
Active Flow Expansion for Out-of-Distribution Discovery: from Theory to Molecules

Anonymous Authors¹

Abstract

Standard flow and diffusion pre-training matches the distribution of available data (e.g., molecules), which often covers only a small fraction of the valid design space. In generative discovery, however, one aims to sample valid new-to-nature designs, assigned negligible probability under, and thus inaccessible to, standard models fitted to the observed data. To overcome this limitation, we depart from data distribution matching and view a generative model through its *generable set*: the region it covers with non-negligible probability. This allows to introduce a new learning principle for *out-of-distribution flow modeling*: enlarging a model’s generable set to increase coverage of the valid design space. We propose *Active Flow Expansion* (ACTFLOW), a continued pre-training method that employs verifier feedback to expand a pre-trained model over new valid regions by iteratively adapting to synthetic data generated through active exploration in the learned flow representation. Theoretically, we establish to our knowledge first-of-their-kind statistical learning guarantees for out-of-distribution flow modeling, analyzing generable set expansion as a local-to-global reachability process over a learned representation. Empirically, we assess ACTFLOW with suitable out-of-distribution generative modeling metrics across small organic molecules, mid-sized drug-like molecules, therapeutic peptides, and protein sequence design tasks. Results show that ACTFLOW expands valid coverage far beyond the region modeled by the initial pre-trained model, significantly outperforming widely adopted synthetic flow pre-training methods.

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

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

1. Introduction

Large-scale generative modeling has advanced rapidly in recent years, with flow (Lipman et al., 2022; 2024) and diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020) emerging as powerful approaches for generating high-fidelity samples across domains including chemistry (Hoogetboom et al., 2022), biology (Corso et al., 2022), and robotics (Chi et al., 2025). Alongside this progress, over the last years, a growing literature has introduced test-time reward-tilting techniques (e.g., Uehara et al., 2025; Perez Jensen et al., 2026; Uehara et al., 2024; De Santi et al., 2025b). These advances have expanded the use of generative models in scientific discovery, enabling applications such as molecular property optimization (Gutjahr et al., 2025), or functional enzyme design (Rector-Brooks et al., 2026).

A Core Limitation of Standard Flow Pre-training for Discovery. Despite recent progress, generative discovery remains largely constrained by standard pre-training. Flow and diffusion models are learned by matching the distribution of available data (Song et al., 2020; Lipman et al., 2024), which typically covers only a limited portion of the valid design space. This objective is natural when the goal is to reproduce observed examples, e.g., in robotics for manufacturing, where behavioral cloning is often sufficient (Pearce et al., 2023), but it is poorly suited to out-of-distribution discovery, where one typically seeks valid new-to-nature designs with negligible probability under the data distribution, and arbitrarily far from available data acquired via nature’s evolution and prior discoveries. In other words, the goal of out-of-distribution discovery is fundamentally opposed to standard generative modeling as distribution matching. Nonetheless, pre-trained generative models remain, to date, arguably the most promising way to access the complex design spaces in which scientific discovery takes place. This raises the central question of this work:

How can we adapt a pre-trained flow or diffusion model in a task-agnostic way to enable out-of-distribution discovery?

Answering this question would advance the algorithmic-theoretical foundations of generative discovery, with direct implications for high-impact scientific discovery

| Method | GEOM-Drugs Molecules | | Therapeutic Peptides | | Protein Sequences | |
|-------------|--------------------------------------|-------------------------------------|--------------------------------------|-----------------------------------|--------------------------------------|-------------------------------------|
| | Coverage \uparrow | Diversity \uparrow | Coverage \uparrow | Diversity \uparrow | Coverage \uparrow | Diversity \uparrow |
| Pre-trained | 35.89 \pm 2.05 | 255.03 \pm 1.09 | 133.30 \pm 11.1 | 5.70 \pm 0.11 | 66.50 \pm 5.63 | 12.87 \pm 0.63 |
| REC-NF | 44.67 \pm 2.36 | 267.10 \pm 2.88 | 0.00 \pm 0.00 | 0.00 \pm 0.00 | 63.75 \pm 11.45 | 11.85 \pm 0.34 |
| REC-F | 89.33 \pm 11.56 | 284.30 \pm 2.62 | 45.00 \pm 37.22 | 6.03 \pm 1.19 | 49.50 \pm 9.60 | 11.67 \pm 0.40 |
| ACTFLOW | 144.30 \pm 19.28 | 303.10 \pm 5.71 | 197.00 \pm 66.25 | 7.45 \pm 6.17 | 102.75 \pm 18.36 | 42.14 \pm 10.85 |

Table 1. ACTFLOW expands model coverage (i.e., valid clusters) and diversity (i.e., Vendi (Friedman & Dieng, 2022) of valid samples) across domains from de novo molecular design to protein sequence design, significantly outperforming widely adopted baselines.

applications. We take a step toward this goal by moving beyond data distribution matching and *rethinking flow model learning* for out-of-distribution generative modeling. Concretely, we make the following contributions.

Our contributions.

- We introduce the notion of *generable set*: the region of design space a model can access with non-negligible probability. This yields a rigorous mathematical framework for *out-of-distribution flow modeling*, and a new learning principle for continued pre-training: expanding a model generable set to cover a larger portion of the valid design space, which we name *generable set expansion*.
- We propose Active Flow Expansion (ACTFLOW), a continued pre-training method that expands a pre-trained flow model through recurrent adaptation on self-generated synthetic data – as illustrated in Fig. 1 (right). ACTFLOW actively explores in the *learned diffusion representation* at intermediate noise levels, where the geometry is more favorable to global exploration.
- We establish, to our knowledge, first-of-their-kind statistical learning guarantees for out-of-distribution flow modeling. These rely on an energy-based modeling abstraction, and analyze active generable set expansion as a local-to-global reachability process in a learned representation.
- We employ domain-agnostic metrics for out-of-distribution generative modeling and assess ACTFLOW across small organic molecules, mid-sized drug-like molecules, therapeutic peptides, and protein design. ACTFLOW significantly increases valid coverage beyond the initial pre-trained model, vastly outperforming widely used recursive data generation baselines. We further show that the same algorithmic principles and empirical gains extend to discrete (diffusion) models.

2. From Data Distribution Matching to Out-of-Distribution Flow Modeling

We consider a design space $\mathcal{X} \subseteq \mathbb{R}^d$, where each $x \in \mathcal{X}$ is a candidate design, such as a molecule or protein. Let

p_{data} denote the distribution of the available data, assumed to consist of i.i.d. samples that satisfy the domain-specific notion of validity, e.g., valid molecules. A continuous-time flow model is defined by a time-dependent velocity field $u_\theta : \mathcal{X} \times [0, 1] \rightarrow \mathcal{X}$ and the ordinary differential equation

$$\frac{dx_t}{dt} = u_\theta(x_t, t), \quad x_0 \sim p_0, \quad (1)$$

where p_0 is a simple base distribution, typically Gaussian (Lipman et al., 2024). Let p_t^θ denote the marginal density induced at time t by the flow, and let p_1^θ denote its terminal density. In standard flow matching pre-training, one specifies conditional interpolation paths $p_t(\cdot | x_0, x_1)$ between $x_0 \sim p_0$ and $x_1 \sim p_{data}$, with associated target velocity field $u_t^*(\cdot | x_0, x_1)$, and learns u_θ by minimizing the loss:

$$\mathbb{E}[\|u_\theta(x_t, t) - u_t^*(x_t | x_0, x_1)\|_2^2], \quad (2)$$

$$t \sim U[0, 1], \quad x_t \sim p_t(\cdot | x_0, x_1). \quad (3)$$

This vector-field regression objective allows to match the terminal distribution p_1^θ to the available-data distribution p_{data} , as by Equation 4. Beyond flow matching, the same distribution-matching learning principle underlies also standard pre-training of continuous (Song et al., 2020) and discrete (Lou et al., 2023) diffusion models.

