models is self-consistency. Paired with folding-model confidence, it is termed **co-designability**. ESMFold (Lin et al., 2023) is the most commonly used folding model, while DISCO (Rector-Brooks et al., 2026) uses Chai-1 (team et al., 2024) to incorporate ligand input for enzyme design. Crucially, folding models must not require multiple sequence alignments (MSAs) as input, because de novo sequences generally lack meaningful MSAs. This rules out direct use of structure or ensemble predictors that require MSAs for strong performance, such as AlphaFold2/3 (Jumper et al., 2021; Abram-

son et al., 2024) and BioEmu (Lewis et al., 2025). When evaluating the models that only generate protein backbones, sequences are first designed with an inverse folding model such as ProteinMPNN (Dauparas et al., 2022), and the metric is termed **designability**. We discuss the implications of ProteinMPNN redesign in Section 6, but our evaluation scope is the self-consistency component of co-designability in sequence-structure co-design models.

**Distributional metrics for generated protein sets.** Evaluation of protein generative models has largely focused on distributional metrics (Geffner et al., 2025b; Faltings et al., 2026; Lu et al., 2025), which compare generated set to reference set of native proteins. These metrics are motivated by the observation that even highly (co-)designable models can exhibit limited fold-level or secondary-structure diversity (Lin et al., 2024), such as overproducing helical bundles. Distributional metrics are important but orthogonal to our work. Because distributional metrics are often computed after self-consistency filtering, biases in the per-sample filter will propagate to distributional metrics (Sec. 5.3).

**Self-consistency failures in PDB and AFDB** Prior work has noted that self-consistency is imperfect and that optimizing for it can distort generated protein distributions, especially secondary structure and tertiary motifs composition (Lu et al., 2025). We show that the problem is more fundamental; Two sets of proteins can closely match the distribution of secondary structure composition while having completely different flexibility (Sec. 4.2). Distortion in secondary structure is one possible artifact of optimizing for self-consistency, but it does not explain the flexibility bias we focus on in this work. (Reidenbach et al., 2025) report that only 26% of the AlphaFold database (AFDB) is co-designable, and attributes this to a mismatch between native sequences and synthetic structures, motivating consistency-distilled datasets built through ProteinMPNN sequence design and refolding. We show self-consistency failures persist even for native structures paired with their native sequences. Thus, the bias is not a property of the synthetic dataset; it is built into the metric itself. Constructing metric-distilled datasets can only reinforce it.

## 3. Formalizing self-consistency

In this section, we first specify what an ideal *ensemble-level* model and *metric* would mean, and how this relates to the self-consistency test. Building on this perspective, we then formally characterize the self-consistency test’s two failure modes.

**Ideal ensemble-level criterion.** For a sequence  $x$ , let  $\rho_x$  denote the equilibrium density over structural coordinates  $z$  that the sequence can adopt. After marginalizing over unmodeled degrees of freedom, such as solvent, ligands,

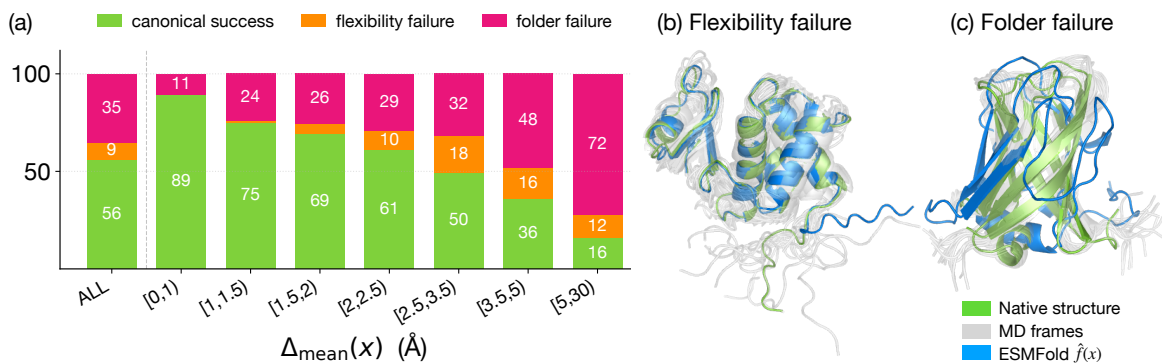


Figure 1. (a) Flexibility-binned native protein recall of ESMFold on ATLAS database. (b) Example of flexibility failure, where ESMFold prediction is close to a plausible conformation produced by MD but is far from the native structure. (c) Example of folder failure, where MD ensemble is tight but ESMFold prediction (optimally aligned to native structure) is incorrect.

and binding partners, this structural ensemble can be written in Boltzmann form:  $\rho_{\mathbf{x}}(\mathbf{z}) \propto \exp(-\mathbf{u}_{\mathbf{x}}(\mathbf{z}))$ , where  $\mathbf{u}_{\mathbf{x}}(\mathbf{z})$  is a (dimensionless) free energy.

This ensemble is the biophysical object that an ideal design model would desire to match. Thus, an ideal design model would not merely output a single structure for a given sequence  $\mathbf{x}$ , but instead output a distribution over structures, which we denote by  $\hat{\rho}_{\mathbf{x}}$ . The ideal ensemble-level criterion is that this designed ensemble be close to the true equilibrium ensemble. For example, using the forward KL divergence, we desire  $\mathcal{D}_{\text{KL}}(\hat{\rho}_{\mathbf{x}}, \rho_{\mathbf{x}}) < \varepsilon$ , or equivalently,

$$\mathbb{E}_{z' \sim \hat{\rho}_{\mathbf{x}}}[-\log \rho_{\mathbf{x}}(z')] = \mathcal{D}_{\text{KL}}(\hat{\rho}_{\mathbf{x}}, \rho_{\mathbf{x}}) + H(\hat{\rho}_{\mathbf{x}}) < \varepsilon + H(\hat{\rho}_{\mathbf{x}}). \quad (1)$$

where  $H(\hat{\rho}_{\mathbf{x}})$  is the entropy of the designed ensemble. Small  $\mathcal{D}_{\text{KL}}$  implies that structures sampled from the designed ensemble have high probability under the true sequence-conditioned equilibrium ensemble.

**Single-structure generations.** For models that generate a single structure, we can view the generated ensemble  $\hat{\rho}_{\mathbf{x}}$  as *coarse-grained* point mass distribution at  $\mathbf{z}$ . Under this limit, (1) reduces to

$$-\log \rho_{\mathbf{x}}(\mathbf{z}) = \mathbb{E}_{z' \sim \hat{\rho}_{\mathbf{x}}}[-\log \rho_{\mathbf{x}}(z')] < \varepsilon + H(\hat{\rho}_{\mathbf{x}}).$$

Therefore, after absorbing the constant and the tolerance, an ideal plausibility criterion for a single generated structure becomes  $-\log \rho_{\mathbf{x}}(\mathbf{z}) < \gamma$  for some  $\gamma > 0$ .

**Canonical self-consistency as an approximation.** In practice, as  $\rho_{\mathbf{x}}$  is unknown, the self-consistency test replaces this equilibrium ensemble with a deterministic folding-model prediction, such as ESMFold, namely  $\hat{f}_{\theta}(\cdot)$ . The standard self-consistency test accepts a generated pair  $(\mathbf{x}, \mathbf{z})$  if  $\mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) < \tau$ , where  $\mathbf{d}$  denotes root mean square deviation (RMSD). In other words, self-consistency uses a geometric comparison as a proxy for the ideal probability-

based criterion:

$$-\log \rho_{\mathbf{x}}(\mathbf{z}) < \gamma \Rightarrow \mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) < \tau, \quad \tau = 2 \text{ \AA}.$$

**Two failure modes.** This approximation suffers from an apparent *false rejection*: a given pair  $(\mathbf{x}, \mathbf{z})$  is acceptable under the point-to-ensemble criterion (blue), but rejected by empirical self-consistency (green) i.e.  $-\log \rho_{\mathbf{x}}(\mathbf{z}) < \gamma$  but  $\mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) \geq \tau$ . We identify the source of such *false rejections* through two scenarios.

