

---

# Improving the Efficacy of Test-Time Steering in Masked Diffusion Models with Parallel Tempering

---

Anonymous Authors<sup>1</sup>

## Abstract

Masked Diffusion Models (MDMs) provide expressive generative priors for discrete biological sequences such as proteins and DNA. However, many downstream tasks require steering these models at inference time to optimize arbitrary, external reward functions. Existing test-time steering methods face a fundamental exploration–exploitation trade-off in multimodal reward landscapes: they either collapse into suboptimal or unrealistic modes or require massive sampling budgets to find rare, high-reward states. We address this tension with Parallel Tempering for MDMs (PT-MDM). To adapt parallel tempering to MDM’s generation, we tightly couple the reward temperature to the model’s sequence remasking fraction. Hot replicas apply aggressive remasking for global exploration, while cold replicas use conservative remasking for targeted refinement, and periodic replica exchange ensures rare discoveries propagate across replica chains. This framework enables both global exploration and local exploitation without compromising sequence plausibility. Experiments on inverse protein folding and regulatory DNA design show that PT-MDM consistently outperforms test-time baselines, approaches fine-tuned reward performance on key metrics, and preserves sample fidelity without any training.

## 1. Introduction

Masked Diffusion Models (MDMs) provide expressive generative priors for discrete biological sequences such as proteins, DNA, and other structured biomolecules (Nie et al., 2026; de Groot et al., 2025; Shi et al., 2024). Nevertheless, in many downstream problems, matching the data distribu-

tion is only a starting point to ensure sample fidelity, and generated candidates must also optimize external objectives such as human preferences, biological activity, structural constraints, or task-specific reward models.

The goal can be naturally formalized as a regularized reward optimization problem: seeking a distribution  $q$  that maximizes the expected reward while remaining close to the pretrained MDM prior  $p$  to retain sample fidelity:

$$\max_q \mathbb{E}_{\mathbf{x} \sim q} [R(\mathbf{x})] - \frac{1}{\beta} D_{\text{KL}}(q \parallel p), \quad (1)$$

where  $R(\mathbf{x})$  denotes the reward function, and the temperature  $\frac{1}{\beta}$  controls how strongly  $q$  is constrained to remain close to the prior  $p$ , as quantified by the KL divergence  $D_{\text{KL}}$ .

The solution to Eq. 1 is known as the reward-tilted distribution (Korbak et al., 2022; Gheshlaghi Azar et al., 2024; Go et al., 2023; Rafailov et al., 2023):

$$\tilde{p}(\mathbf{x}) = \frac{1}{Z} p(\mathbf{x}) \exp(\beta R(\mathbf{x})), \quad (2)$$

where  $Z := \sum p(\mathbf{x}) \exp(\beta R(\mathbf{x}))$  is the normalizing constant.

A natural approach to sample from the reward-tilted distribution  $\tilde{p}(\mathbf{x})$  is to fine-tune the model so that its induced distribution approximates the reward-tilted target  $\tilde{p}$ , e.g., DRAKES (Wang et al., 2025). While effective for a fixed objective, fine-tuning couples model parameters tightly to a particular reward and typically requires retraining whenever the reward function changes. This limits their applicability in settings where objectives are frequently updated or changed (Yang et al., 2026). These settings, therefore, demand methods that can adapt immediately to changing rewards without repeated optimization of the base model.

Test-time steering addresses this need by keeping the pretrained MDM fixed and instead modifying the sampling dynamics at inference time to approximately sample from  $\tilde{p}$  (Schiff et al., 2025). Despite this flexibility, effective test-time steering remains challenging because the reward-tilted distribution is often sharp, rugged, and highly multimodal. High-reward solutions may lie in narrow regions assigned low probability under the pretrained prior, while overly aggressive reward optimization can drive samples away from

---

<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.

realistic sequences. Thus, successful steering must balance two competing goals: broad exploration to discover diverse reward modes, and local refinement to concentrate samples within those modes.

Existing particle-based methods expose this tension. Best-of- $K$  explores broadly by drawing independent samples and selecting the highest-reward candidate, but uses reward information only at the final selection stage (Mudgal et al., 2024). As a result, it often requires a large sampling budget to find high-reward regions when they are rare under the prior, and even when such regions are occasionally reached, Best-of- $K$  lacks any mechanism to refine samples further once they are identified (Wu et al., 2023). Methods based on Sequential Monte Carlo (SMC) (Kim et al., 2025; Ou et al., 2026) incorporate reward signals to the transition kernel during generation and can refine promising trajectories more efficiently, yet resampling may reduce diversity and trap particles in early-discovered modes. Consequently, current methods often emphasize either exploration or refinement, but struggle to achieve both simultaneously (He et al., 2026).

### 1.1. Our contribution

In this work, we develop a Parallel Tempering (PT)-based framework for test-time steering of MDMs. Rather than relying on a single inverse temperature  $\beta$  of reward guidance, we maintain a collection of replicas evolving at different temperatures  $\{\beta_k\}_{k=1}^K$  with  $\beta_K = \beta$ . High-temperature (small  $\beta_k$ ) replicas stay close to the pretrained model and promote broad exploration over sequence space, while low-temperature (large  $\beta_k$ ) replicas place stronger emphasis on the reward to refine high-quality candidates. Periodic replica exchange enables information sharing across temperatures, allowing discoveries made in exploratory regimes to propagate to more exploitative ones, thereby mitigating the collapse often observed in single-chain reward-guided sampling.

To support both broad exploration and fine-grained refinement in PT, we augment the native remask-and-denoise sampler of the pretrained MDM with a controllable masking fraction parameter. At each candidate update within a replica, a chosen fraction of low-confidence tokens is remasked and resampled conditioned on the remaining context, yielding a model-aligned transition  $\mathbf{x} \rightarrow \mathbf{x}_{\text{masked}} \rightarrow \mathbf{x}'$ . Larger masking fractions induce broader moves that explore distant regions of sequence space, whereas smaller fractions produce targeted local edits for refinement. A temperature-dependent acceptance step further biases proposals toward reward-improving moves, while replica exchange transfers discoveries across exploratory and exploitative regimes.

As illustrated in Figure 1, PT-MDM achieves a more favorable reward–structure trade-off in inverse protein folding: it attains substantially higher Pred-ddG rewards while keeping

reasonable scRMSD, indicating that the reward improvement does not come from severe structural drift. In contrast, Best-of- $K$  and SMC can preserve plausible structures, but their rewards remain limited and do not reach the same high-reward regime as PT-MDM. This supports our central motivation: effective test-time steering must provide a mechanism for discovering and refining rare, high-reward candidates.

In summary, the main contributions of this work are the introduction of a PT framework for test-time steering in discrete generation. Specifically, we design an MDM-induced remask-and-denoise proposal that does not modify the sampling interface. Our empirical results show that PT consistently improves reward steering on DNA and protein sequence design while preserving prior-aware sample quality, approaching fine-tuned reward performance on key metrics without any parameter updates.

### 1.2. Related work

**Remasking-based discrete diffusion models.** Masked diffusion models introduce a remasking-and-denoising paradigm for discrete generation, enabling iterative refinement through partial corruption and reconstruction. Methods such as LLaDA (Nie et al., 2026) and recent remasking-based scaling approaches (Wang et al., 2026; Schiff et al., 2026; Rector-Brooks et al., 2025) demonstrate that selective remasking improves controllability and sample quality at inference time. While these methods provide a natural mechanism for local refinement, they do not explicitly address global exploration in multimodal reward landscapes.