Distribution Matching: Standard Flow Model Pre-Training

$$\text{learn } \theta \text{ such that } p_1^\theta \approx p_{data}. \quad (4)$$

2.1. Generable Set Expansion: Flow Modeling Beyond Data Distribution Matching

In this work, we depart from viewing generative models as data distribution approximators and instead cast them as *valid design space approximators*. Let $v : \mathcal{X} \rightarrow \{0, 1\}$ denote whether a design is valid according to the domain one wishes to model, e.g., v might express molecular physiochemical validity, protein sequence foldability, etc. We indicate via Ω^* the entire valid design space:

$$\Omega^* := \{x \in \mathcal{X} : v(x) = 1\}, \quad (5)$$

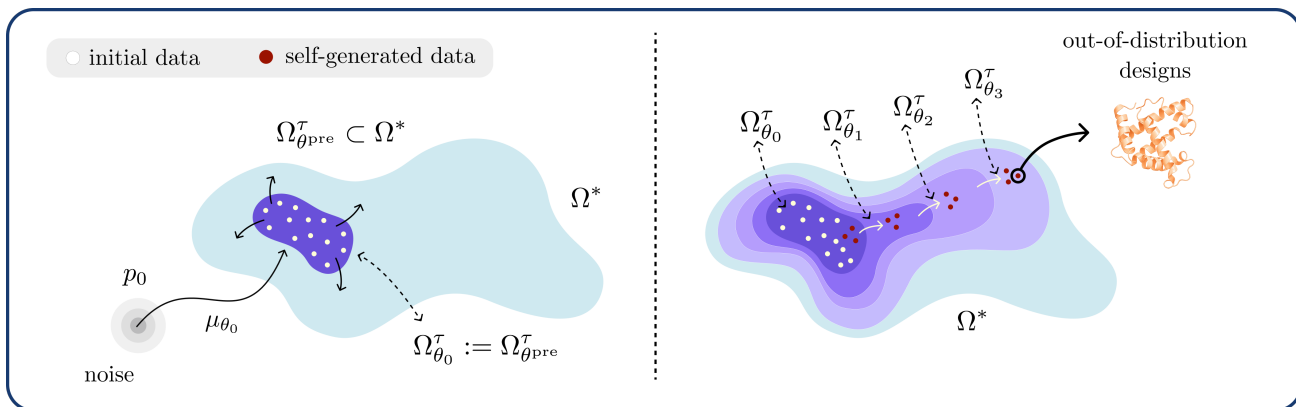


Figure 1. (left) Pre-trained flow model μ_{θ_0} generates with sufficient probability a set $\Omega_{\theta_0}^\tau$ poorly covering the valid design space Ω^* , i.e., $\Omega_{\theta_0}^\tau \subset \Omega^*$. (right) Active Flow Expansion (ACTFLOW) increases model coverage over new valid regions of Ω^* , enabling out-of-distribution flow modeling.

which we wish our generative model to cover well. To capture the data regions sufficiently covered by a pre-trained flow model μ_θ , we introduce the following notion of *generable set*.

Definition 1 (Generable Set). Let μ_θ be a flow model inducing terminal density p_1^θ on \mathcal{X} . For $\tau > 0$, we define its τ -level generable set as

$$\Omega_\theta^\tau := \{x \in \mathcal{X} \mid p_1^\theta(x) \geq \tau\}. \quad (6)$$

For $\tau \rightarrow 0$ the generable set approximates the model support, i.e., $\Omega_\theta^0 := \text{Supp}(p_1^\theta)$. However, the support can include regions assigned only infinitesimal, yet positive probability, and therefore may substantially overstate the model’s effective coverage. Instead, for sufficiently small $\tau > 0$, the generable set Ω_τ^τ captures the region likely to be sampled under a finite budget. With this notion of generable set, we can state the following limitation of standard pre-trained flows for out-of-distribution discovery¹:

Central Limitation of Standard Flow Pre-training

$$\Omega_{\theta^{\text{pre}}}^\tau \subset \Omega^* \quad \text{with} \quad \text{Vol}(\Omega_{\theta^{\text{pre}}}^\tau) \ll \text{Vol}(\Omega^*). \quad (7)$$

Under exact distribution matching, Eq. (7) allows to define the valid out-of-distribution (OOD) region:

$$\text{Valid OOD region: } \bar{\Omega}_{\theta^{\text{pre}}}^\tau := \Omega^* \setminus \Omega_{\theta^{\text{pre}}}^\tau \quad (8)$$

This is the set of valid designs not reliably covered by the pre-trained model $\mu_{\theta^{\text{pre}}}$ and typically not represented in p_{data} , such as new-to-nature designs. To overcome this limitation, we introduce *generable set expansion* (see Eq. 9),

¹For clarity, here we consider a model generating only valid designs. In general $\Omega_{\theta^{\text{pre}}}^\tau \cap (\mathcal{X} \setminus \Omega^*)$ is non-empty (Apx. B)

a new learning principle for continued pre-training: given a standard pre-trained flow model, expand its generable set to cover a larger portion of the valid design space Ω^* .

Generable Set Expansion: Out-of-Distribution Flow Modeling

$$\text{learn } \theta \text{ such that } \Omega_\theta^\tau \approx \Omega^*. \quad (9)$$

2.2. Generable Set Expansion via Iterative Continued Pre-training

Given a pre-trained flow model $\mu_{\theta^{\text{pre}}}$, we approach generable set expansion through continued pre-training on self-generated data. Beyond the data distribution, validity must be obtained through black-box verifier queries, e.g., molecular validity (De Santi et al., 2026; Guo & Schwaller, 2024) or protein foldability (Jumper et al., 2021; Watson et al., 2023). Since such queries can only be made on generable designs, expansion is intrinsically recursive, as shown in Fig. 1:

$$\theta^{\text{pre}} = \theta_0 \rightarrow \theta_1 \rightarrow \dots \rightarrow \theta_T, \quad (10)$$

$$\Omega_{\theta_0}^\tau \subseteq \Omega_{\theta_1}^\tau \subseteq \dots \subseteq \Omega_{\theta_T}^\tau \approx \Omega^*. \quad (11)$$

The central algorithmic question is therefore how to self-generate data that expands the model’s generable set, rather than merely reinforcing regions already assigned high density. Standard sampling is poorly suited to this goal: it preferentially draws from dominant modes, and may therefore neglect valid frontier regions that enable valid expansion. We report an illustrative Gaussian analysis in Apx. C.1 to isolate this failure mode in a minimal setting. The analysis shows that, even when expansion-enabling samples remain within the current generable set, passive self-generation can be exponentially unlikely to sample them when valid directions of expansion are sparse. This highlights the core algorithmic requirement: self-generation must steer

sampling toward generable and valid *frontier* regions, rather than merely reinforce dominant modes. In the next section, we introduce ACTFLOW, which implements this principle and yields provable generable set expansion guarantees.