- **Flexibility failure:** Both the evaluated structure and the folder prediction are plausible conformations, yet they are separated by more than the RMSD threshold:

$$-\log \rho_{\mathbf{x}}(\mathbf{z}) < \gamma, \quad -\log \rho_{\mathbf{x}}(\hat{f}_{\theta}(\mathbf{x})) < \gamma, \quad \mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) \geq \tau.$$

The rejection is caused by the collapse of a broad or multimodal ensemble to a single structural representative.

- **Folder failure:** The evaluated structure is plausible under the sequence-conditioned ensemble, but the deterministic folder prediction is not:

$$-\log \rho_{\mathbf{x}}(\mathbf{z}) < \gamma, \quad -\log \rho_{\mathbf{x}}(\hat{f}_{\theta}(\mathbf{x})) \geq \gamma, \quad \mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) \geq \tau.$$

In this case, the refolded structure lies in a low-probability region of the sequence-conditioned ensemble. The rejection, therefore, reflects an error in the folding surrogate.

## 4. Two failure modes of self-consistency

In Sec. 3, we defined two failure modes of canonical self-consistency in terms of the unobserved sequence-conditioned equilibrium ensemble  $\rho_{\mathbf{x}}$ . In this section, we examine how these failures appear in native proteins. Native proteins provide a positive-control setting, for which the deposited experimental structure should be accepted by a plausibility metric. Therefore, a rejection of this pair is a false negative of the evaluation pipeline.

Our goal in this section is to turn the formal definitions from Sec. 3 into empirical labels. We use short-time molecular dynamics (MD) frames as a finite empirical approximation to the sufficiently low-energy basin of each sequence, allowing us to distinguish *flexibility failures* from *folder failures*.

#### 4.1. Empirically labeling failure modes in ATLAS

We use the ATLAS database (Vander Meersche et al., 2024), a large collection of all-atom MD trajectories initialized from experimentally characterized protein structures. Each ATLAS protein is simulated for three independent 100ns replicates. For each native sequence  $\mathbf{x}$ , let  $\mathcal{Z}_{\mathbf{x}} = \{\mathbf{z}_{\mathbf{x}}^i\}$  denote the set of MD frames, and let  $\mathbf{z}_{\mathbf{x}}^0$  denote the deposited experimental structure used to initialize the simulations, that is,  $-\log \rho_{\mathbf{x}}(\mathbf{z}_{\mathbf{x}}^0) < \gamma$  presumably holds.

At the *dataset* level, ATLAS is *not* enriched for flexible proteins. Its filtering is based on structure quality, completeness, and fold-level diversity, yielding a broad and relatively non-redundant sample of characterized protein space. Therefore, we treat ATLAS as a representative set of the structurally characterized protein space.

At the *protein* level, the finite MD frames can serve as an *empirical proxy* of the high-probability region of the ensemble, which we instantiate as follows.

$$-\log \rho_{\mathbf{x}}(\hat{f}_{\theta}(\mathbf{x})) < \gamma \Rightarrow \min_{\mathbf{z} \in \mathcal{Z}_{\mathbf{x}}} \mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) < \tau_{\text{MD}}.$$

This empirical approximation is plausible: ATLAS trajectories are short, intended to estimate the flexibility profile of proteins in solution, not to exhaustively sample rare events, unfolding, or large conformational rearrangements. Thus, ATLAS provides a useful basis to study plausible conformations. (See App. F for the caveats.)

In this regard, two failure modes reduce to:

$$\begin{aligned} \text{Canonical success:} \quad & \mathbf{d}(\mathbf{z}_{\mathbf{x}}^0, \hat{f}_{\theta}(\mathbf{x})) < \tau, \\ \text{Flexibility failure:} \quad & \mathbf{d}(\mathbf{z}_{\mathbf{x}}^0, \hat{f}_{\theta}(\mathbf{x})) \geq \tau \quad \text{and} \\ & \min_{\mathbf{z} \in \mathcal{Z}_{\mathbf{x}}} \mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) < \tau_{\text{MD}} \\ \text{Folder failure:} \quad & \mathbf{d}(\mathbf{z}_{\mathbf{x}}^0, \hat{f}_{\theta}(\mathbf{x})) \geq \tau \quad \text{and} \\ & \min_{\mathbf{z} \in \mathcal{Z}_{\mathbf{x}}} \mathbf{d}(\mathbf{z}, \hat{f}_{\theta}(\mathbf{x})) \geq \tau_{\text{MD}}. \end{aligned}$$

where  $\mathbf{d}$  denotes RMSD and  $\tau_{\text{MD}} = \tau = 2\text{\AA}$ .

We quantify the flexibility of each protein by the mean RMSD from the initial structure  $\mathbf{z}_{\mathbf{x}}^0$  to the MD frames, denoted  $\Delta_{\text{mean}}(\mathbf{x})$ . This quantity often exceeds  $2\text{\AA}$  (App. B), already showing the caveat of using a fixed distance cutoff regardless of the varying flexibility of proteins. We then bin ATLAS proteins by  $\Delta_{\text{mean}}(\mathbf{x})$  and report the fraction of each success/failure class within each bin. As shown in Fig. 1a, the failure ratio of canonical self-consistency

Table 1. Distributional distances and structural diversity for ATLAS subsets median-split along each axis. Rg denotes radius of gyration.

Axis	FPD	MMD <sup>2</sup>	Diversity (%)
$\alpha$ -helix %	115.56	0.3206	72.6% / 77.8%
$\beta$ -strand %	105.93	0.2984	75.2% / 78.4%
coil %	64.94	0.1620	74.5% / 81.3%
length	79.93	0.1916	73.0% / 71.1%
Rg	25.19	0.1074	72.1% / 72.7%
$\Delta_{\text{mean}}$	10.30	0.0047	79.4% / 79.4%
random	1.79	0.00	86.4% / 85.2%

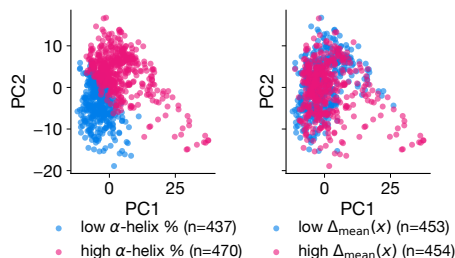


Figure 2. PCA of ESM3 embeddings for median splits by helix content and by  $\Delta_{\text{mean}}$ .

decreases sharply with increasing flexibility. Flexibility failures begin to appear once  $\Delta_{\text{mean}}(\mathbf{x}) > 1.5\text{\AA}$  and increase with protein flexibility, as expected. More importantly, folder failures also increase with flexibility, indicating that flexible proteins are harder for the folding surrogate to predict.

We refer to this binned decomposition as *flexibility-binned native protein recall*, and use it as the diagnostic throughout the rest of the paper.

#### 4.2. Distributional metrics can miss flexibility collapse

The failure-mode analysis shows that self-consistency penalizes proteins along the flexibility axis. A generative model optimized for self-consistency can therefore improve its score by avoiding flexible proteins. A natural question is whether existing distributional metrics would detect this collapse.

We partition the self-consistent ATLAS subset into two equally sized groups using the median  $\Delta_{\text{mean}}$ , and then compute distributional metrics between the two groups to test whether these metrics detect a large difference in flexibility. Because the scale of distributional metrics is often difficult to interpret, we contextualize this split by comparing it to analogous median splits along other axes: secondary structure composition ( $\alpha$ -helix,  $\beta$ -strand, and coil fraction), sequence length, radius of gyration, and a random control split. For each split, we compare the resulting protein sets using standard distributional statistics from the protein de-

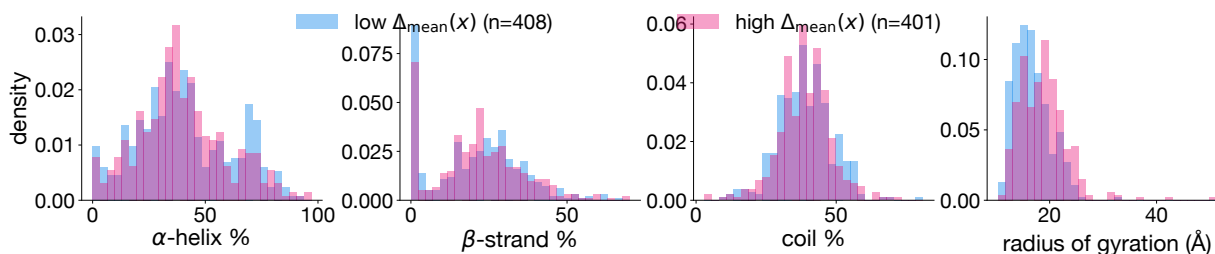


Figure 3. Distribution of secondary structure composition and radius of gyration for the low- and high- $\Delta_{\text{mean}}$  splits.