**Test-time steering for diffusion models.** Recent PT-based methods such as CREPE (He et al., 2026) run parallel diffusion trajectories with replica exchange across timesteps, improving diversity along the sampling path but without explicitly separating exploration from reward-driven exploitation. In contrast, we adapt parallel tempering to directly address this trade-off by maintaining replicas at different reward temperatures ( $\beta$ ). Coupled with temperature-dependent remasking fractions, this design separates broad, prior-driven exploration in hot replicas from focused, reward-guided refinement in cold replicas, reducing collapse into local reward modes.

## 2. Background

**Masked Diffusion Models.** Masked Diffusion Models (MDMs) (Ou et al., 2025) define a generative model over discrete sequences  $\mathbf{x} = (x_1, \dots, x_L) \in \mathcal{V}^L$ , where  $\mathcal{V}$  is a finite vocabulary, i.e., the collection of all possible discrete tokens, and  $L$  is the sequence length. Unlike continuous diffusion models that inject Gaussian noise, MDMs operate through *masking* (Sahoo et al., 2024). A special token

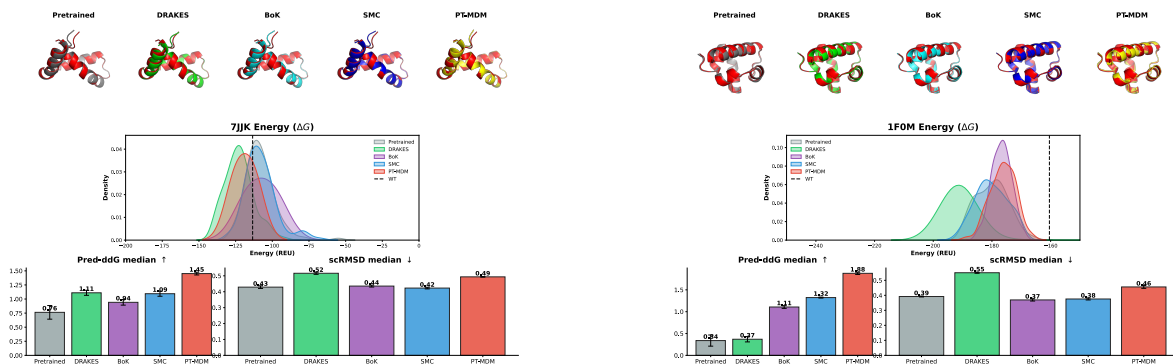


Figure 1. **Protein inverse-folding contrast on representative targets.** The panels compare generated protein designs for 7JJK (left) and 1FOM (right), illustrating how test-time steering changes the reward–structure trade-off relative to the target backbone. PT steers samples toward higher predicted  $\Delta\Delta G$  while retaining structural fidelity through the pretrained denoising model.

[MASK] is introduced, and generation is learned by recovering clean sequences from partially masked inputs (Ou et al., 2025; Zheng et al., 2025).

Formally, let  $\mathbf{x}^{(0)} \sim p_{\text{data}}$  denote a clean sequence. A forward corruption process constructs progressively noisier states  $\{\mathbf{x}^{(t)}\}_{t=1}^T$  via stochastic masking:

$$\mathbf{x}^{(t)} = \mathcal{M}_t(\mathbf{x}^{(0)}), \quad t \in \{1, \dots, T\}, \quad (3)$$

where  $\mathcal{M}_t$  increases the masking level with  $t$ . A standard formulation uses an increasing masking schedule  $\sigma_t \in [0, 1]$  with  $\sigma_1 = 1$ , where  $\sigma_t$  denotes the probability that a token is masked at step  $t$  (Sahoo et al., 2024; Nie et al., 2026). For each position  $i$ , we define

$$x_i^{(t)} = \begin{cases} [\text{MASK}] & \text{if } m_{t,i} = 1, \\ x_i^{(0)} & \text{otherwise,} \end{cases} \quad (4)$$

where

$$m_{t,i} \sim \text{Bernoulli}(\sigma_t). \quad (5)$$

This construction ensures that  $\mathbf{x}^{(t)}$  preserves a subset of the original context while progressively increasing the proportion of masked tokens as  $t$  grows.

The reverse process is parameterized by a denoising model  $p_\theta$  that predicts clean tokens conditioned on the masked sequence:

$$p_\theta(\mathbf{x}^{(0)} | \mathbf{x}^{(t)}) = \prod_{i \in \mathcal{M}(\mathbf{x}^{(t)})} p_\theta(x_i^{(0)} | \mathbf{x}^{(t)}), \quad (6)$$

where  $\mathcal{M}(\mathbf{x}^{(t)})$  denotes the masked positions. Training minimizes the expected cross-entropy over masked tokens (Ou et al., 2025):

$$\mathcal{L}_{\text{MDM}}(\theta) = \mathbb{E}_{t, \mathbf{x}^{(0)}, \mathbf{x}^{(t)}} \left[ - \sum_{i \in \mathcal{M}(\mathbf{x}^{(t)})} \log p_\theta(x_i^{(0)} | \mathbf{x}^{(t)}) \right]. \quad (7)$$

**Remark-and-Denoise Generation.** Sampling from an MDM  $p_\theta$  is started from a fully masked sequence  $\mathbf{x}_{\text{full}} = [\text{MASK}, \dots, \text{MASK}]$ , with the initial state sampled as  $\hat{\mathbf{x}}^{(T)} \sim p_\theta(\cdot | \mathbf{x}_{\text{full}})$ .

Generation then proceeds backward from  $t = T$  to 1 through an iterative remark-and-denoise process (Nie et al., 2026; Wang et al., 2026). Given the current state  $\hat{\mathbf{x}}^{(t)}$ , a remarking operator first produces a partially masked intermediate state:

$$\hat{\mathbf{x}}_{\text{mask}}^{(t)} \sim \rho(\cdot | \hat{\mathbf{x}}^{(t)}, m), \quad (8)$$

where  $\rho$  selects  $\lfloor mL \rfloor$  tokens to reset to [MASK] and  $m \in [0, 1]$  controls the expected proportion of tokens to be remarking. Conditioned on this input, the denoiser samples the next state:

$$\hat{\mathbf{x}}^{(t-1)} \sim p_\theta(\cdot | \hat{\mathbf{x}}_{\text{mask}}^{(t)}). \quad (9)$$

Together, these define the one-step transition kernel

$$g_\theta(\hat{\mathbf{x}}^{(t-1)} | \hat{\mathbf{x}}^{(t)}) = \sum_{\mathbf{x}_{\text{mask}}^{(t)}} p_\theta(\hat{\mathbf{x}}^{(t-1)} | \hat{\mathbf{x}}_{\text{mask}}^{(t)}) \rho(\hat{\mathbf{x}}_{\text{mask}}^{(t)} | \hat{\mathbf{x}}^{(t)}, m). \quad (10)$$

Notice that the generated states  $\{\hat{\mathbf{x}}^{(t)}\}_{t=0}^T$  are clean sequences except in implementations that allow residual mask tokens at intermediate steps. In typical implementations, the final output  $\hat{\mathbf{x}}^{(0)}$  is required to be a clean sequence with no remaining mask tokens.