3. Algorithm: Active Flow Expansion

We now introduce Active Flow Expansion (ACTFLOW) (Alg. 1), a continued pre-training method for *generable set expansion*: it expands a pre-trained flow model’s generable set to cover a larger portion of the valid design space. Concretely, ACTFLOW fine-tunes the model on self-generated data acquired via inference-time active exploration in a learned flow representation at intermediate noising levels.

High-level Algorithm Summary. At round t , ACTFLOW fits a verifier uncertainty estimate $\sigma_t(\cdot)$ from a buffer $\mathcal{D}_t = \{(x_i, y_i)\}_{i=1}^t$, self-generates a new candidate, or a batch, x_{t+1} via Eq. (12), queries the verifier to obtain $y_{t+1} = \tilde{v}(x_{t+1})$, appends (x_{t+1}, y_{t+1}) to \mathcal{D}_t , and updates the flow into $\mu_{\theta_{t+1}}$.

Active Exploration over a Noised Flow Representation. ACTFLOW performs active self-generation in the learned representation \mathcal{Z}_s , given by hidden features of the velocity network μ_θ at noising level $s \in (0, 1)$. At each iteration, Algorithm 1 steers the current model at inference-time toward high verifier-uncertainty regions in \mathcal{Z}_s , while remaining close to the current model density:

(Inference-Time) Active Exploration over Noised Flow Representation \mathcal{Z}_s

$$x_{t+1} \sim \tilde{p}_t \in \arg \max_q \mathbb{E}_{x \sim q} [\sigma_t(\phi_s^t(x))] - \beta \text{KL}(q \| p_1^{\theta_t}) \quad (12)$$

Here, the first term favors informative queries about the verifier v , and thus about the valid design space Ω^* . The KL term regularizes search toward the current generative prior, biasing exploration toward likely valid regions. The parameter β controls the exploration–prior trade-off: as $\beta \rightarrow \infty$, Eq. (12) recovers standard sampling, $x_{t+1} \sim p_1^{\theta_t}$; as $\beta \rightarrow 0$, it becomes pure uncertainty maximization, which may target remote regions where the model prior might no longer provide a useful validity bias.

To instantiate Eq. (12), we model verifier uncertainty directly in \mathcal{Z}_s . Let $\phi_s^t : \mathcal{X} \rightarrow \mathcal{Z}_s$ denote the representation extracted from the velocity network at noising level $s \in (0, 1)$. We then view the verifier labels as noisy observations $y_t = \tilde{v}(x_t)$ of an unknown validity function over the representation space \mathcal{Z}_s . To represent the uncertainty about this function, we use a linear kernel over the learned representation space: $k_{\phi_s^t}(x, x') := \langle \phi_s^t(x), \phi_s^t(x') \rangle$. Then, for a set of queried designs x_1, \dots, x_t , let $X_t := (x_1, \dots, x_t)$, let $K_{t, \phi_s^t} \in \mathbb{R}^{t \times t}$ be

Algorithm 1 Active Flow Expansion (ACTFLOW)

Require: Flow model μ_{θ_0} , black-box verifier \tilde{v} , representation time-step $s \in (0, 1)$, iterations T

- 1: $\mathcal{D}_0 \leftarrow \emptyset$
- 2: **for** $t = 0, 1, \dots, T - 1$ **do**
- 3: Update surrogate uncertainty σ_t from \mathcal{D}_t
- 4: Self-generate x_{t+1} according to:

$$p_t \in \arg \max_q \mathbb{E}_{x \sim q} [\sigma_t(\phi_s^t(x))] - \beta \text{KL}(q \| p_1^{\theta_t})$$
- 5: Query verifier: $y_{t+1} \leftarrow \tilde{v}(x_{t+1})$
- 6: $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(x_{t+1}, y_{t+1})\}$
- 7: $\theta_{t+1} \leftarrow \text{UPDATEFLOW}(\theta_t, \mathcal{D}_{t+1})$
- 8: **end for**
- 9: **return** μ_{θ_T}

the kernel matrix with entries $(K_t)_{ij} = k_{\phi_s^t}(x_i, x_j)$, and write $k_{\phi_s^t}(x, X_t) := (k_{\phi_s^t}(x, x_1), \dots, k_{\phi_s^t}(x, x_t))$. We then express our uncertainty about the unknown verifier via the following closed-form expression:

$$\sigma_t^2(x) = k_{\phi_s^t}(x, x) - k_{\phi_s^t}(x, X_t) (K_t, \phi_s^t + \lambda I)^{-1} k_{\phi_s^t}(X_t, x) \quad (13)$$

$$\sigma_t(x) := \sqrt{\sigma_t^2(x)}.$$

This corresponds to the posterior variance in Bayesian linear regression under Gaussian observations.

Model update. The uncertainty estimator is fit on the pairs in \mathcal{D}_t . The flow model is then updated by replay-based continued pre-training. Let $\mathcal{D}_t^+ = \{x : (x, 1) \in \mathcal{D}_t\}$ and $\mathcal{D}_t^- = \{x : (x, 0) \in \mathcal{D}_t\}$ denote accepted and rejected samples, and let $U_t^\pm \subseteq \mathcal{D}_t^\pm$ be minibatches and $\hat{L}_t^\pm(\theta)$ their respective standard flow-matching loss terms. We use the signed update $g_t = \nabla \hat{L}_t^+(\theta_t) - \alpha_t \nabla \hat{L}_t^-(\theta_t)$, where rejected samples are seen as an unlearning signal (Alberti et al., 2025); see Apx. D.1. In practice, our results often hold with $\alpha_t = 0$, i.e., standard flow matching (Lipman et al., 2024) on verifier-accepted samples.

Across iterations, ACTFLOW reallocates model mass from dominant pre-trained modes toward newly discovered valid regions. It remains to show that this process truly expands the model’s generable set, rather than only redistributing density within it. The next section answers this question affirmatively by establishing statistical guarantees for out-of-distribution generative modeling via ACTFLOW.

4. Statistical Learning Guarantees for Out-of-Distribution Flow Modeling

We now theoretically analyze Active Flow Expansion (ACTFLOW), thereby providing a first guarantee for out-of-distribution generative modeling in terms of generable set

expansion rather than distribution matching of p_{data} . Formal theorem and detailed derivations are reported in Apx. F.

As a first step, since ACTFLOW effectively performs active exploration in the learned representation \mathcal{Z}_s , we introduce the following notion of *generable representation set* $\Omega_{\pi, \phi}^\tau$ for a fixed map $\phi := \phi_s$.