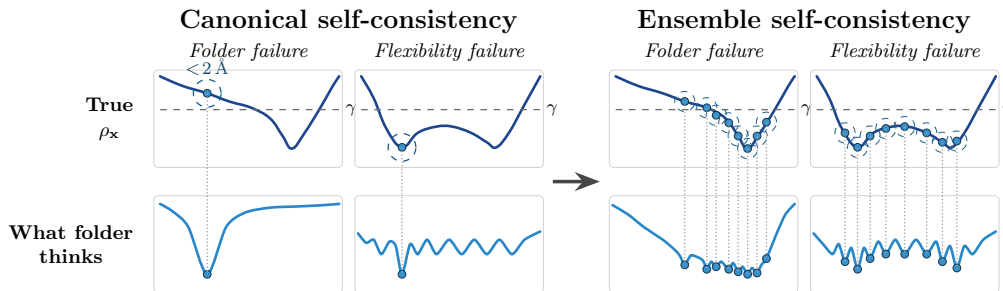


Figure 4. Canonical vs. ensemble self-consistency under the two failure modes. Above denotes the true sequence-conditioned ensemble  $\rho_{\mathbf{x}}$ ; circles of radius  $\varepsilon$  around each folder prediction show the  $2 \text{ \AA}$  acceptance neighborhood. **Canonical (left):** a single deterministic prediction either lands above the energy cutoff  $\gamma$  (folder failure) or covers only part of a broad basin below  $\gamma$  (flexibility failure). **Ensemble (right):** the folder produces a distribution; an acceptance neighborhood is taken around each sample, providing coverage of multiple modes and rescuing both failure types.

sign literature, including structural diversity and Fréchet Protein Distance (FPD) (Lu et al., 2025).

Structural diversity is computed as the percentage of distinct structural clusters under a TM-score cutoff of 0.5, using FoldSeek (van Kempen et al., 2022); we report the diversity of each side of the split. FPD measures the distance between two protein sets as the Fréchet distance between Gaussian fits in ESM3 (Hayes et al., 2025) embedding space, following Lu et al. (2025). Because our sample size is modest and Fréchet distance can be unstable in this regime, we also report maximum mean discrepancy (MMD) (Gretton et al., 2012), a kernel two-sample statistic that compares distributions without assuming Gaussianity.

Table 1 shows that the split by  $\Delta_{\text{mean}}$  produces much smaller distances than the split by other properties. Diversity is also always similar in two sets; Protein sets can exhibit fold-level diversity while suppressing one of the property axes. Figure 2 is consistent with this result: ESM3 embeddings separate proteins split by helix content much more clearly than proteins split by flexibility. Protein generative models already exhibit high FPD to PDB, often in the range of 20-200 (Lu et al., 2025), largely due to distortions in secondary structure composition. In the presence of such dominant distortions, flexibility collapse may be difficult to detect with current distributional metrics. We further compare secondary structure composition and radius of gyration

between the low- and high-flexibility sets in Fig. 3. The two sets are nearly indistinguishable under helix fraction, strand fraction, and coil fraction, and the radius of gyration differs only marginally.

If a generator systematically loses flexible proteins by optimizing for self-consistency, standard distributional metrics may not flag the failure. These results motivate the need for a better self-consistency criterion that does not penalize flexibility, as well as a method for estimating the flexibility of generated proteins.

## 5. Ensemble self-consistency

The flexibility failure analysis in Sec. 4 suggests that the deterministic structure estimate of canonical self-consistency is restrictive by construction: a single point cannot represent a basin whose plausible conformations span more than  $\tau = 2 \text{ \AA}$ . In light of this, we consider an ensemble-aware version of self-consistency (Fig. 4). Concretely, an ensemble conformation predictor  $\theta$  produces a distribution  $\mathbf{z} \sim \hat{\rho}_{\mathbf{x},\theta}$ . Given  $N$  predicted samples  $\hat{\mathbf{z}}_1 \dots, \hat{\mathbf{z}}_N \sim \hat{\rho}_{\mathbf{x},\theta}$ , we define ensemble self-consistency as  $\min_i \mathbf{d}(\mathbf{z}, \hat{\mathbf{z}}_i) < \tau$ , i.e., ensemble self-consistency examines whether  $\mathbf{z}$  is close to at least one sample from the predicted conformation ensemble.

In the limit as  $N \rightarrow \infty$ , ensemble self-consistency mitigates the flexibility failure mode. Intuitively, if  $\hat{\rho}_{\mathbf{x},\theta}$  assigns

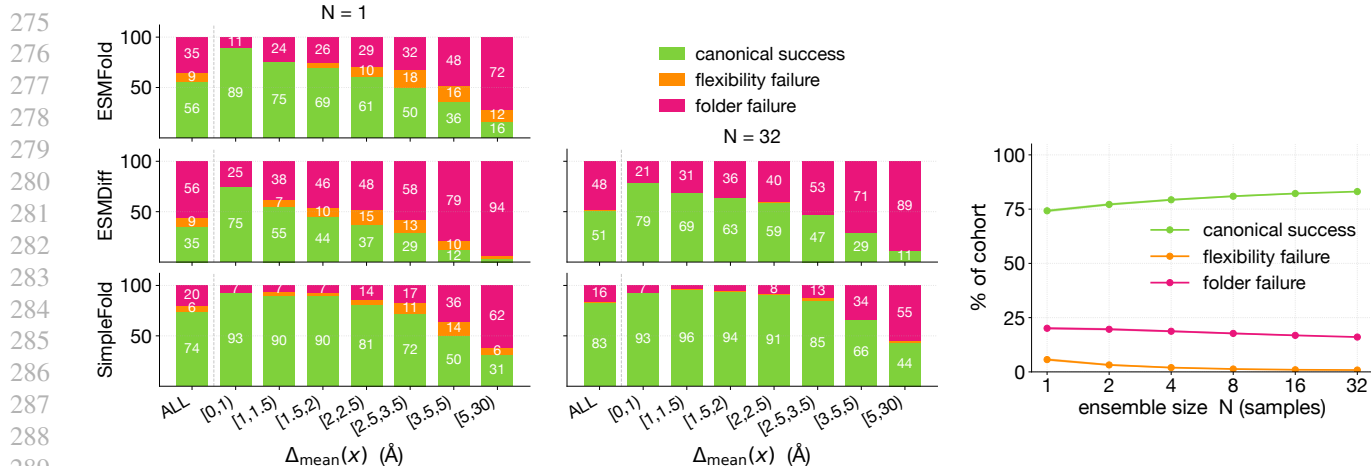


Figure 5. (a) Failure-mode decomposition for canonical ESMFold and ensemble self-consistency using ESMDiff and SimpleFold. (b) Effect of SimpleFold ensemble size on canonical success and failure modes in the full ATLAS set.

sufficient probability mass to regions that are high probable under  $\rho_x$ , then eventually at least one sampled structure will fall in such a region, regardless of how broad or multi-modal the plausible set is.

At finite  $N$ , ensemble self-consistency can still fail even when the predicted samples are individually plausible: the predicted ensemble may sample from the correct conformation, yet the finite sample may not include a conformation close to the particular evaluated structure. Therefore, we extend the diagnosis from Sec. 4 to distinguish *finite-sampling* flexibility failures and folder failures.