For the remarking distribution  $\rho(\hat{\mathbf{x}}_{\text{mask}}^{(t)} | \hat{\mathbf{x}}^{(t)}, m)$ , we consider two representative choices:

1. **Uniform:** A subset of tokens is selected uniformly at random by sampling an index set

$$\mathcal{S} \sim \text{Unif}(\{\mathcal{S} \subseteq \mathcal{I} : |\mathcal{S}| = \lfloor mL \rfloor\}), \quad (11)$$

and defining the masked sequence by

$$\hat{x}_{\text{mask},i} = \begin{cases} [\text{MASK}] & \text{if } i \in \mathcal{S}, \\ \hat{x}_i^{(t)} & \text{otherwise.} \end{cases} \quad (12)$$

2. **Low-Confidence (Nie et al., 2026):** Tokens are selected based on the predictive density of  $p_\theta$ . Let

$$c_i = p_\theta(x_i^{(t)} = \hat{x}_i^{(t)} \mid \hat{\mathbf{x}}_{\text{mask}}^{(t-1)}), \quad (13)$$

denote the confidence of token  $i$  given  $\hat{\mathbf{x}}_{\text{mask}}^{(t-1)}$ . We define the index set  $\mathcal{S}$  as the indices of the  $\lfloor mL \rfloor$  smallest confidence values  $c_i$  and mask all positions in  $\mathcal{S}$ .

### 3. Method

We introduce Parallel-Tempering for MDMs (PT-MDM), a test-time steering method for MDMs inspired by the PT algorithm (Earl & Deem, 2005) that jointly enables *global exploration* and *fine-grained refinement* within a unified sampling framework. The method maintains a set of interacting replicas operating at different reward temperatures, allowing simultaneous exploration of diverse high-reward regions and refinement under stronger reward pressure. PT’s exchange mechanism mitigates the exploration–exploitation tension by propagating high-reward discoveries from exploratory replicas to exploitative ones, while injecting diverse states in the opposite direction to prevent mode collapse.

To operationalize this idea in masked diffusion models, we introduce a temperature-dependent remasking schedule, where each replica is assigned a remasking fraction that controls the expected proportion of tokens reset at each update. This remasking fraction acts as a *discrete step size* in sequence space: higher values induce larger structural edits for broad exploration, while lower values restrict updates to localized refinements around high-reward candidates. By coupling this step size to the PT temperature ladder, we obtain a single mechanism that simultaneously governs the scale of local proposals and the strength of reward pressure across replicas. The full procedure is summarized in Algorithm 1.

**Temperature ladder and target distributions.** We maintain  $K$  replicas indexed by increasing inverse temperatures

$$0 \leq \beta_1 < \beta_2 < \dots < \beta_K = \beta, \quad (14)$$

where  $\beta_K = \beta$  is the user-specified target reward temperature of interest. Each replica targets the tempered distribution

$$p_{\beta_k}(\mathbf{x}) = \frac{1}{Z_{\beta_k}} p_\theta(\mathbf{x}) \exp(\beta_k R(\mathbf{x})). \quad (15)$$

In particular, small  $\beta_k$  favors exploration under the prior, while large  $\beta_k$  increasingly emphasizes reward preference.

**Local updates via temperature-scaled remasking.** To construct local proposals within each replica, we adapt the standard remask-and-denoise procedure (Section 2) by introducing temperature-scaled masking fractions  $\{m_k\}_{k=1}^K$  such that  $1 \geq m_1 > m_2 > \dots > m_K > 0$ .

Given the current sequence  $\mathbf{x}^{[k]}$  in replica  $k$ , a candidate state is generated via:

$$\mathbf{x}_{\text{mask}}^{[k]} \sim \rho(\cdot \mid \mathbf{x}^{[k]}, m_k), \quad \tilde{\mathbf{x}}^{[k]} \sim p_\theta(\cdot \mid \mathbf{x}_{\text{mask}}^{[k]}), \quad (16)$$

where  $\rho$  remasks exactly  $\lfloor m_k L \rfloor$  tokens of the length- $L$  sequence. The masking fraction  $m_k$  is tightly coupled to the replica’s inverse temperature  $\beta_k$ : high-temperature replicas (small  $\beta_k$ ) use large  $m_k$  to encourage aggressive, global exploration of the sequence space, while low-temperature replicas (large  $\beta_k$ ) use small  $m_k$  for conservative, fine-grained refinement. This unified mechanism directly aligns the physical size of each proposal edit with the replica’s intended role in the exploration–exploitation trade-off.

**Within-replica acceptance.** We decide whether to accept the proposal  $\tilde{\mathbf{x}}^{[k]}$  as the next state of replica  $k$  using a Metropolis–Hastings update. The exact acceptance probability is

$$A_k^{\text{MH}}(\mathbf{x}^{[k]}, \tilde{\mathbf{x}}^{[k]}) = \min \left\{ 1, \frac{p_\theta(\tilde{\mathbf{x}}^{[k]}) g_k(\mathbf{x}^{[k]} \mid \tilde{\mathbf{x}}^{[k]})}{p_\theta(\mathbf{x}^{[k]}) g_k(\tilde{\mathbf{x}}^{[k]} \mid \mathbf{x}^{[k]})} \exp\left(\beta_k [R(\tilde{\mathbf{x}}^{[k]}) - R(\mathbf{x}^{[k]})]\right) \right\}. \quad (17)$$

However, for masked diffusion proposals, the model ratio  $p_\theta$  is intractable and the proposal density  $g_k(\mathbf{x}^{[k]} \mid \tilde{\mathbf{x}}^{[k]})$  is computationally expensive, as it requires marginalizing over all possible masking configurations. We therefore adopt a simplified proxy:

$$A_k(\mathbf{x}^{[k]}, \tilde{\mathbf{x}}^{[k]}) = \min \left\{ 1, \exp\left(\beta_k [R(\tilde{\mathbf{x}}^{[k]}) - R(\mathbf{x}^{[k]})]\right) \right\}. \quad (18)$$

In particular, we can show that when remasking follows the uniform masking scheme discussed in Section 2, the exact acceptance probability reduces exactly to the simplified proxy:

**Proposition 3.1.** *Assume the remasking distribution  $\rho$  selects a mask of fixed size  $\lfloor m_k L \rfloor$  uniformly at random. Then the exact Metropolis–Hastings acceptance probability  $A_k^{\text{MH}}(\mathbf{x}^{[k]}, \tilde{\mathbf{x}}^{[k]})$  for the remask-and-denoise proposal kernel  $g_k$  simplifies exactly to:*

$$A_k^{\text{MH}}(\mathbf{x}^{[k]}, \tilde{\mathbf{x}}^{[k]}) = A_k(\mathbf{x}^{[k]}, \tilde{\mathbf{x}}^{[k]}) = \min \left\{ 1, \exp\left(\beta_k [R(\tilde{\mathbf{x}}^{[k]}) - R(\mathbf{x}^{[k]})]\right) \right\}. \quad (19)$$

The proof is provided in Appendix A. This result guarantees that under uniform remasking, our simplified acceptance

rule exactly preserves detailed balance with respect to the tempered target distribution. While low-confidence remarking introduces asymmetries that break this guarantee in practice, we find in our experiments that the simplified acceptance rule remains an effective proxy, yielding stable sampling and strong empirical performance across all tasks.