Definition 2 (Generable Representation Set). *Let $\phi : \mathcal{X} \rightarrow \mathcal{Z}$ be fixed, and let π be a diffusion model inducing design-space density p_1^π on \mathcal{X} . Let $p_1^{\pi, \phi}$ denote the induced density on \mathcal{Z} obtained by pushing forward p_1^π through ϕ . For $\tau > 0$, we define its τ -level generable representation set as*

$$\Omega_{\pi, \phi}^\tau := \{z \in \mathcal{Z} \mid p_1^{\pi, \phi}(z) \geq \tau\}. \quad (14)$$

Formal feedback model: logistic bandit feedback We now formalize the verifier as a probabilistic model over $Z = \phi(X)$. At each round t , ACTFLOW generates a sample $z_t := \phi(x_t)$ and observes a binary label $y_t \in \{0, 1\}$, generated according to a logistic model:

$$\Pr[y_t = 1 \mid z_t] = s(g(z_t)) \quad (15)$$

where $g : Z \rightarrow \mathbb{R}$ is an unknown latent *validity score*, and $s : \mathbb{R} \rightarrow [0, 1]$ is the sigmoid function. We assume there exists a threshold $h \in [0, 1]$, such that we can formally define the valid design space as:

$$\Omega_* := \{z \in Z : s(g(z)) \geq h\}. \quad (16)$$

For the sake of analysis, we assume g lies in an RKHS \mathcal{H}_k with $\|g\|_k \leq B$, $\int_{\mathcal{Z}} \exp(g(z)) \leq \bar{Z}$, and is L_g -Lipschitz in a metric d over \mathcal{Z} . This model allows us to construct *any-time confidence sequences* for logistic prediction (Pásztor et al., 2024). Next, we model ACTFLOW’s uncertainty-tilted sampling procedure.

Uncertainty-tilted sampling oracle modeling We model the inference-time uncertainty sampling in Eq. (12) as approximately maximizing verifier uncertainty within the current generable set Ω_t^τ for some $\tau > 0$. Concretely, at round t it returns x_t such that $z_t = \phi(x_t) \in \Omega_t^\tau$ satisfies:

$$\sigma_t(z_t) \geq \frac{1}{\alpha} \max_{z \in \Omega_t^\tau} \sigma_t(z), \quad \alpha \geq 1. \quad (17)$$

We depart from standard bandit analyses, which assume global search over Ω^* by instead restricting the (generative) sampler to approximate uncertainty maximization over the generable set. We suppose that the τ -level generable set of the pre-trained model contains a non-empty set S_0 of valid designs. Further regions can be reached only after intermediate expansion steps. We formalize a local expansion step

via the following *one-step reachability operator over the learned representation*:

$$R_\epsilon(S) := \{z \in \mathcal{Z} : \exists z' \in S : s(g(z')) - L_s L_g d(z, z') - \epsilon \geq h\}$$

Here, $R_\epsilon(S)$ contains the representations whose validity is certifiable from S up to accuracy ϵ . Let $R_\epsilon^H(S_0)$ denote its H -fold recursive application starting from S_0 . This enables viewing of ACTFLOW as executing a local-to-global expansion process toward the reachable valid set $R_\epsilon^H(S_0)$.

We can now present the main theorem, presented formally in Theorem F.5. It allows to state set-theoretic guarantees for out-of-distribution generative modeling.

(Informal) Assumptions:

- **(Well-specified verifier).** Verifier labels follow a calibrated logistic model in the learned representation space, $\mathbb{P}(y = 1 \mid z) = s(g(z))$, where the latent validity score g has bounded RKHS norm and is Lipschitz continuous in the learned representation.
- **(Approximate and local uncertainty sampling).** At each round, the sampler approximately maximizes verifier uncertainty over the current generable set Ω_t^τ , as by Eq. (17).
- **(EBM generative update).** We employ an energy-based model (EBM) abstraction, detailed in Assump. F.1, to abstract flow model endpoint density updates as learning an implicit energy functional over the flow-learned representation Z . This captures that flow models preserve high density on regions that our verifier regards as valid, and maintain low density on points far from this set.

Theorem 4.1 ((Informal) Generable representation set covers reachable set). *Fix $\epsilon > 0$ and an integer $H \geq 1$. Let τ and T^* satisfy*

$$\tau \leq \frac{h}{(1-h)\bar{Z}} \quad \text{and} \quad T^* \gtrsim \left(\frac{\alpha \gamma_{HT^*}^{\Omega_*}}{\epsilon} \right)^2, \quad (18)$$

with up to problem-dependent constants in \gtrsim . Define the maximum information gain from t samples by

$$\gamma_t^A \triangleq \max_{z_1, \dots, z_t \in A} \frac{1}{2} \log \det \left(I_t + (\lambda \kappa)^{-1} K_t \right)$$

By running ACTFLOW, it holds with probability at least $1 - \delta$ that after T^ verified samples,*

$$R_\epsilon^H(S_0) \subseteq \Omega_{T^*}^\tau. \quad (19)$$

Theorem 4.1 is stated in the learned representation space Z . Under a map ϕ with measure-preserving regularity conditions stated in Assumption F.4, the same reachability guarantee transfers to the design-space generable set. For $\tau_X > 0$, define the design-space generable set at round t by

$$\Omega_t^{X, \tau_X} := \{x \in \mathcal{X} : p_1^{\pi_t}(x) \geq \tau_X\}. \quad (20)$$