We recall the notation from Sec. 4, including  $\mathcal{Z}_x = \{z_x^j\}$  and  $\mathbf{z}_x^0$ . For an ensemble predictor with samples  $\hat{z}_1, \dots, \hat{z}_N$ , we recast two failure mode as:

**Ensemble success:**

$$\min_i d(\mathbf{z}_x^0, \hat{z}_i) < \tau$$

**Finite-sample flexibility failure:**

$$\min_i d(\mathbf{z}_x^0, \hat{z}_i) \geq \tau \text{ and } \forall i, \min_{z' \in \mathcal{Z}_x} d(z', \hat{z}_i) < \tau_{\text{MD}},$$

**Folder failure:**

$$\min_i d(\mathbf{z}_x^0, \hat{z}_i) \geq \tau \text{ and } \exists i, \min_{z' \in \mathcal{Z}_x} d(z', \hat{z}_i) \geq \tau_{\text{MD}}.$$

The above definition follows our intuition. In finite-sample flexibility failure, the predictor is sampling within the MD-sampled plausible basin, but the finite sample set does not cover the particular reference conformation. In folder failure, the predicted distribution assigns mass to conformations that are not supported by the MD ensemble. Note that this decomposition falls into that in Sec. 4.1 when  $N = 1$ .

### 5.1. Gains in flexibility-binned native protein recall

In this section, using our definition of ensemble self-consistency, we examine where current MSA-free ensemble

models stand in native protein recall. We evaluate two representative MSA-free conformation ensemble predictors, ESMDiff (Lu et al., 2024) and SimpleFold (Wang et al., 2025). For each ATLAS sequence  $x$ , we draw  $N = 32$  samples from each ensemble predictor and evaluate ensemble self-consistency against the deposited native structure  $\mathbf{z}_x^0$ . We then apply the failure mode decomposition defined above within flexibility bins of  $\Delta_{\text{mean}}(x)$  used in Sec. 4.

As shown in Figure 5a, increasing the ensemble size to  $N = 32$  with either ESMDiff or SimpleFold nearly eliminates the flexibility failures observed under the  $N = 1$  criterion. This is the failure mode ensemble self-consistency is designed to address, and the decomposition makes clear that an ensemble is what resolves them. In Fig. 5b, SimpleFold’s flexibility failures are nearly eliminated by  $N = 16$ , implying that a moderate number of samples is usually sufficient to cover the plausible basin of a protein.

A substantial fraction of targets still exhibit folder failure. More samples increase the chance that at least one prediction lands near the native structure (see Fig. 4 for illustration), but this only helps when the folding model assigns non-negligible probability mass to the correct plausible basin. When the model’s predicted landscape is biased, additional samples provide limited benefit, so folder failure decays only slowly with  $N$ .

SimpleFold is the strongest **current** choice, outperforming both ESMFold and ESMDiff even at  $N = 1$ . Increasing the SimpleFold ensemble to  $N = 16$  further raises overall recall from 74% to 82%, and recall in the most flexible bin from 50% to 65%. We therefore use SimpleFold  $N = 16$  for the remaining analyses, while emphasizing that the choice should be revisited as ensemble predictors improve—*flexibility-binned native protein recall* provides the diagnostic for that comparison.

Table 2. Self-consistency and structural diversity of co-design models under canonical ( $N=1$ ) versus ensemble ( $N=16$ ) using SimpleFold. **SC- $N$** : percentage of generations accepted by SF- $N$ . **Div- $N$** : structural diversity. LP=La-Proteina.

Model	SC-1	SC-16	Div-1	Div-16
DISCO	41.4%	66.4%	13.8%	18.0%
LP ( $\eta = 0.1$ )	30.6%	47.8%	17.2%	27.6%
LP ( $\eta = 0.2$ )	26.8%	45.0%	8.0%	10.4%
LP ( $\eta = 0.3$ )	22.4%	35.6%	8.4%	9.8%
LP ( $\eta = 0.4$ )	15.6%	25.4%	7.0%	8.8%
RFdiffusion3	14.6%	23.4%	5.4%	6.2%

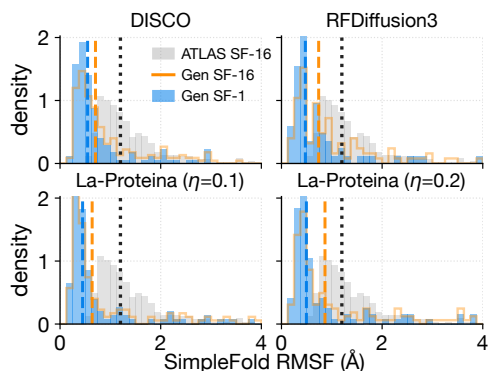


Figure 6. Distribution of SimpleFold RMSF of ATLAS and generated (Gen) self-consistent subset under SF-16 and -1. Dotted lines denote median.

## 5.2. Predicted ensembles as a flexibility proxy

The ideal distributional metric for detecting flexibility collapse would compare the flexibility distribution of generated proteins to that of native proteins. Since even short-time MD simulation is costly, we ask whether SimpleFold ensembles can predict flexibility computed via MD. For each protein, we compute the mean per-residue root mean square fluctuation (RMSF), which is the standard deviation of  $C_\alpha$  position around its mean across an ensemble, from the  $N = 16$  SimpleFold samples, and compare it to the RMSF computed over the ATLAS trajectory. We find a moderate correlation (Pearson  $r = 0.431$ , Spearman  $\rho = 0.550$ , though Fig. 10 shows the absolute scale is compressed), letting us use SimpleFold-derived RMSF as a proxy for detecting protein flexibility collapse in the analyses that follow.

## 5.3. Re-evaluating co-design models using ensemble self-consistency

We re-evaluate three co-design models: DISCO (Rector-Brooks et al., 2026), La-Proteina (Geffner et al., 2025a) at four noise scales  $\eta \in \{0.1, 0.2, 0.3, 0.4\}$ , and RFDiffusion3 (Watson et al., 2023) under ensemble self-consistency with SimpleFold. For each model, we generate 100 samples at each of five lengths  $\{100, 200, 300, 400, 500\}$ , for

a total of 500 generations per model, and report percentage of self-consistent samples and structural diversity under two settings: **SF-1**, canonical self-consistency using a single SimpleFold prediction, and **SF-16**, ensemble self-consistency over  $N = 16$  SimpleFold samples. Table 2 shows all six model variants gain both self-consistency and structural diversity under SF-16 relative to SF-1.

Self-consistent generations under SF-16 span a broader RMSF range than those under SF-1, and their median RMSF is closer to that of the self-consistent ATLAS subset (Fig. 6). However, a gap remains even under SF-16. This is partly because both DISCO and La-Proteina were evaluated using sampling settings chosen to maximize co-designability, which already enriches the sampled distribution for less flexible proteins. Raising La-Proteina’s noise scale does increase median RMSF, but at a steep cost to self-consistency and structural diversity (Table 2). The same limitation we observed for native proteins appears here: flexible proteins are harder to predict, and harder to generate.

Figure 7a shows SimpleFold predicting diverse C-terminal loop conformations for a La-Proteina-generated sequence, with one refold matching the generated loop at RMSD = 1.7Å. In Figure 7b, the 16 SimpleFold refold predictions form two clusters (13 and 3 members), separated by a hinge motion in loop 3 that ultimately reorients helix 4; one of these clusters closely matches the La-Proteina generation (RMSD=0.9Å). This example illustrates that modeling conformational heterogeneity can have effects beyond local fluctuations: flexibility in loops connecting secondary structure elements can drive larger-scale motions that easily exceed a global 2Å RMSD cutoff.

## 5.4. Ensemble self-consistency preserves high specificity in Rosetta decoy dataset

A natural concern with ensemble self-consistency is that taking a minimum of  $N$  predictions could make the criterion too permissive, leading to false positives. We test this using hard negative samples derived from the Rosetta decoy set (Park et al., 2016). These decoys are non-native or near-native alternative conformations generated for proteins with known native structures via local fragment perturbations. Thus, they are much more challenging negatives than randomly mismatched sequence-structure pairs. To avoid labeling plausible alternative conformations as negatives, we intersect this set with ATLAS MD trajectories and retain only decoys that are at least 2Å away from every MD frame of the corresponding protein. This yields 935 decoys across 10 targets (App. E).