**Replica exchange for global exploration.** To enable communication across replicas, we periodically attempt swaps between adjacent replicas. This mechanism allows high-reward candidates discovered by exploratory (low- $\beta$ ) replicas to propagate to exploitative (high- $\beta$ ) replicas, while those same exploitative replicas receive diverse states that help them escape local modes. For replicas  $k$  and  $k + 1$ , we propose the exchange

$$(\mathbf{x}^{[k]}, \mathbf{x}^{[k+1]}) \rightarrow (\mathbf{x}^{[k+1]}, \mathbf{x}^{[k]}), \quad (20)$$

accepted with probability

$$A_{\text{swap}} = \min \left\{ 1, \frac{p_{\beta_k}(\mathbf{x}^{[k+1]})p_{\beta_{k+1}}(\mathbf{x}^{[k]})}{p_{\beta_k}(\mathbf{x}^{[k]})p_{\beta_{k+1}}(\mathbf{x}^{[k+1]})} \right\} \\ = \min \left\{ 1, \exp \left( (\beta_{k+1} - \beta_k) [R(\mathbf{x}^{[k]}) - R(\mathbf{x}^{[k+1]})] \right) \right\}. \quad (21)$$

Notably, the prior  $p_\theta$  cancels exactly in the ratio, so the swap criterion depends only on reward differences – making it both theoretically exact and cheap to evaluate in practice.

## 4. Experiments

We evaluate PT-MDM across inverse protein folding and regulatory DNA sequence design.

### 4.1. Experimental setup

We compare PT-MDM against test-time steering baselines and, for biological sequence design, reward-fine-tuned DRAKES variants as references for performance and over-optimization. PT-MDM reports both uniform remarking, which matches Proposition 3.1, and low-confidence remarking, which is a practical refinement heuristic. We use  $K$  as a method-specific search-width parameter, so compute-normalized comparisons report denoising-model function evaluations (NFE); Appendix B gives the full baseline, metric, and compute-accounting details.

### 4.2. Protein sequence design

**Dataset, settings, and metrics.** We follow the DRAKES inverse-folding protocol (Wang et al., 2025), using the pre-trained protein MDM as a fixed generative prior and evaluating test-time methods without generator updates. We report predicted stability (Pred-ddG), structural fidelity (scRMSD), validity thresholds, and the combined success rate.

---

### Algorithm 1 Discrete Parallel Tempering for Masked Diffusion Models

---

- 1: **Input:** Pretrained MDM  $p_\theta$ , reward function  $R(\mathbf{x})$ , number of replicas  $K$ , number of iterations  $N$ , swap interval  $\tau_{\text{swap}}$ .
- 2: **Input:** Temperature ladder  $0 \leq \beta_1 < \beta_2 < \dots < \beta_K = \beta$ .
- 3: **Input:** Masking-aggressiveness ladder  $m_1 > m_2 > \dots > m_K$ .
- 4: **Initialize:** Generate initial clean sequences  $\mathbf{x}_0^{[k]} \sim p_\theta$  from fully masked sequences for all  $k = 1, \dots, K$ .
- 5: **for**  $n = 1$  to  $N$  **do**
- 6:   # 1. Within-replica remark-and-denoise updates
- 7:   **for**  $k = 1$  to  $K$  **do**
- 8:     **Remark:** Construct  $\mathbf{x}_{\text{mask}}^{[k]}$  by remarking a fraction  $m_k$  of positions in  $\mathbf{x}_{n-1}^{[k]}$ .
- 9:     **Denoise:** Propose  $\tilde{\mathbf{x}}^{[k]} \sim p_\theta(\cdot | \mathbf{x}_{\text{mask}}^{[k]})$ .
- 10:     Compute acceptance rate

$$A_k = \min \left\{ 1, \exp(\beta_k [R(\tilde{\mathbf{x}}^{[k]}) - R(\mathbf{x}_{n-1}^{[k]})]) \right\}.$$

- 11:     Sample  $u \sim \text{Uniform}(0, 1)$ .
- 12:     **if**  $u < A_k$  **then**
- 13:        $\mathbf{x}_n^{[k]} \leftarrow \tilde{\mathbf{x}}^{[k]}$  {Accept}
- 14:     **else**
- 15:        $\mathbf{x}_n^{[k]} \leftarrow \mathbf{x}_{n-1}^{[k]}$  {Reject}
- 16:     **end if**
- 17:   **end for**
- 18:   # 2. Replica exchange
- 19:   **if**  $n \bmod \tau_{\text{swap}} = 0$  **then**
- 20:     **for**  $k = 1$  to  $K - 1$  **do**
- 21:       Compute acceptance rate

$$A_{\text{swap}} = \min \left\{ 1, \exp \left( (\beta_{k+1} - \beta_k) [R(\mathbf{x}_n^{[k]}) - R(\mathbf{x}_n^{[k+1]})] \right) \right\}.$$

- 22:     Sample  $u \sim \text{Uniform}(0, 1)$ .
  - 23:     **if**  $u < A_{\text{swap}}$  **then**
  - 24:       Swap states:  $\mathbf{x}_n^{[k]} \leftrightarrow \mathbf{x}_n^{[k+1]}$ .
  - 25:     **end if**
  - 26:   **end for**
  - 27: **end if**
  - 28: **end for**
  - 29: **Output:** Final steered sample  $\mathbf{x}_{\text{PT}} = \mathbf{x}_N^{[k]}$  from the coldest replica.
- 

**Results.** Table 1 summarizes results on inverse protein folding. PT-MDM substantially improves over existing test-time steering methods in predicted stability while maintaining high structural validity. Although DRAKES w/o KL obtains high predicted stability, it suffers from severe structural degradation, as reflected by its much larger scRMSD and low fraction of structurally valid sequences. This confirms

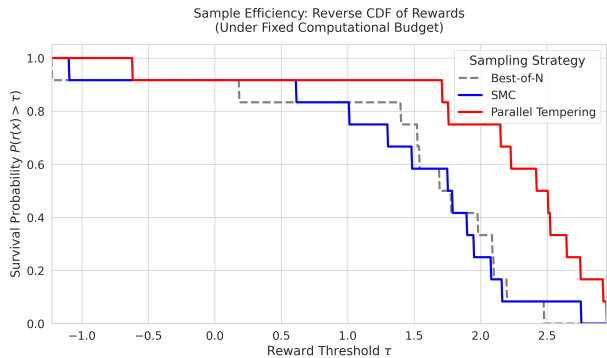


Figure 2. Discovery efficiency on inverse protein folding. The reverse-CDF curves compare the fraction of generated sequences exceeding each Pred-ddG reward threshold under matched denoising-model NFE.

that aggressive reward optimization without sufficient distributional constraints can lead to over-optimization. Low-confidence remasking technically breaks detailed balance, but its gains in reward and structural fidelity make this a useful empirical trade-off. Appendix C visualizes the reward-structure trade-off directly.

**Investigations.** Figure 2 compares the high-reward tail under matched denoising-model NFE and the same number of final evaluated sequences. PT-MDM assigns more mass to high Pred-ddG thresholds, indicating that it converts the same denoising budget into rare high-reward candidates more efficiently. The representative targets in Figure 1 show the same qualitative pattern: PT-MDM raises predicted  $\Delta\Delta G$  while retaining structural fidelity.

### 4.3. DNA sequence design

**Dataset, settings, and metrics.** We follow the DRAGES regulatory enhancer protocol (Wang et al., 2025), using the pretrained DNA MDM as a fixed prior over fixed-length sequences. We report Pred-Activity, ATAC-Acc, 3-mer and JASPAR correlations, approximate log-likelihood, and a diversity diagnostic.