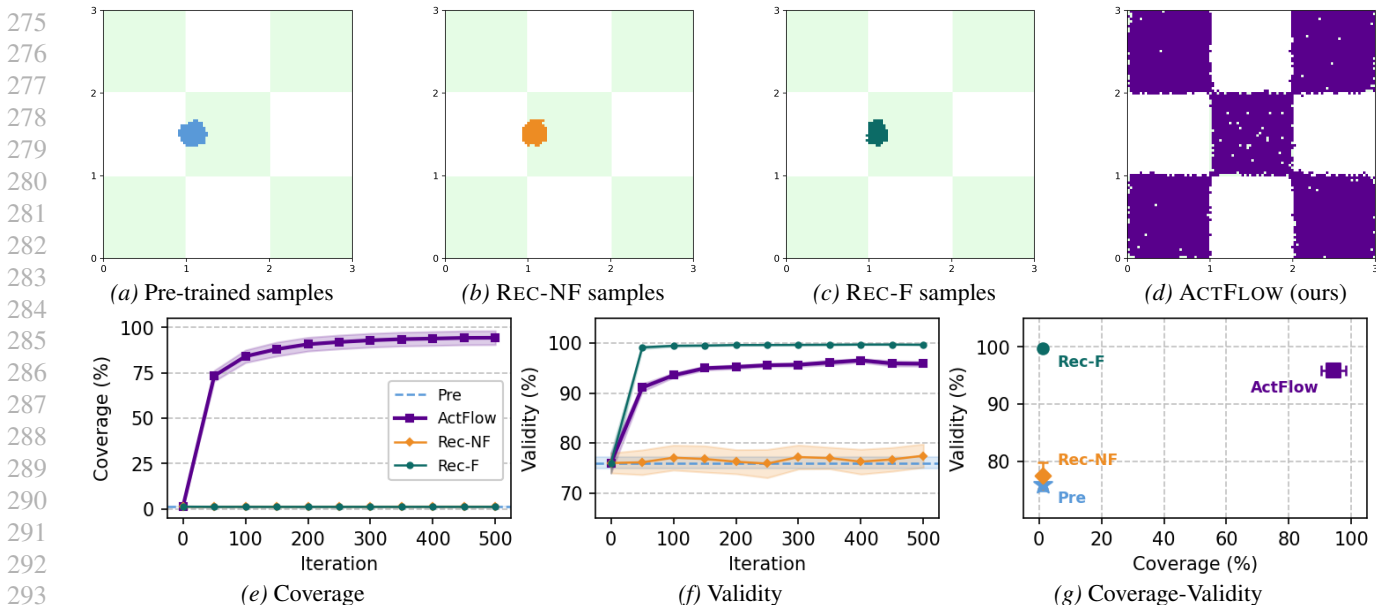


Figure 2. (2a) The valid design space (green) and the estimated generable set ($\tau = 0.01$) of a pre-trained model (blue). (2a-2d) REC-NF and REC-F fail to expand the pre-trained model’s generable set, while ACTFLOW is able to discover and expand over the entire valid design space. (2e-2g) ACTFLOW increases both model coverage (1.16% to 94.27%) and validity (76% to 95.9%), REC-NF fails at increasing both, while REC-F increases validity while even decreasing coverage (1.16% to 1.1%).

We also define the induced one-step reachability operator on \mathcal{X} by

$$R_{\epsilon}^{X,\phi}(A) := \left\{ x \in \mathcal{X} : \exists x' \in A \text{ s.t. } s(g(\phi(x'))) - L_s L_g d(\phi(x), \phi(x')) - \epsilon \geq h \right\}, \quad (21)$$

for any $A \subseteq \mathcal{X}$, and let $(R_{\epsilon}^{X,\phi})^H(A)$ denote its H -fold iterate.

Corollary 4.2 (Design-space coverage of the induced reachable valid set). *Assume the conditions of Theorem 4.1 and Assumption F.4, with $j_{\min} := \inf_{x \in \mathcal{X}} |\det J_{\phi}(x)| > 0$. Let $S_0^X := \phi^{-1}(S_0)$ and $\tau_X := j_{\min} \tau$. Then, after the same number T^* of verified samples as in Theorem 4.1, with probability at least $1 - \delta$,*

$$(R_{\epsilon}^{X,\phi})^H(S_0^X) \subseteq \Omega_{T^*}^{X,\tau_X}. \quad (22)$$

From reachable-set to full coverage. The guarantee is reachability-based: if Ω^* contains components that are disconnected from S_0^X in the learned geometry, local expansion provides no guarantee of covering those components. Conversely, if every $x \in \Omega^*$ can be connected to S_0^X by at most H local valid expansions in representation space, then Corollary F.10 yields full valid-space coverage,

$$\Omega^* \subseteq \Omega_{T^*}^{X,\tau_X}.$$

Theorem 4.1 shows that ACTFLOW expands the model’s τ -generable set for any sufficiently small τ , effectively describ-

ing finite-sample coverage. Corollary F.8 controls model validity, i.e., how tight the expanded generable set is w.r.t. Ω^* – the opposite set-inequality for *generable set expansion*.

5. Experimental Evaluation of ActFlow on Molecules, Peptides, and Proteins

We evaluate ACTFLOW for expanding the valid generable set of pre-trained generative models, and compare it against widely adopted self-generation baselines: continued pre-training on unfiltered model samples (REC-NF) (Shumailov et al., 2024; Alemohammad et al., 2023) and on verifier-filtered samples (REC-F) (Dong et al., 2023; Gulcehre et al., 2023). We propose two types of experiments: (i) an illustrative visually interpretable setting, and (ii) high-dimensional biochemical design tasks over molecules, therapeutic peptides, and protein sequences. Since standard generative modeling evaluation metrics are misaligned with OOD generative modeling, we employ the OOD modeling criteria presented in Apx. H.1, namely *coverage* via number of valid clusters, *diversity* via Vendi (Friedman & Dieng, 2022), and overall model validity. Additional experimental details are provided in Apx. H.6.

Illustrative Visually Interpretable Setting. We first evaluate ACTFLOW in a two-dimensional illustrative setting with valid design space shown in green in Fig. 2a–2d, and visualize the models’ generable sets via discretization. All methods run for $T = 500$ iterations and generate $B = 64$ samples per iteration. ACTFLOW uses $\alpha_t = 0.005$,

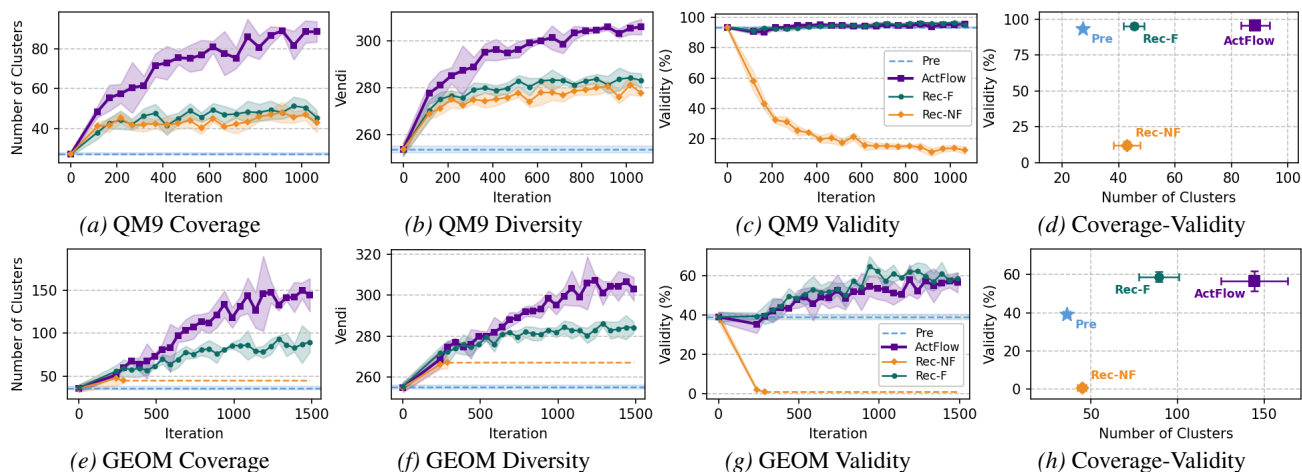


Figure 3. (3a-3d) Molecular design on QM9 (Ramakrishnan et al., 2014). ACTFLOW expands valid coverage substantially more than REC-NF and REC-F, reaching 88.40 valid clusters versus 42.80 and 45.40 respectively (Fig. 3a). It also achieves the highest diversity (3b) while preserving high validity (3d). (3e-3h) Design of drug-like molecules on GEOM-Drugs (Axelrod & Gomez-Bombarelli, 2022). ACTFLOW increases the pre-trained model coverage, diversity, and validity, vastly outperforming both REC-F and REC-NF baselines.

$\beta = 1/13$, and flow representation at timestep $s = 0.9$. Fig. 2a–2d show that ACTFLOW (violet) expands the pre-trained generable set ($\tau = 0.01$, Fig. 2a) to near-optimally cover the valid design space (Fig. 2c), whereas both baselines fail to expand it. Fig. 2e–2g show that ACTFLOW increases both coverage, from 1.16% to 94.27%, and validity, from 76.00% to 95.89%. By contrast, REC-NF leaves coverage unchanged and barely improves validity, while REC-F attains higher validity but decreases coverage. Thus, even in this simple two-dimensional setting, ACTFLOW is the only method that achieves OOD generable-set expansion, substantially improving both coverage and validity, expanding through sparse directions (corners) to new valid regions.