On this specificity test, ensemble self-consistency produces only a single additional false positive out of 935: SimpleFold-32 accepts 2 decoys, SimpleFold-1 accepts 1. The min-of-32 RMSD is also only 0.1-0.3Å lower than the

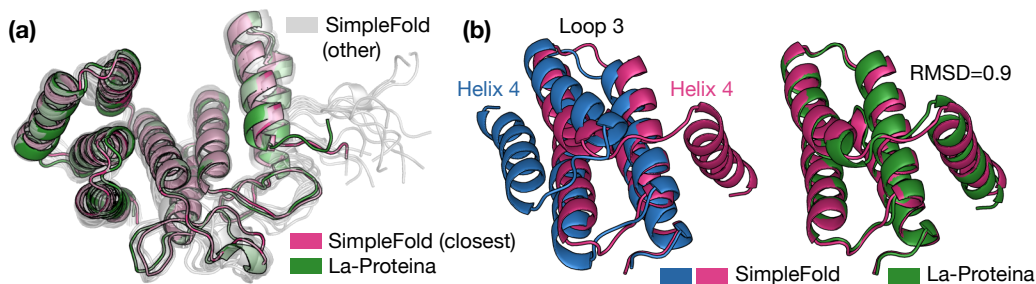


Figure 7. Example La-Proteina generations rescued by ensemble self-consistency. (a) Structural ensemble of 16 SimpleFold refold predictions (gray) for a La-Proteina-generated structure (green), highlighting variability in the C-terminal loop. The closest SimpleFold member (pink) aligns with the La-Proteina at RMSD = 1.7Å. (b) Two SimpleFold refold predictions, one from each of two clusters, distinguished by a hinge motion in loop 3 that reorients helix 4. The smaller cluster representative (pink) matches the La-Proteina conformation (green) at RMSD = 0.9Å.

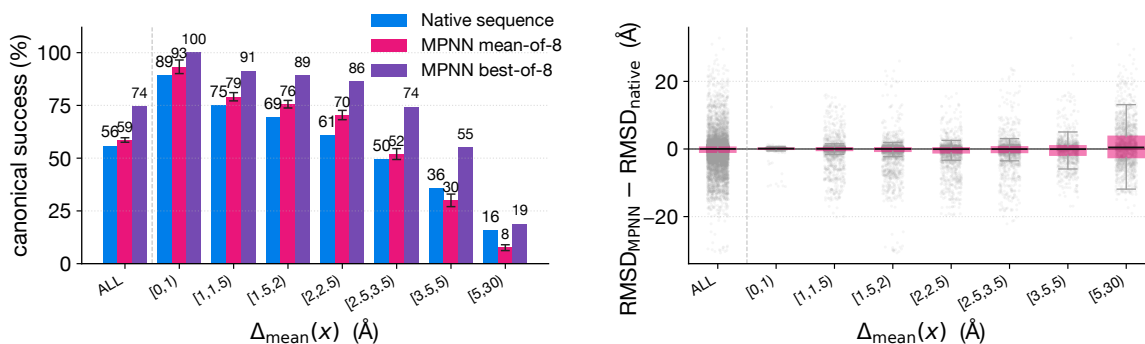


Figure 8. Canonical self-consistency on ATLAS proteins, comparing native sequences against ProteinMPNN redesigns. (a) Acceptance rate per flexibility bin. For MPNN mean-of-8, bars show the mean and error bars the bootstrapped success rate over 100 trials, each sampling a different MPNN redesign per protein. (b) Per-protein RMSD difference  $d(z_0, \hat{f}(\tilde{x}_i)) - d(z_0, \hat{f}(x))$ .

single sample RMSD (Fig. 11), suggesting that the SimpleFold ensembles remain tight for these rigid targets. Thus, ensemble self-consistency improves native protein recall while preserving specificity, even on a hard negative set of near-native decoys.

## 6. ProteinMPNN redesign is not a remedy

The folder failure of Sec. 3 is sequence-specific: for the same plausible structure, some sequences are easy for the folder and some are not. We test this directly using ProteinMPNN redesign. For each ATLAS protein  $x$  with deposited structure  $z_x^0$ , we generate  $N = 8$  MPNN-redesigned sequences  $\tilde{x}_1, \dots, \tilde{x}_8 \sim \text{ProteinMPNN}(\cdot | z_x^0)$  and refold each with ESMFold.

Fig. 8a shows that MPNN redesigns yield only a marginal increase in canonical self-consistency overall, and the improvement is concentrated on rigid proteins; in the most flexible bin canonical success rate falls from 16% to 8%. Fig. 8b shows that even where aggregate acceptance rates are similar, the per-protein RMSD differences are large, with a spread that widens with flexibility.

A natural response is that one could let generative models output  $(x, z)$  pairs but replace  $x$  with an MPNN redesign  $\tilde{x}$  which agrees most with  $z$ , for both evaluation and for downstream design. We argue against this for two reasons. First, sequence selection reflects the model’s choice of sequence for the desired interaction or property, where flexibility may itself be important; it is not merely a way to express a structure. Replacing model-generated sequences with ProteinMPNN redesigns can change the underlying chemical properties, potentially destroying both the flexibility and the functional property the model designed. Second, MPNN inherits the same flexibility-correlated failure. The gain from MPNN redesign diminishes sharply with flexibility: in the most flexible bin, mean acceptance is similar to native sequence. Therefore, redesign provides no rescue precisely where canonical self-consistency fails the most.

## 7. Conclusion

In this work, we disentangle two primary failure modes of the self-consistency pipeline and provide both conceptual and empirical analyses to quantify their respective impacts. Building on this understanding, we extend the standard self-

440 consistency framework to explicitly account for conforma-  
441 tional heterogeneity, resulting in a more robust, ensemble-  
442 based self-consistency pipeline. This approach improves  
443 sensitivity to flexible proteins (i.e., reduces false negatives)  
444 without a loss in precision. Importantly, the ensemble-based  
445 formulation is particularly valuable for downstream design  
446 tasks such as motif scaffolding and binder design, as it  
447 enables the selection of sequences exhibiting diverse conforma-  
448 tions, rather than favoring only rigid structures as in the  
449 standard pipeline. Such conformational diversity is critical  
450 for functional design, facilitating improved signal modu-  
451 lation and access to otherwise occluded or buried binding  
452 pockets.

453 Despite these advances, several limitations remain. First,  
454 our failure-mode decomposition relies on the assumption  
455 that short ATLAS MD trajectories provide sufficient cov-  
456 erage of the plausible conformational basin for each se-  
457 quence. While this approximation is necessary for scaling  
458 the analysis to thousands of proteins, some folder failures  
459 may correspond to plausible conformations not observed in  
460 the short MD. Future work could validate this decomposi-  
461 tion on proteins with longer MD simulations or enhanced  
462 sampling trajectories. Second, our analysis focuses on self-  
463 consistency and its failure modes in isolation. In practice,  
464 downstream design pipelines combine self-consistency with  
465 folding-model confidence metrics to enrich for experimen-  
466 tally viable designs. Understanding how ensemble self-  
467 consistency interacts with confidence measures, and how  
468 the two should be jointly calibrated, remains an important  
469 direction for future work.