**Results.** Table 2 summarizes the regulatory DNA design results. DRAGES w/o KL achieves high Pred-Activity but shows weaker agreement with natural sequence statistics than KL-regularized DRAGES, suggesting that unconstrained reward optimization can exploit the reward oracle. PT-MDM achieves the highest Pred-Activity while maintaining competitive sequence-statistic correlations and improving approximate log-likelihood relative to the fine-tuned baselines. Its ATAC-Acc remains below DRAGES, suggesting that the Pred-Activity reward used for steering does not fully capture chromatin accessibility; DRAGES, by contrast, may implicitly absorb this aspect during fine-tuning

by adapting the base model itself to the data distribution. Overall, PT-MDM improves reward-oriented performance without additional training while preserving stronger prior-aware sample quality than reward-only fine-tuning.

## 5. Conclusion

In this work, we present Parallel Tempering for Masked Diffusion Models (PT-MDM), a test-time framework designed to steer MDMs toward high-reward regions without the need for expensive, reward-specific fine-tuning. By explicitly coupling the reward temperatures of parallel replicas to their sequence remasking fractions, PT-MDM resolves the inherent exploration-exploitation trade-off in multimodal reward landscapes. Hot replicas naturally execute broad, prior-driven structural exploration, while cold replicas perform targeted, reward-guided refinement. Through periodic replica exchange, our framework seamlessly propagates high-reward discoveries. Empirical evaluations across inverse protein folding and regulatory DNA sequence design demonstrate that PT-MDM consistently outperforms existing test-time steering methods and approaches fine-tuned reward performance on key metrics while preserving sequence plausibility.

Despite its strong empirical performance, we identify a few limitations that motivate future work. Most notably, the current temperature-scaled proposal mechanism selects tokens for remasking largely independently. While this risks disrupting higher-order interactions or coordinated motifs—such as long-range amino-acid contacts in folded proteins—we find that the pretrained MDM’s dense conditioning context naturally helps preserve these global structures during the denoising step. Consequently, sequence validity remains high. Nevertheless, future work exploring token-level, motif-aware, or structure-informed remasking strategies could enable correlated proposals that accelerate reward improvement. In addition, maintaining  $K$  parallel replicas introduces an approximately  $K$ -fold increase in memory usage relative to a single chain, which may become a practical bottleneck for very large models or long sequences. Adaptive temperature ladders and dynamic remasking schedules may help mitigate this cost by allocating computation more efficiently while improving exploration. Finally, extending PT-MDM to hybrid continuous-discrete diffusion models could broaden its applicability across a wider range of scientific design and discovery problems.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning, with applications to generative modeling for biological sequence design. There are many potential societal consequences of our work, none which we feel

## Parallel Tempering for Test-Time Steering in Masked Diffusion Models

Table 1. (Protein sequence design) Model performance on inverse protein folding. Method\* indicates our implementation method.

Group	Method	Pred-ddG $\uparrow$	%(ddG > 0) $\uparrow$	scRMSD $\downarrow$	%(scRMSD < 2) $\uparrow$	Success Rate $\uparrow$
Reference	Pretrained	-0.544 $\pm$ 0.037	36.6 $\pm$ 1.0	0.849 $\pm$ 0.013	90.9 $\pm$ 0.6	34.4 $\pm$ 0.5
Fine-tuned	DRAKES w/o KL	<b>1.108 <math>\pm</math> 0.004</b>	<b>100.0 <math>\pm</math> 0.0</b>	7.307 $\pm$ 0.054	34.1 $\pm$ 0.2	34.1 $\pm$ 0.2
	DRAKES	1.095 $\pm$ 0.026	86.4 $\pm$ 0.2	0.918 $\pm$ 0.006	91.8 $\pm$ 0.5	<b>78.6 <math>\pm</math> 0.7</b>
Test-time guidance	CG	-0.561 $\pm$ 0.045	36.9 $\pm$ 1.1	0.839 $\pm$ 0.012	90.9 $\pm$ 0.6	34.7 $\pm$ 0.9
	SMC	0.659 $\pm$ 0.044	68.5 $\pm$ 3.1	0.841 $\pm$ 0.006	93.8 $\pm$ 0.4	63.6 $\pm$ 4.0
	TDS	0.674 $\pm$ 0.086	68.2 $\pm$ 2.4	0.834 $\pm$ 0.001	94.4 $\pm$ 1.2	62.9 $\pm$ 2.8
	Best-of- $K$ ( $K=10$ )*	0.361 $\pm$ 0.031	59.288 $\pm$ 0.246	0.528 $\pm$ 0.004	93.859 $\pm$ 0.188	56.228 $\pm$ 0.433
Ours	PT-MDM ( $K=10$ , low-confidence)*	0.918 $\pm$ 0.003	74.38 $\pm$ 0.49	0.532 $\pm$ 0.002	<b>94.76 <math>\pm</math> 0.03</b>	71.19 $\pm$ 0.42
	PT-MDM ( $K=10$ , random)*	0.919 $\pm$ 0.015	74.80 $\pm$ 0.55	<b>0.523 <math>\pm</math> 0.003</b>	94.60 $\pm$ 0.72	71.16 $\pm$ 0.61

Table 2. (DNA sequence design) Model performance on regulatory DNA sequence design.

Group	Method	Pred-Activity $\uparrow$	ATAC-Acc (%) $\uparrow$	3-mer Corr $\uparrow$	JASPAR Corr $\uparrow$	App-Log-Lik $\uparrow$
Reference	Pretrained	0.17 $\pm$ 0.04	1.5 $\pm$ 0.2	-0.061 $\pm$ 0.034	0.249 $\pm$ 0.015	-261 $\pm$ 0.6
Fine-tuned	DRAKES w/o KL	6.44 $\pm$ 0.04	82.5 $\pm$ 2.8	0.307 $\pm$ 0.001	0.557 $\pm$ 0.015	-281 $\pm$ 0.6
	DRAKES	5.61 $\pm$ 0.07	<b>92.5 <math>\pm</math> 0.6</b>	<b>0.887 <math>\pm</math> 0.002</b>	<b>0.911 <math>\pm</math> 0.002</b>	-264 $\pm$ 0.6
Test-time guidance	CG	3.30 $\pm$ 0.00	0.0 $\pm$ 0.0	-0.065 $\pm$ 0.001	0.212 $\pm$ 0.035	-266 $\pm$ 0.6
	SMC	4.15 $\pm$ 0.33	39.9 $\pm$ 8.7	0.840 $\pm$ 0.045	0.756 $\pm$ 0.068	-259 $\pm$ 2.5
	TDS	4.64 $\pm$ 0.21	45.3 $\pm$ 16.4	0.848 $\pm$ 0.008	0.846 $\pm$ 0.044	-257 $\pm$ 1.5
	Best-of- $K$ ( $K=10$ )*	1.787 $\pm$ 0.077	7.6 $\pm$ 1.0	0.530 $\pm$ 0.033	0.603 $\pm$ 0.022	-254.399 $\pm$ 0.701
Ours	PT-MDM ( $K=10$ , low-confidence)*	<b>6.854 <math>\pm</math> 0.023</b>	52.6 $\pm$ 0.5	0.870 $\pm$ 0.003	0.818 $\pm$ 0.005	-204.541 $\pm$ 0.278
	PT-MDM ( $K=10$ , random)*	5.876 $\pm$ 0.017	20.4 $\pm$ 3.2	0.848 $\pm$ 0.005	0.787 $\pm$ 0.009	<b>-189.199 <math>\pm</math> 0.978</b>

must be specifically highlighted here. As with any reward-driven generative method applied to biology, practitioners should track biological plausibility and reward-model over-optimization when using these methods to inform downstream experimental decisions.