Molecular design on QM9. We evaluate ACTFLOW on FlowMol Gaussian (Dunn & Koes, 2024), pre-trained on QM9 (Ramakrishnan et al., 2014). All methods run for 1000 iterations after 66 initial iterations without fine-tuning. ACTFLOW uses the flow representation at timestep $s = 0.9$ and $\beta = 1/10$. Fig. 3a–3d show that, relative to the pre-trained model, ACTFLOW substantially expands valid molecular coverage, reaching 88.40 valid clusters, compared to 45.40 for REC-F and 42.80 for REC-NF. ACTFLOW also achieves the highest diversity, with Vendi 306.08, while preserving high validity (95.90%). In contrast, REC-F preserves validity (95.24%) but expands significantly less, whereas REC-NF severely degrades validity (12.26%). Thus, on QM9, ACTFLOW significantly outperforms both recursive sampling baselines, achieving a significantly stronger combination of valid coverage, diversity, and validity.

Molecular design on GEOM-Drugs. We further evaluate ACTFLOW on FlowMol Gaussian (Dunn & Koes, 2024), pre-trained on GEOM-Drugs (Axelrod & Gomez-Bombarelli,

2022), a substantially larger, and more chemically relevant dataset of drug-like molecules. All methods use 2000 fine-tuning steps after a warm-up period, where 4096 samples are acquired. ACTFLOW uses flow representation timestep $s = 0.8$, $\beta = 1/7$, $\alpha_t = 0$. Fig. 3e–3h show that ACTFLOW expands valid molecular coverage substantially beyond the baselines, reaching 144.3 valid clusters, compared to 89.33 for REC-F and 44.67 for REC-NF. ACTFLOW also achieves the highest diversity, with Vendi 303.1, while maintaining comparable validity to REC-F (56.6% versus 58.83%). In contrast, REC-F expands substantially less, whereas REC-NF collapses in validity, as reported in Fig. 3h. Thus, on GEOM-Drugs, ACTFLOW provides a stronger expansion profile than both baselines, ultimately increasing the pre-trained model coverage by 144.3% and its validity by 56.6%.

Therapeutic Peptide Design We assess ACTFLOW on biological sequence generation via discrete diffusion models, as detailed in Apx. G.3. We consider the task of therapeutic peptide design (Wang et al., 2022), and employ the pre-trained SMILES discrete diffusion model from PepTune (Tang et al., 2025). As shown in Table 2 and Fig 5, ACTFLOW significantly increases the number of clusters (i.e., coverage) from 133.3 to 197, while further increasing validity, determined via the SMILES2PEPTIDE verifier (Tang et al., 2025), from 31.9% to 52.0%. Similarly, ACTFLOW achieves the highest diversity of 7.45 Vendi

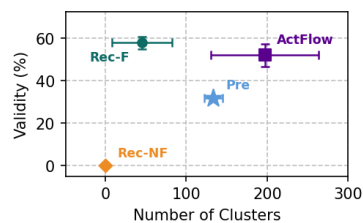


Figure 5. Peptides coverage-validity.

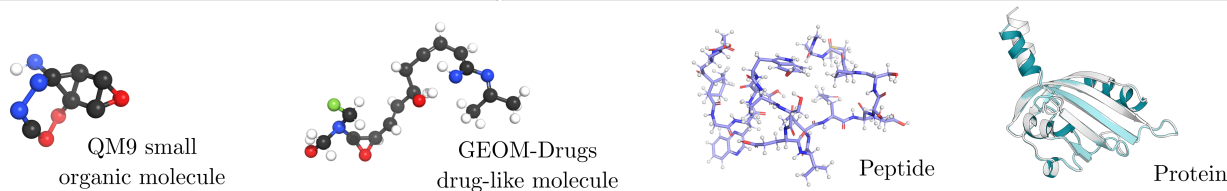


Figure 4. 3D structures directly generated or string-converted by models expanded via ACTFLOW.

| Method | Therapeutic Peptide Design | | | | Protein Sequence Design | | | |
|---------|----------------------------|----------------------|-----------------|-------------------------|-------------------------|----------------------|-----------------|-------------------------|
| | Coverage \uparrow | Diversity \uparrow | FID | Validity (%) \uparrow | Coverage \uparrow | Diversity \uparrow | FID | Validity (%) \uparrow |
| Pre | 133.30 \pm 11.1 | 5.70 \pm 0.11 | 0.00 \pm 0.00 | 31.92 \pm 0.77 | 66.50 \pm 5.63 | 12.87 \pm 0.63 | 0.05 \pm 0.01 | 70.81 \pm 1.12 |
| REC-NF | 0.00 \pm 0.00 | 0.00 \pm 0.00 | 0.00 \pm 0.00 | 0.00 \pm 0.00 | 63.75 \pm 11.45 | 11.85 \pm 0.34 | 0.27 \pm 0.06 | 67.81 \pm 3.11 |
| REC-F | 45.00 \pm 37.22 | 6.03 \pm 1.19 | 4.99 \pm 1.53 | 57.84 \pm 2.97 | 49.50 \pm 9.60 | 11.67 \pm 0.40 | 0.25 \pm 0.04 | 88.12 \pm 1.42 |
| ACTFLOW | 197.00 \pm 66.25 | 7.45 \pm 6.17 | 4.49 \pm 6.44 | 52.02 \pm 5.34 | 102.75 \pm 18.36 | 42.14 \pm 10.85 | 5.45 \pm 3.34 | 83.74 \pm 2.70 |

Table 2. ACTFLOW significantly expands valid model coverage (i.e., number of valid clusters) and diversity (i.e., Vendi (Friedman & Dieng, 2022)), while also increasing validity across therapeutic peptide and protein sequence design tasks. We report in red models where coverage or diversity decreased from initial model. These results show that ACTFLOW vastly outperforms widely adopted synthetic pre-training baselines.

score over PeptideCLM embeddings (Feller & Wilke, 2025). On the contrary, REC-F decreases the pre-trained model coverage to 45 clusters, while increasing validity, and REC-NF collapses entirely to 0.0% validity, as a single misplaced token can result in an invalid peptide – thus inducing no valid clusters. These results show that ACTFLOW’s strong empirical performance extends beyond continuous flows to discrete diffusion models (see Apx. G.3), significantly increasing jointly model coverage, diversity, and validity, and outperforming baselines.