470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494

## References

- 495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549
- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, 2024.
- Adams, P. D., Afonine, P. V., Bunkóczi, G., Chen, V. B., Davis, I. W., Echols, N., Headd, J. J., Hung, L.-W., Kapral, G. J., Grosse-Kunstleve, R. W., et al. Phenix: a comprehensive python-based system for macromolecular structure solution. *Biological crystallography*, 66(2):213–221, 2010.
- Chen, V. B., Arendall, W. B., Headd, J. J., Keedy, D. A., Immormino, R. M., Kapral, G. J., Murray, L. W., Richardson, J. S., and Richardson, D. C. Molprobity: all-atom structure validation for macromolecular crystallography. *Biological crystallography*, 66(1):12–21, 2010.
- Dauparas, J., Anishchenko, I., Bennett, N., Bai, H., Ragotte, R. J., Milles, L. F., Wicky, B. I., Courbet, A., de Haas, R. J., Bethel, N., et al. Robust deep learning-based protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- Didi, K., Reidenbach, D., Penner, M., Ravi, S., Case, M., Nichols, M., Swanson, E., Reis, A., Prescott, M., Qian, Y., Qian, D., Yang, J., Li, W., Li, L., Shonai, D., Gay, S., Mallik, B. B., Chim, H. Y., Chen, L., Juantay, M. A., Klein, H., Macintyre, A., Secor, M., Granata, D., Cao, Z., Zhou, G., Geffner, T., Chen, X., Livne, M., Zhang, Z., Zhang, T., Bronstein, M. M., Steinegger, M., Deibler, K., Soderling, S., Khmelinskaia, A., Hollfelder, F., Dallago, C., Kucukbenli, E., Vahdat, A., Ogden, P., and Kreis, K. Latent generative search unlocks de novo design of untapped biomolecular interactions at scale. [https://research.nvidia.com/labs/genair/proteina-complexa/assets/proteina\\_complexa\\_validation.pdf](https://research.nvidia.com/labs/genair/proteina-complexa/assets/proteina_complexa_validation.pdf), 2026a.
- Didi, K., Zhang, Z., Zhou, G., Reidenbach, D., Cao, Z., Cha, S., Geffner, T., Dallago, C., Tang, J., Bronstein, M. M., Steinegger, M., Kucukbenli, E., Vahdat, A., and Kreis, K. Scaling atomistic protein binder design with generative pretraining and test-time compute. In *The Fourteenth International Conference on Learning Representations (ICLR)*, 2026b.
- Faltings, F., Stark, H., Jaakkola, T., and Barzilay, R. Protein fid: Improved evaluation of protein structure generative models. *Bioinformatics*, 42(4), 2026.
- Geffner, T., Didi, K., Cao, Z., Reidenbach, D., Zhang, Z., Dallago, C., Kucukbenli, E., Kreis, K., and Vahdat, A. La-proteina: Atomistic protein generation via partially latent flow matching. *arXiv preprint arXiv:2507.09466*, 2025a.
- Geffner, T., Didi, K., Zhang, Z., Reidenbach, D., Cao, Z., Yim, J., Geiger, M., Dallago, C., Kucukbenli, E., Vahdat, A., et al. Proteina: Scaling flow-based protein structure generative models. *arXiv preprint arXiv:2503.00710*, 2025b.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *The journal of machine learning research*, 13(1):723–773, 2012.
- Hayes, T., Rao, R., Akin, H., Sofroniew, N. J., Oktay, D., Lin, Z., Verkuil, R., Tran, V. Q., Deaton, J., Wiggert, M., et al. Simulating 500 million years of evolution with a language model. *Science*, 387(6736):850–858, 2025.
- Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek, A., Potapenko, A., et al. Highly accurate protein structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- Lewis, S., Hempel, T., Jiménez-Luna, J., Gastegger, M., Xie, Y., Foong, A. Y., Satorras, V. G., Abidin, O., Veeling, B. S., Zaporozhets, I., et al. Scalable emulation of protein equilibrium ensembles with generative deep learning. *Science*, 389(6761):eadv9817, 2025.
- Lin, Y., Lee, M., Zhang, Z., and AlQuraishi, M. Out of many, one: Designing and scaffolding proteins at the scale of the structural universe with genie 2. *arXiv preprint arXiv:2405.15489*, 2024.
- Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W., Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- Lu, J., Chen, X., Lu, S. Z., Shi, C., Guo, H., Bengio, Y., and Tang, J. Structure language models for protein conformation generation. *arXiv preprint arXiv:2410.18403*, 2024.
- Lu, T., Liu, M., Chen, Y., Kim, J., and Huang, P.-S. Assessing generative model coverage of protein structures with shapes. *Cell Systems*, 16(8), 2025.
- Park, H., Bradley, P., Greisen Jr, P., Liu, Y., Mulligan, V. K., Kim, D. E., Baker, D., and DiMaio, F. Simultaneous optimization of biomolecular energy functions on features from small molecules and macromolecules. *Journal of chemical theory and computation*, 12(12):6201–6212, 2016.

- 550 Rector-Brooks, J., Lambert, T., Skreta, M., Roth, D., Long,  
551 Y., Li, Z.-Q., Zhang, X., Cretu, M., Li, F.-Z., Ganap-  
552 athy, T., Jin, E., Bose, A. J., Yang, J., Neklyudov, K.,  
553 Bengio, Y., Tong, A., Arnold, F. H., and Liu, C.-H. Gen-  
554 eral multimodal protein design enables dna-encoding of  
555 chemistry. 2026. URL [https://arxiv.org/abs/  
556 2604.05181](https://arxiv.org/abs/2604.05181).
- 557 Reidenbach, D., Cao, Z., Zhang, Z., Didi, K., Geffner, T.,  
558 Zhou, G., Tang, J., Dallago, C., Vahdat, A., Kucukbenli,  
559 E., et al. Consistent synthetic sequences unlock structural  
560 diversity in fully atomistic de novo protein design. *arXiv*  
561 *preprint arXiv:2512.01976*, 2025.
- 563 Stark, H., Faltings, F., Choi, M., Xie, Y., Hur, E., O’Donnell,  
564 T., Bushuiev, A., Uçar, T., Passaro, S., Mao, W., et al.  
565 Boltzgen: Toward universal binder design. *bioRxiv*, pp.  
566 2025–11, 2025.
- 568 team, C. D., Boitreaud, J., Dent, J., McPartlon, M., Meier,  
569 J., Reis, V., Rogozhonikov, A., and Wu, K. Chai-1: De-  
570 coding the molecular interactions of life. *BioRxiv*, pp.  
571 2024–10, 2024.
- 572 van Kempen, M., Kim, S. S., Tumescheit, C., Mirdita, M.,  
573 Gilchrist, C. L., Söding, J., and Steinegger, M. Foldseek:  
574 fast and accurate protein structure search. *Biorxiv*, pp.  
575 2022–02, 2022.
- 577 Vander Meersche, Y., Cretin, G., Gheeraert, A., Gelly, J.-C.,  
578 and Galochkina, T. Atlas: protein flexibility description  
579 from atomistic molecular dynamics simulations. *Nucleic*  
580 *acids research*, 52(D1):D384–D392, 2024.
- 582 Wang, Y., Lu, J., Jaitly, N., Susskind, J., and Bautista, M. A.  
583 Simplefold: Folding proteins is simpler than you think.  
584 *arXiv preprint arXiv:2509.18480*, 2025.
- 585 Watson, J. L., Juergens, D., Bennett, N. R., Trippe, B. L.,  
586 Yim, J., Eisenach, H. E., Ahern, W., Borst, A. J., Ragotte,  
587 R. J., Milles, L. F., et al. De novo design of protein struc-  
588 ture and function with rfdiffusion. *Nature*, 620(7976):  
589 1089–1100, 2023.
- 591 Yim, J., Campbell, A., Foong, A. Y., Gastegger, M.,  
592 Jiménez-Luna, J., Lewis, S., Satorras, V. G., Veeling,  
593 B. S., Barzilay, R., Jaakkola, T., et al. Fast protein back-  
594 bone generation with se (3) flow matching. *arXiv preprint*  
595 *arXiv:2310.05297*, 2023.
- 597 Zambaldi, V., La, D., Chu, A. E., Patani, H., Danson, A. E.,  
598 Kwan, T. O., Frerix, T., Schneider, R. G., Saxton, D.,  
599 Thillaisundaram, A., et al. De novo design of high-  
600 affinity protein binders with alphaproteo. *arXiv preprint*  
601 *arXiv:2409.08022*, 2024.
- 602  
603  
604

## A. Settings

### A.1. Inference

We used the official repositories, checkpoints, and default inference hyperparameters for ESMFold, SimpleFold, ProteinMPNN, DISCO, La-Proteina, and RFDiffusion3 unless otherwise specified. For ProteinMPNN, we used the  $C\alpha$ -only model with the default noise scale of 0.1. For ESMFold, we used the 3B checkpoint. For SimpleFold, we used the 3B checkpoint and sampled with temperature of 0.8, following the setting used in the SimpleFold paper for multi conformation sampling. For DISCO, we used the full-effort setting, corresponding to 4 recycles and 200 diffusion steps. For La-Proteina, we used the unconditional checkpoint intended for proteins up to 500 residues and sampled with noise scales  $\eta = \eta_x = \eta_z$ , using the default value of 0.1 and additional values of 0.2, 0.3, and 0.4. All other sampling parameters were left at their defaults.