## References

- de Groot, L., Kuiper, R. J. A., and Bagheri, A. A survey of discrete diffusion for text and genomic sequence generation. In *The 37th Benelux Conference on Artificial Intelligence and the 34th Belgian Dutch Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=ubOzVt7oMw>.
- Earl, D. J. and Deem, M. W. Parallel tempering: Theory, applications, and new perspectives. *Physical Chemistry Chemical Physics*, 7(23):3910–3916, 2005.
- Gheshlaghi Azar, M., Daniel Guo, Z., Piot, B., Munos, R., Rowland, M., Valko, M., and Calandriello, D. A general theoretical paradigm to understand learning from human preferences. In Dasgupta, S., Mandt, S., and Li, Y. (eds.), *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pp. 4447–4455. PMLR, 02–04 May 2024. URL <https://proceedings.mlr.press/v238/gheshlaghi-azar24a.html>.
- Go, D., Korbak, T., Kruszewski, G., Rozen, J., Ryu, N., and Dymetman, M. Aligning language models with preferences through f-divergence minimization. In *Proceedings of the 40th International Conference on Machine Learning*, ICML’23. JMLR.org, 2023.
- He, J., Jeha, P., Potapchik, P., Zhang, L., Hernández-Lobato, J. M., Du, Y., Syed, S., and Vargas, F. CREPE: Controlling diffusion with REPLICA exchange. In *The Fourteenth International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=uOZRWcbiZl>.
- Kim, S., Kim, M., and Park, D. Test-time alignment of diffusion models without reward over-optimization. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=vi3DjUHFVm>.
- Korbak, T., Elshahar, H., Kruszewski, G., and Dymetman, M. On reinforcement learning and distribution matching for fine-tuning language models with no catastrophic forgetting. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=XvI6h-s4un>.
- Mudgal, S., Lee, J., Ganapathy, H., Li, Y., Wang, T., Huang, Y., Chen, Z., Cheng, H.-T., Collins, M., Strohmaier, T., Chen, J., Beutel, A., and Beirami, A. Controlled decoding from language models. In *Forty-first International*

- 385 *Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=bV1cZb7Qa0>.
- 386
- 387
- 388 Nie, S., Zhu, F., You, Z., Zhang, X., Ou, J., Hu, J., ZHOU, J.,  
389 Lin, Y., Wen, J.-R., and Li, C. Large language diffusion  
390 models. In *The Thirty-ninth Annual Conference on Neural*  
391 *Information Processing Systems*, 2026. URL <https://openreview.net/forum?id=KnqiC0znVF>.
- 392
- 393
- 394 Ou, J., Nie, S., Xue, K., Zhu, F., Sun, J., Li, Z., and Li,  
395 C. Your absorbing discrete diffusion secretly models the  
396 conditional distributions of clean data. In *The Thirteenth*  
397 *International Conference on Learning Representations*,  
398 2025. URL <https://openreview.net/forum?id=sMyXP8Tanm>.
- 399
- 400
- 401 Ou, Z., Pani, C., and Li, Y. Inference-time scaling of discrete  
402 diffusion models via importance weighting and optimal  
403 proposal design. In *The Fourteenth International Confer-*  
404 *ence on Learning Representations*, 2026. URL <https://openreview.net/forum?id=7wbrFQvfdH>.
- 405
- 406 Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D.,  
407 Ermon, S., and Finn, C. Direct preference optimization:  
408 Your language model is secretly a reward model. In *Thirty-*  
409 *seventh Conference on Neural Information Processing*  
410 *Systems*, 2023. URL <https://openreview.net/forum?id=HPuSIXJaa9>.
- 411
- 412
- 413 Rector-Brooks, J., Hasan, M., Peng, Z., Liu, C.-H., Mittal,  
414 S., Dziri, N., Bronstein, M. M., Chatterjee, P., Tong, A.,  
415 and Bose, J. Steering masked discrete diffusion models  
416 via discrete denoising posterior prediction. In *The Thir-*  
417 *teenth International Conference on Learning Represent-*  
418 *ations*, 2025. URL <https://openreview.net/forum?id=Ombm8S40zN>.
- 419
- 420
- 421 Sahoo, S. S., Arriola, M., Gokaslan, A., Marroquin, E. M.,  
422 Rush, A. M., Schiff, Y., Chiu, J. T., and Kuleshov, V.  
423 Simple and effective masked diffusion language mod-  
424 els. In *The Thirty-eighth Annual Conference on Neural*  
425 *Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=L4uaAR4ArM>.
- 426
- 427
- 428 Schiff, Y., Sahoo, S. S., Phung, H., Wang, G., Boshar, S.,  
429 Dalla-torre, H., de Almeida, B. P., Rush, A. M., PIER-  
430 ROT, T., and Kuleshov, V. Simple guidance mecha-  
431 nisms for discrete diffusion models. In *The Thirteenth*  
432 *International Conference on Learning Representations*,  
433 2025. URL <https://openreview.net/forum?id=i5MrJ6g5G1>.
- 434
- 435
- 436 Schiff, Y., Belhasin, O., Uziel, R., Wang, G., Arriola, M.,  
437 Turok, G., Elad, M., and Kuleshov, V. Learn from your  
438 mistakes: Self-correcting masked diffusion models, 2026.  
439 URL <https://arxiv.org/abs/2602.11590>.
- Shi, J., Han, K., Wang, Z., Doucet, A., and Titsias, M. Simplified and generalized masked diffusion for discrete data. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=xcqSOfHt4g>.
- Wang, C., Uehara, M., He, Y., Wang, A., Lal, A., Jaakkola, T., Levine, S., Regev, A., Hanchen, and Biancalani, T. Fine-tuning discrete diffusion models via reward optimization with applications to DNA and protein design. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=G328D1xt4W>.
- Wang, G., Schiff, Y., Sahoo, S. S., and Kuleshov, V. Re-masking discrete diffusion models with inference-time scaling. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2026. URL <https://openreview.net/forum?id=IJryQAOy0p>.
- Wu, L., Trippe, B. L., Naesseth, C. A., Cunningham, J. P., and Blei, D. Practical and asymptotically exact conditional sampling in diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=eWKqrlzcrv>.
- Yang, J., Chu, W., Khalil, D., Astudillo, R., Wittmann, B. J., Arnold, F. H., and Yue, Y. Steering generative models with experimental data for protein fitness optimization. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2026. URL <https://openreview.net/forum?id=Ice2BHIumz>.
- Zheng, K., Chen, Y., Mao, H., Liu, M.-Y., Zhu, J., and Zhang, Q. Masked diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=CTC7CmirNr>.

### A. Proof of Proposition 3.1

*Proof.* Let  $\mathbf{x}^{[k]} = (x_1^{[k]}, \dots, x_L^{[k]})$  and  $\tilde{\mathbf{x}}^{[k]} = (\tilde{x}_1^{[k]}, \dots, \tilde{x}_L^{[k]})$  denote the current and proposed discrete sequences of length  $L$ , respectively. For any subset of indices  $A \subseteq \{1, \dots, L\}$ , we let  $\mathbf{x}_A^{[k]}$  denote the sub-sequence of tokens at the positions specified by  $A$ , and let  $\mathbf{x}_{\setminus A}^{[k]}$  denote the unmasked context—the sub-sequence of tokens at the complementary positions  $\{1, \dots, L\} \setminus A$ .