Protein Sequence Design. We finally evaluate ACTFLOW on protein sequence design, via a continuous ESM diffusion model from SGPO (Yang et al., 2025), pre-trained on the CreiLOV fluorescence dataset (Chen et al., 2023). We use 512 iterations with 1000 fine-tuning steps each. ACTFLOW employs flow representation timestep at $t = 0.8$. ACTFLOW substantially expands valid coverage, computed over token-level ESM embeddings, increasing the number of clusters from 66.50 to 102.75, compared to 63.75 for REC-NF and 49.50 for REC-F, which lead to significant decrease in both coverage and validity metrics, see Table 2 and Fig. 6. ACTFLOW also achieves the highest diversity, with Vendi 42.14 vs 12.87 for the pre-trained model and 11.85/11.67 for the baselines, while improving validity from 70.81% to 83.74%. In contrast, REC-F reaches higher validity (88.12%) but col-

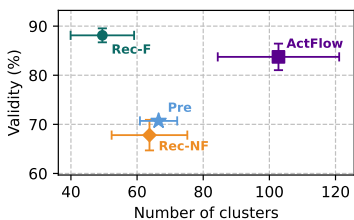


Figure 6. Proteins coverage-validity.

lapses in coverage and diversity; REC-NF reduces coverage and diversity relative to the pre-trained model. ACTFLOW achieves the strongest coverage–diversity–validity profile (see Fig. 6), significantly expanding the pre-trained protein sequence model to new valid regions of the design space.

6. Related Work

We present in Apx. A a discussion of related works spanning diffusion and flow model reward adaptation, synthetic data generation for model adaptation, flow-based design space exploration, and safe active exploration theory.

7. Conclusion

We depart from data distribution matching and formulate *out-of-distribution flow modeling* via *generable set expansion*, i.e., expanding the valid region of design space a model samples with non-negligible probability. We introduced ACTFLOW, a continued-pretraining scheme that uses verifier feedback on self-generated data to actively expand in the learned flow representation – leading to first-of-their-kind reachability-based guarantees for out-of-distribution flow modeling. Across small organic molecules, drug-like molecules, therapeutic peptides, and protein sequences, ACTFLOW consistently improves coverage, diversity, and validity over widely adopted recursive self-generation baselines. As for limitations, while our framework allows for task-agnostic expansions toward valid, previously inaccessible regions where new-to-nature discoveries may reside – future work will need to assess whether this form of exploration yields concrete gains (i.e., discoveries) in specific real-world applications.

Impact Statement

This paper presents work whose goal is to advance the field of Generative Modeling. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

Alberti, S., Hasanaliyev, K., Shah, M., and Ermon, S. Data unlearning in diffusion models. *arXiv preprint arXiv:2503.01034*, 2025.

Alemohammad, S., Casco-Rodriguez, J., Luzi, L., Humayun, A. I., Babaei, H., LeJeune, D., Siahkoochi, A., and Baraniuk, R. Self-consuming generative models go mad. In *The Twelfth International Conference on Learning Representations*, 2023.

Alemohammad, S., Humayun, A. I., Agarwal, S., Colomosse, J., and Baraniuk, R. Self-improving diffusion models with synthetic data. *arXiv preprint arXiv:2408.16333*, 2024.

Axelrod, S. and Gomez-Bombarelli, R. Geom, energy-annotated molecular conformations for property prediction and molecular generation. *Scientific Data*, 9(1):185, 2022.

Azizi, S., Kornblith, S., Saharia, C., Norouzi, M., and Fleet, D. J. Synthetic data from diffusion models improves imagenet classification. *arXiv preprint arXiv:2304.08466*, 2023.

Berkenkamp, F., Turchetta, M., Schoellig, A., and Krause, A. Safe model-based reinforcement learning with stability guarantees. *Advances in neural information processing systems*, 30, 2017.

Boucheron, S., Lugosi, G., and Massart, P. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.

Celik, O., Li, Z., Blessing, D., Li, G., Palenicek, D., Peters, J., Chalvatzaki, G., and Neumann, G. Dime: Diffusion-based maximum entropy reinforcement learning. *arXiv preprint arXiv:2502.02316*, 2025.

Chen, Y., Hu, R., Li, K., Zhang, Y., Fu, L., Zhang, J., and Si, T. Deep mutational scanning of an oxygen-independent fluorescent protein creilov for comprehensive profiling of mutational and epistatic effects. *ACS Synthetic Biology*, 12(5):1461–1473, 2023.

Chi, C., Xu, Z., Feng, S., Cousineau, E., Du, Y., Burchfiel, B., Tedrake, R., and Song, S. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025.

Corso, G., Stärk, H., Jing, B., Barzilay, R., and Jaakkola, T. Diffdock: Diffusion steps, twists, and turns for molecular docking. *arXiv preprint arXiv:2210.01776*, 2022.

De Santi, R., Vlastelica, M., Hsieh, Y.-P., Shen, Z., He, N., and Krause, A. Provable maximum entropy manifold exploration via diffusion models. In *International Conference on Machine Learning*, 2025a.

De Santi, R., Vlastelica, M., Hsieh, Y.-P., Shen, Z., He, N., and Krause, A. Flow density control: Generative optimization beyond entropy-regularized fine-tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025b.

De Santi, R., Protopapas, K., Hsieh, Y.-P., and Krause, A. Verifier-constrained flow expansion for discovery beyond the data. In *International Conference on Learning Representations (ICLR)*, April 2026.

Dong, H., Xiong, W., Goyal, D., Zhang, Y., Chow, W., Pan, R., Diao, S., Zhang, J., Shum, K., and Zhang, T. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.

Dunn, I. and Koes, D. R. Mixed continuous and categorical flow matching for 3d de novo molecule generation. *ArXiv*, pp. arXiv–2404, 2024.

Feller, A. L. and Wilke, C. O. Peptide-aware chemical language model successfully predicts membrane diffusion of cyclic peptides. *Journal of chemical information and modeling*, 65(2):571–579, 2025.

Friedman, D. and Dieng, A. B. The vendi score: A diversity evaluation metric for machine learning. *arXiv preprint arXiv:2210.02410*, 2022.

Gulcehre, C., Paine, T. L., Srinivasan, S., Konyushkova, K., Weerts, L., Sharma, A., Siddhant, A., Ahern, A., Wang, M., Gu, C., et al. Reinforced self-training (rest) for language modeling. *arXiv preprint arXiv:2308.08998*, 2023.

Guo, J. and Schwaller, P. It takes two to tango: Directly optimizing for constrained synthesizability in generative molecular design. *arXiv preprint arXiv:2410.11527*, 2024.

Gutjahr, S., De Santi, R., Schaufelberger, L., Jorner, K., and Krause, A. Constrained molecular generation via sequential flow model fine-tuning. In *ICML 2025 Generative AI and Biology (GenBio) Workshop*, 2025.

Halgren, T. A. Merck molecular force field. i. basis, form, scope, parameterization, and performance of mmff94. *Journal of computational chemistry*, 17(5-6):490–519, 1996.