### A.2. Evaluation

We compute RMSD after Kabsch alignment using  $C\alpha$  atoms only. Our official codebase includes all scores computed on the ATLAS database and generated structures, together with scripts to reproduce every figure and table reported in this paper.

### A.3. Compute

All ProteinMPNN, ESMDiff, and ESMFold inference was performed on NVIDIA A100 GPUs (80GB), using under 300 GPU-hours in total. All SimpleFold inference and co-design model generation was performed on AMD MI300A GPUs with the cuBLAS backend, using approximately 1,200 GPU-hours: roughly 1,000 GPU-hours for SimpleFold ( $N = 32$ ) predictions on  $1,520 + 500 \times 6 = 4,520$  samples, and around 200 GPU-hours for co-design model sampling.

## B. Statistics of ATLAS database

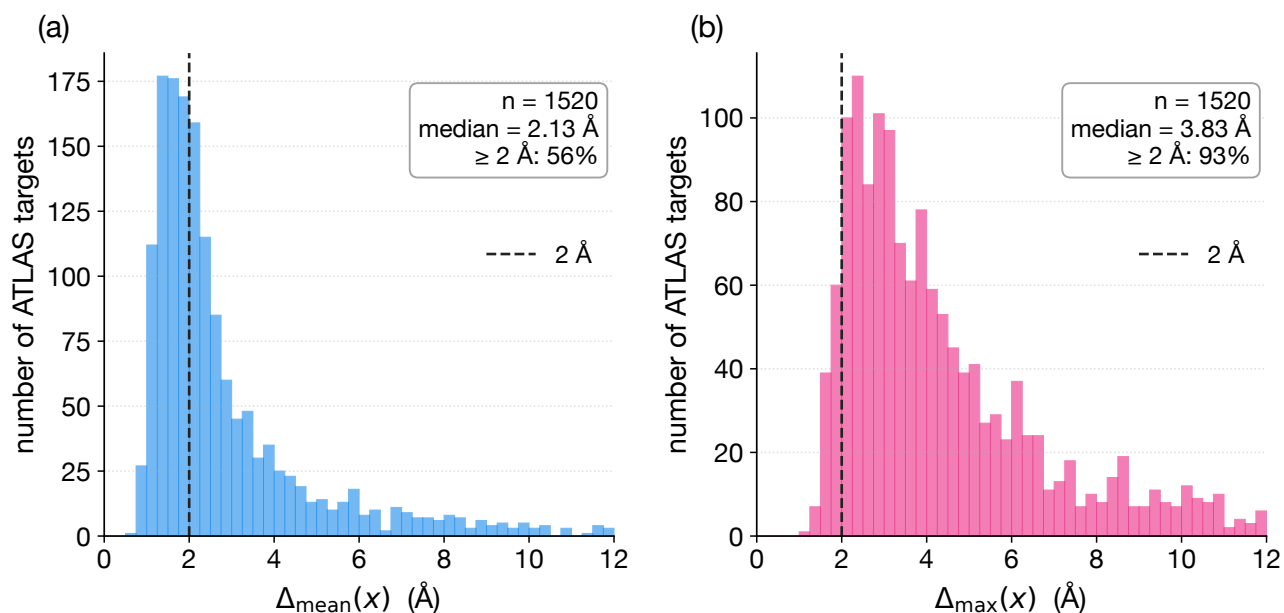


Figure 9. **Distribution of MD displacement in ATLAS.** (a) Mean RMSD to deposited structure across MD frames,  $\Delta_{\text{mean}}$ . (b) Maximum RMSD to deposited structure across MD frames,  $\Delta_{\text{max}}$ .

Figure 9 shows that ATLAS proteins often exhibit conformations beyond the fixed 2 Å threshold used in canonical self-consistency. The median  $\Delta_{\text{mean}}$  is 2.13 Å, and 56% of proteins have mean MD displacement above 2 Å. The maximum displacement is even larger, with median  $\Delta_{\text{max}} = 3.83$  Å and 93% of proteins exceeding 2 Å. This supports the use of proposed ensemble self-consistency: a fixed RMSD cutoff can reject native proteins simply because their plausible conformational ensemble is broad.

### C. SimpleFold flexibility vs. MD flexibility

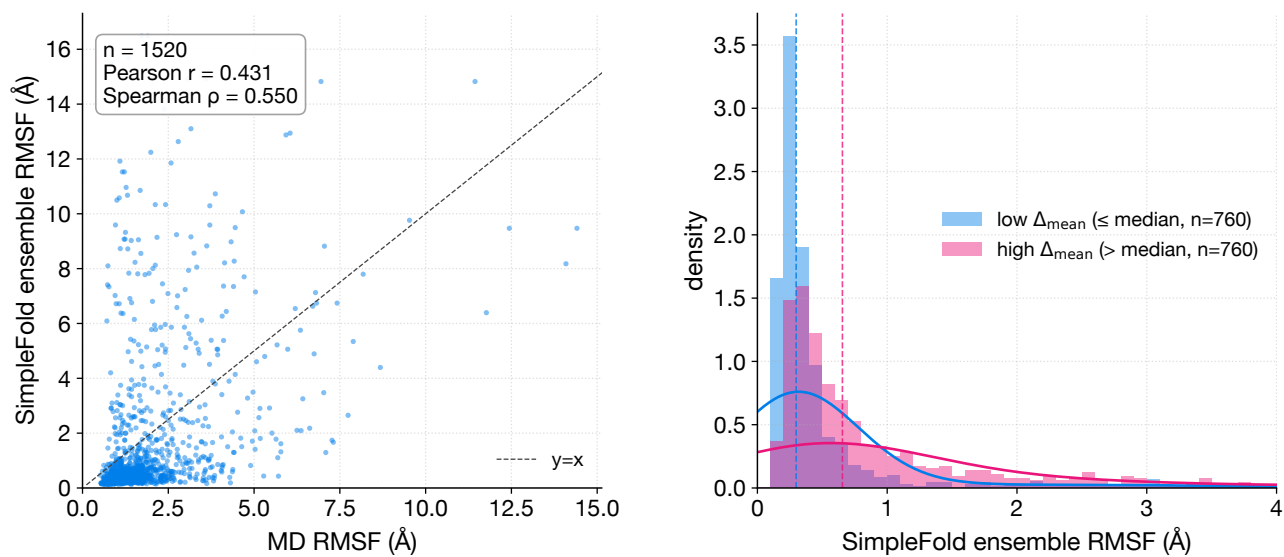


Figure 10. **Comparison between SimpleFold-predicted flexibility and ATLAS MD flexibility.** Per-residue  $C\alpha$  RMSF is computed from SimpleFold ensembles and from ATLAS MD trajectories. SimpleFold captures the relative flexibility of proteins with moderate correlation to MD, but underestimates the magnitude of fluctuations.

Figure 10 compares flexibility estimated from SimpleFold ensembles with flexibility measured from ATLAS MD. SimpleFold RMSF is moderately correlated with MD RMSF, suggesting that ensemble predictors can provide a useful proxy for flexibility. However, SimpleFold systematically underestimates the magnitude of MD fluctuations, so we use it as a diagnostic rather than as a calibrated replacement for MD.

## D. Self-consistency of co-design models by generation length

Table 3–7 report the length-wise breakdown of the results in Table 2. LP denotes La-Proteina.

Table 3. Results for  $L = 100$ .