Let  $M \subset \{1, \dots, L\}$  denote the set of indices chosen to be masked. Under the assumption of uniform remasking with a fixed masking fraction  $m_k$ , the size of the mask is strictly defined as  $|M| = \lfloor m_k L \rfloor$ . The probability of selecting any specific mask  $M$  of this size is a combinatorial constant independent of the actual token values in the sequence:

$$\rho(\mathbf{x}_{\setminus M}^{[k]} | \mathbf{x}^{[k]}, m_k) = \rho(\tilde{\mathbf{x}}_{\setminus M}^{[k]} | \tilde{\mathbf{x}}^{[k]}, m_k) = \frac{1}{\binom{L}{\lfloor m_k L \rfloor}} \equiv C, \quad (22)$$

where we slightly abuse notation and use  $\mathbf{x}_{\setminus M}^{[k]}$  to denote the remasked version of  $\mathbf{x}^{[k]}$ .

Next, let  $D = \{i \in \{1, \dots, L\} : x_i^{[k]} \neq \tilde{x}_i^{[k]}\}$  denote the set of indices where the current and proposed sequences differ. For the forward proposal kernel  $g_k(\tilde{\mathbf{x}}^{[k]} | \mathbf{x}^{[k]})$  to have a non-zero probability of generating  $\tilde{\mathbf{x}}^{[k]}$ , the sampled mask  $M$  must cover all differing tokens; otherwise, the unmasked tokens would prevent the transition. Therefore, we require  $D \subseteq M$ . Consequently, for any valid mask  $M$ , the unmasked context remains identical between the two sequences:  $\mathbf{x}_{\setminus M}^{[k]} = \tilde{\mathbf{x}}_{\setminus M}^{[k]}$ .

Using this property, we can write the forward proposal kernel as a sum over all valid masks that contain the differing tokens:

$$g_k(\tilde{\mathbf{x}}^{[k]} | \mathbf{x}^{[k]}) = \sum_{M \supseteq D} p_\theta(\tilde{\mathbf{x}}_M^{[k]} | \mathbf{x}_{\setminus M}^{[k]}) \rho(\mathbf{x}_{\setminus M}^{[k]} | \mathbf{x}^{[k]}, m_k) = C \sum_{M \supseteq D} p_\theta(\tilde{\mathbf{x}}_M^{[k]} | \mathbf{x}_{\setminus M}^{[k]}). \quad (23)$$

By the standard definition of conditional probability, the joint probability of the proposed sequence under the pretrained model factorizes into the masked and unmasked components:  $p_\theta(\tilde{\mathbf{x}}^{[k]}) = p_\theta(\tilde{\mathbf{x}}_M^{[k]} | \tilde{\mathbf{x}}_{\setminus M}^{[k]}) p_\theta(\tilde{\mathbf{x}}_{\setminus M}^{[k]})$ . Rearranging this expression yields  $p_\theta(\tilde{\mathbf{x}}_M^{[k]} | \tilde{\mathbf{x}}_{\setminus M}^{[k]}) = \frac{p_\theta(\tilde{\mathbf{x}}^{[k]})}{p_\theta(\tilde{\mathbf{x}}_{\setminus M}^{[k]})}$ . Substituting this ratio into the forward proposal, and applying the established equality  $\mathbf{x}_{\setminus M}^{[k]} = \tilde{\mathbf{x}}_{\setminus M}^{[k]}$ , we obtain:

$$g_k(\tilde{\mathbf{x}}^{[k]} | \mathbf{x}^{[k]}) = C \sum_{M \supseteq D} \frac{p_\theta(\tilde{\mathbf{x}}^{[k]})}{p_\theta(\mathbf{x}_{\setminus M}^{[k]})} = C p_\theta(\tilde{\mathbf{x}}^{[k]}) \sum_{M \supseteq D} \frac{1}{p_\theta(\mathbf{x}_{\setminus M}^{[k]})} \equiv C p_\theta(\tilde{\mathbf{x}}^{[k]}) C', \quad (24)$$

where  $C' = \sum_{M \supseteq D} [p_\theta(\mathbf{x}_{\setminus M}^{[k]})]^{-1}$  is constant across all valid masks.

Applying the exact same logic and derivation to the reverse proposal kernel yields a symmetric result:

$$g_k(\mathbf{x}^{[k]} | \tilde{\mathbf{x}}^{[k]}) = C \sum_{M \supseteq D} \frac{p_\theta(\mathbf{x}^{[k]})}{p_\theta(\tilde{\mathbf{x}}_{\setminus M}^{[k]})} = C p_\theta(\mathbf{x}^{[k]}) \sum_{M \supseteq D} \frac{1}{p_\theta(\tilde{\mathbf{x}}_{\setminus M}^{[k]})} = C p_\theta(\mathbf{x}^{[k]}) C'. \quad (25)$$

We now substitute these symmetric proposal functions into the model-ratio and proposal-ratio term of the exact Metropolis-Hastings acceptance probability:

$$\frac{p_\theta(\tilde{\mathbf{x}}^{[k]}) g_k(\mathbf{x}^{[k]} | \tilde{\mathbf{x}}^{[k]})}{p_\theta(\mathbf{x}^{[k]}) g_k(\tilde{\mathbf{x}}^{[k]} | \mathbf{x}^{[k]})} = \frac{p_\theta(\tilde{\mathbf{x}}^{[k]}) [C p_\theta(\mathbf{x}^{[k]}) C']}{p_\theta(\mathbf{x}^{[k]}) [C p_\theta(\tilde{\mathbf{x}}^{[k]}) C']} = 1. \quad (26)$$

Because the model priors and the forward/reverse proposal probabilities perfectly cancel one another, the exact acceptance probability simplifies purely to the energy difference defined by the reward-tilted target distribution:

$$A_k^{\text{MH}}(\mathbf{x}^{[k]}, \tilde{\mathbf{x}}^{[k]}) = \min \left\{ 1, \exp \left( \beta_k [R(\tilde{\mathbf{x}}^{[k]}) - R(\mathbf{x}^{[k]})] \right) \right\}. \quad (27)$$

This concludes the proof.  $\square$

## B. Experimental protocol details

**Baselines and remasking variants.** For biological sequence design, we report reward-fine-tuned DRAKES and DRAKES without KL regularization as references for assessing performance and over-optimization, rather than as direct test-time competitors. We compare with classifier guidance (CG), sequential Monte Carlo (SMC), and twisted diffusion sampling (TDS) as reported in DRAKES, and we additionally implement Best-of- $K$  sampling. For PT-MDM, uniform remasking matches the detailed-balance setting of Proposition 3.1, while low-confidence remasking targets uncertain tokens for stronger local refinement.

**Terminology.** We use *candidate sequence* for any completed biological sequence that is ultimately evaluated by the reward and validation metrics. For Best-of- $K$ , the  $K$  objects are independently generated candidate sequences. For SMC, the  $K$  objects are particles that are reweighted and resampled during generation. For PT, the  $K$  objects are persistent replicas, each associated with a single inverse temperature on the temperature ladder. Equivalently, a PT replica is the MCMC trajectory at a fixed inverse temperature; we reserve the term particle for SMC.

**Search width and compute budget.** The symbol  $K$  denotes the internal search width across Best-of- $K$ , SMC, and PT, but its algorithmic role is method-specific. The same value of  $K$  therefore does not necessarily imply the same computational cost. For sample-efficiency comparisons, we normalize by the total number of denoising model forward evaluations, which we report as the number of function evaluations (NFE). For PT, every denoising forward pass performed by any replica at any MCMC proposal step is included in the reported NFE count. Reward oracle evaluations can differ across methods, so they should be reported separately when comparing practical evaluation cost.