Model	SC-1	SC-16	Div-1	Div-16
DISCO	66.0%	83.0%	34.0%	37.0%
LP ( $\eta = 0.1$ )	71.0%	83.0%	40.0%	47.0%
LP ( $\eta = 0.2$ )	62.0%	78.0%	33.0%	41.0%
LP ( $\eta = 0.3$ )	60.0%	78.0%	38.0%	45.0%
LP ( $\eta = 0.4$ )	46.0%	65.0%	29.0%	38.0%
RFdiffusion3	44.0%	58.0%	16.0%	18.0%

Table 4. Results for  $L = 200$ .

Model	SC-1	SC-16	Div-1	Div-16
DISCO	56.0%	75.0%	16.0%	18.0%
LP ( $\eta = 0.1$ )	52.0%	76.0%	28.0%	46.0%
LP ( $\eta = 0.2$ )	40.0%	72.0%	15.0%	27.0%
LP ( $\eta = 0.3$ )	29.0%	47.0%	10.0%	15.0%
LP ( $\eta = 0.4$ )	23.0%	38.0%	9.0%	14.0%
RFdiffusion3	21.0%	42.0%	5.0%	6.0%

Table 5. Results for  $L = 300$ .

Model	SC-1	SC-16	Div-1	Div-16
DISCO	38.0%	67.0%	9.0%	14.0%
LP ( $\eta = 0.1$ )	22.0%	45.0%	12.0%	26.0%
LP ( $\eta = 0.2$ )	20.0%	45.0%	7.0%	13.0%
LP ( $\eta = 0.3$ )	13.0%	28.0%	6.0%	8.0%
LP ( $\eta = 0.4$ )	4.0%	15.0%	3.0%	6.0%
RFdiffusion3	5.0%	12.0%	4.0%	4.0%

Table 6. Results for  $L = 400$ .

Model	SC-1	SC-16	Div-1	Div-16
DISCO	31.0%	66.0%	9.0%	20.0%
LP ( $\eta = 0.1$ )	6.0%	26.0%	6.0%	15.0%
LP ( $\eta = 0.2$ )	7.0%	23.0%	5.0%	8.0%
LP ( $\eta = 0.3$ )	6.0%	14.0%	5.0%	8.0%
LP ( $\eta = 0.4$ )	3.0%	4.0%	2.0%	2.0%
RFdiffusion3	3.0%	4.0%	2.0%	2.0%

Table 7. Results for  $L = 500$ .

Model	SC-1	SC-16	Div-1	Div-16
DISCO	16.0%	41.0%	6.0%	7.0%
LP ( $\eta = 0.1$ )	2.0%	9.0%	1.0%	7.0%
LP ( $\eta = 0.2$ )	5.0%	7.0%	3.0%	3.0%
LP ( $\eta = 0.3$ )	4.0%	11.0%	3.0%	4.0%
LP ( $\eta = 0.4$ )	2.0%	5.0%	1.0%	3.0%
RFdiffusion3	0.0%	1.0%	–	1.0%

## E. Rosetta decoy specificity test

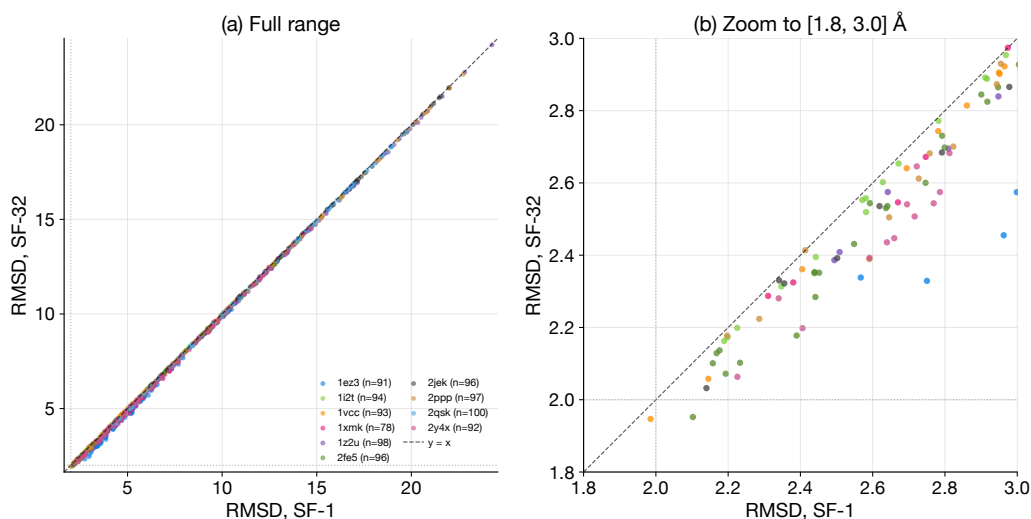


Figure 11. **Ensemble self-consistency preserves specificity on hard negative decoys.** Each point is a decoy structure paired with its native sequence, colored by target. SF-1 denotes the RMSD to the first SimpleFold prediction, corresponding to deterministic self-consistency; SF-32 denotes the minimum RMSD over 32 SimpleFold predictions, corresponding to ensemble self-consistency. The dotted lines mark the 2 Å acceptance threshold.

We identify 10 Rosetta decoy targets that overlap with ATLAS: 1ez3\_B, 1i2t\_A, 1vcc\_A, 1xmk\_A, 1z2u\_A, 2fek\_A, 2ppp\_A, 2qsk\_A, 2y4x\_A. Each target contains 800-1,944 decoys, totaling 13,526 decoys.

Because the decoy set contains low-energy, plausible conformations, we apply an additional filtering procedure to construct negative set. For each decoy, we first compute its  $C\alpha$  RMSD to the experimental reference structure and discard decoys within 2 Å, since these structures are effectively native-like. To obtain a balanced set across decoy difficulty, we bin the remaining decoys for each target into 10 quantiles by RMSD to the native structure and uniformly sample 10 decoys per bin, yielding approximately 100 decoys per target. For each sampled decoy, we then compare it against the ATLAS MD ensemble for the corresponding native sequence. We load the three ATLAS trajectories, restrict all structures to  $C\alpha$  atoms, concatenate the frames, and align the ATLAS  $C\alpha$  sequence to the decoy sequence by exact substring match. We retain a decoy for the specificity test only if it is more than 2 Å from every ATLAS MD frame. This yields 935 decoys across 10 targets.

For each retained negative pair, consisting of a native sequence and a decoy structure, we evaluate SimpleFold self-consistency using 32 SimpleFold samples from the native sequence. We report two scores: SF-1, the RMSD from the decoy to the first SimpleFold sample, and SF-32, the minimum RMSD to any of the 32 samples. A decoy is accepted if its RMSD is below 2 Å.

This experiment therefore measures specificity: whether single sample or ensemble self consistency incorrectly accepts near-native but implausible decoy conformations. As shown in Fig. 11, SF-32 slightly reduces the RMSD for some decoys, as expected when taking the best of 32 samples. However, it rarely moves these hard negatives below the acceptance threshold: SF-32 introduces only one additional false positive out of 935 decoys, accepting 2 decoys compared to 1 accepted by SF-1. Thus, the ensemble self-consistency improves recall for flexible native conformations without sacrificing specificity on this decoy set.

**F. Limitations of the short-time MD-based failure mode analysis**

ATLAS is short trajectories initialized from the experimentally deposited structures; it is not sampling from entire sequence-conditioned equilibrium ensemble  $\rho_x$ . While it justifies us using ATLAS MD frames as plausible (sufficiently low free energy) conformations that protein can adopt in a solution, it also limits the scope of our analysis.

This limitation is most important for conformations separated by large free-energy barriers. Short trajectories are unlikely to capture rare transitions, large-scale domain rearrangements, fold-switching behavior. Consequently, our failure decomposition may overestimate folder failures: a refolded structure that lies outside the ATLAS MD ensemble is not necessarily implausible, but may instead correspond to a conformation that was not sampled within the short trajectory. Finally, ATLAS primarily contains folded proteins with local flexibility and some disordered regions; it does not directly address intrinsically disordered proteins.