**Protein metrics.** Following DRAKES, we report *Pred-ddG*, the predicted stability change of the generated sequence; *scRMSD*, the side-chain RMSD measuring structural fidelity; the fraction of sequences with positive predicted stability, denoted  $\%(ddG > 0)$ ; and the fraction of structurally valid sequences, denoted  $\%(scRMSD < 2)$ . The success rate is the fraction of generated sequences satisfying both  $ddG > 0$  and  $scRMSD < 2$ .

**DNA metrics.** *Pred-Activity* measures predicted enhancer activity under a held-out reward oracle. *ATAC-Acc* evaluates whether generated sequences are classified as chromatin-accessible enhancers by an independent accessibility classifier. *3-mer Corr* is the Pearson correlation between 3-mer frequencies of generated sequences and high-activity natural sequences, and *JASPAR Corr* is the correlation between transcription-factor motif occurrence profiles and those of high-activity natural sequences. *App-Log-Lik* is the approximate sequence log-likelihood under the pretrained MDM, computed as the sequence-level ELBO and averaged over generated sequences. All generated DNA sequences have the same fixed length, making approximate log-likelihood comparisons directly comparable across methods.

**Reverse-CDF analysis.** For sample-efficiency analysis, we compare the empirical reverse CDF of protein *Pred-ddG* rewards,  $\hat{P}(r(x) > \tau)$ , under matched denoising-model NFE and the same number of final evaluated sequences. We estimate the survival probability at threshold  $\tau$  as

$$\hat{P}(r(x) > \tau) = \frac{1}{N_{\text{eval}}} \sum_{i=1}^{N_{\text{eval}}} \mathbf{1}\{r(x_i) > \tau\},$$

where  $N_{\text{eval}}$  is identical across Best-of- $K$ , SMC, and PT-MDM. This prevents high-reward-tail comparisons from being confounded by different numbers of evaluated samples.

**NFE accounting.** Let  $T_{\text{gen}}$  be the number of denoising steps used to generate a completed sequence,  $T_{\text{init}}$  the number of denoising steps used to initialize PT, and  $T_{\text{prop}}$  the number of denoising-model forward evaluations used by one remask-and-denoise proposal. The following table summarizes how the same search-width symbol  $K$  maps to denoising-model NFE for each method. Reward-oracle calls are not included in NFE and should be reported separately.

**Experimental settings** Protein and DNA results are reported over three independent seeds.

**Computational costs** All experiments were run on NVIDIA RTX 6000 Ada GPUs. Per-seed runtime and GPU memory are approximately: Full PT baseline ( $K = 10$ ), 35 minutes and 20 GB; Best-of- $K$ , 30 minutes and 20 GB; uniform

Table 3. Denoising-model NFE accounting for Best-of- $K$ , SMC, and PT. The formulas are written per final reported candidate; multiplying by  $N_{\text{eval}}$  gives the total budget for evaluating  $N_{\text{eval}}$  final candidates.

Method	Role of $K$	Denoising-model NFE per final candidate	Reward-oracle evaluations
Best-of- $K$	$K$ independent completed candidates are sampled, then the best candidate is selected by terminal reward.	$K T_{\text{gen}}$	Typically $K$ terminal reward evaluations per selected candidate.
SMC	$K$ particles are propagated through the denoising process with intermediate reweighting and resampling.	$K T_{\text{gen}}$ for a standard one-forward-pass-per-particle-per-step implementation. Extra corrector or proposal calls should be added explicitly.	Depends on the number of reward-weighting stages; report separately from NFE.
PT	$K$ persistent replicas evolve on an inverse-temperature ladder, with each replica making MCMC proposals and optional replica exchanges.	$T_{\text{init}} + K T_{\text{prop}}$ when one initialized sequence is shared across replicas. Replica swaps add no denoising-model NFE.	Depends on how often rewards are computed for MH acceptance and swap decisions; report separately from NFE.

remasking ( $K = 10$ ), 35 minutes and 20 GB; baseline without guidance, 5 minutes and 5 GB; and SMC, 35 minutes and 20 GB.

### C. Protein reward–structure trade-off

Figure 3 visualizes the inverse-folding trade-off by plotting predicted  $\Delta\Delta G$  against scRMSD. The desired region contains sequences with positive predicted stability and low structural deviation. PT-MDM shifts generated samples toward this high-reward, structurally plausible region, without the severe increase in scRMSD observed under reward-only optimization.

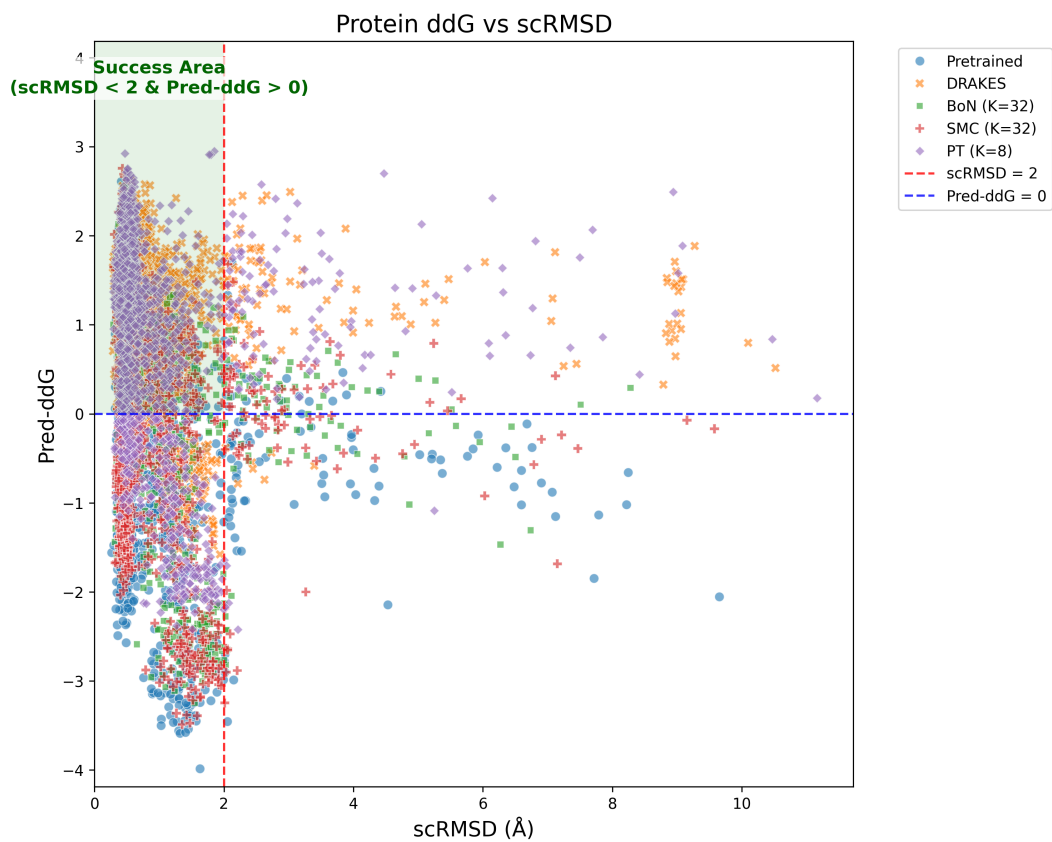


Figure 3. Protein reward–structure trade-off. Predicted  $\Delta\Delta G$  is plotted against scRMSD to show whether methods improve predicted stability while preserving the target fold.