- 495 Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and
496 Hochreiter, S. Gans trained by a two time-scale update
497 rule converge to a local nash equilibrium. *Advances in*
498 *neural information processing systems*, 30, 2017.
499
- 500 Ho, J., Jain, A., and Abbeel, P. Denoising diffusion proba-
501 bilistic models. *Advances in neural information process-*
502 *ing systems*, 33:6840–6851, 2020.
503
- 504 Hoogeboom, E., Satorras, V. G., Vignac, C., and Welling, M.
505 Equivariant diffusion for molecule generation in 3d. In
506 *International conference on machine learning*, pp. 8867–
507 8887. PMLR, 2022.
- 508 Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnoy, M.,
509 Ronneberger, O., Tunyasuvunakool, K., Bates, R., Žídek,
510 A., Potapenko, A., et al. Highly accurate protein structure
511 prediction with alphafold. *nature*, 596(7873):583–589,
512 2021.
513
- 514 Koller, T., Berkenkamp, F., Turchetta, M., and Krause, A.
515 Learning-based model predictive control for safe explo-
516 ration. In *2018 IEEE conference on decision and control*
517 *(CDC)*, pp. 6059–6066. IEEE, 2018.
518
- 519 Laurent, B. and Massart, P. Adaptive estimation of a
520 quadratic functional by model selection. *The Annals*
521 *of Statistics*, 28(5):1302–1338, 2000.
522
- 523 Li, X. and Fourches, D. Smiles pair encoding: a data-driven
524 substructure tokenization algorithm for deep learning.
525 *Journal of chemical information and modeling*, 61(4):
526 1560–1569, 2021.
- 527 Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W.,
528 Smetanin, N., dos Santos Costa, A., Fazel-Zarandi, M.,
529 Sercu, T., Candido, S., et al. Language models of pro-
530 tein sequences at the scale of evolution enable accurate
531 structure prediction. *bioRxiv*, 2022.
532
- 533 Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and
534 Le, M. Flow matching for generative modeling. *arXiv*
535 *preprint arXiv:2210.02747*, 2022.
536
- 537 Lipman, Y., Havasi, M., Holderrieth, P., Shaul, N., Le, M.,
538 Karrer, B., Chen, R. T., Lopez-Paz, D., Ben-Hamu, H.,
539 and Gat, I. Flow matching guide and code. *arXiv preprint*
540 *arXiv:2412.06264*, 2024.
- 541 Liu, J., Liu, G., Liang, J., Li, Y., Liu, J., Wang, X., Wan,
542 P., Zhang, D., and Ouyang, W. Flow-grpo: Training
543 flow matching models via online rl. *arXiv preprint*
544 *arXiv:2505.05470*, 2025.
545
- 546 Lou, A., Meng, C., and Ermon, S. Discrete diffusion model-
547 ing by estimating the ratios of the data distribution. *arXiv*
548 *preprint arXiv:2310.16834*, 2023.
549
- Nguyen, D., Li, J., Zheng, J., and Mirzasoleiman, B. Do we
need all the synthetic data? targeted image augmentation
via diffusion models. In *The Fourteenth International*
Conference on Learning Representations, 2025.
- Ou, J., Nie, S., Xue, K., Zhu, F., Sun, J., Li, Z., and Li,
C. Your absorbing discrete diffusion secretly models the
conditional distributions of clean data. *arXiv preprint*
arXiv:2406.03736, 2024.
- Pásztor, B., Kassraie, P., and Krause, A. Bandits with
preference feedback: A stackelberg game perspective.
Advances in Neural Information Processing Systems, 37:
11997–12034, 2024.
- Pearce, T., Rashid, T., Kanervisto, A., Bignell, D., Sun,
M., Georgescu, R., Macua, S. V., Tan, S. Z., Momenne-
jad, I., Hofmann, K., et al. Imitating human behaviour
with diffusion models. *arXiv preprint arXiv:2301.10677*,
2023.
- Perez Jensen, C., Schaufelberger, L., Santi, R. D., Jorner, K.,
and Krause, A. Value matching: Scalable and gradient-
free reward-guided flow adaptation. In *The Fourteenth*
International Conference on Learning Representations,
2026. URL <https://openreview.net/forum?id=7iXt44Actj>.
- Ramakrishnan, R., Dral, P. O., Rupp, M., and Von Lilienfeld,
O. A. Quantum chemistry structures and properties of
134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- Rector-Brooks, J., Lambert, T., Skreta, M., Roth, D., Long,
Y., Li, Z.-Q., Zhang, X., Cretu, M., Li, F.-Z., Ganapathy,
T., et al. General multimodal protein design enables dna-
encoding of chemistry. *arXiv preprint arXiv:2604.05181*,
2026.
- Sahoo, S., Arriola, M., Schiff, Y., Gokaslan, A., Marroquin,
E., Chiu, J., Rush, A., and Kuleshov, V. Simple and
effective masked diffusion language models. *Advances*
in Neural Information Processing Systems, 37:130136–
130184, 2024.
- Schölkopf, B., Herbrich, R., and Smola, A. J. A general-
ized representer theorem. In *International conference on*
computational learning theory, pp. 416–426. Springer,
2001.
- Shi, J., Han, K., Wang, Z., Doucet, A., and Titsias, M.
Simplified and generalized masked diffusion for discrete
data. *Advances in neural information processing systems*,
37:103131–103167, 2024.
- Shumailov, I., Shumaylov, Z., Zhao, Y., Papernot, N., Ander-
son, R., and Gal, Y. Ai models collapse when trained on
recursively generated data. *Nature*, 631(8022):755–759,
2024.

- 550 Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and
551 Ganguli, S. Deep unsupervised learning using nonequi-
552 librium thermodynamics. In *International conference on*
553 *machine learning*, pp. 2256–2265. PMLR, 2015.
- 554 Song, Y. and Ermon, S. Generative modeling by estimating
555 gradients of the data distribution. *Advances in neural*
556 *information processing systems*, 32, 2019.
- 558 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Er-
559 mon, S., and Poole, B. Score-based generative modeling
560 through stochastic differential equations. *arXiv preprint*
561 *arXiv:2011.13456*, 2020.
- 563 Sui, Y., Gotovos, A., Burdick, J., and Krause, A. Safe
564 exploration for optimization with gaussian processes. In
565 *International conference on machine learning*, pp. 997–
566 1005. PMLR, 2015.
- 567 Sui, Y., Zhuang, V., Burdick, J., and Yue, Y. Stagewise
568 safe bayesian optimization with gaussian processes. In
569 *International conference on machine learning*, pp. 4781–
570 4789. PMLR, 2018.
- 572 Tang, S., Zhang, Y., and Chatterjee, P. Peptune: De novo
573 generation of therapeutic peptides with multi-objective-
574 guided discrete diffusion. *42nd International Conference*
575 *on Machine Learning*, 2025.
- 577 Turchetta, M., Berkenkamp, F., and Krause, A. Safe explo-
578 ration in finite markov decision processes with gaussian
579 processes. *Advances in neural information processing*
580 *systems*, 29, 2016.
- 581 Uehara, M., Zhao, Y., Biancalani, T., and Levine, S. Un-
582 derstanding reinforcement learning-based fine-tuning of
583 diffusion models: A tutorial and review, 2024. URL
584 <https://arxiv.org/abs/2407.13734>.
- 586 Uehara, M., Zhao, Y., Wang, C., Li, X., Regev, A., Levine,
587 S., and Biancalani, T. Inference-time alignment in diffu-
588 sion models with reward-guided generation: Tutorial and
589 review. *arXiv preprint arXiv:2501.09685*, 2025.
- 591 Wang, L., Wang, N., Zhang, W., Cheng, X., Yan, Z., Shao,
592 G., Wang, X., Wang, R., and Fu, C. Therapeutic pep-
593 tides: current applications and future directions. *Signal*
594 *transduction and targeted therapy*, 7(1):48, 2022.
- 595 Watson, J. L., Juergens, D., Bennett, N. R., Trippe, B. L.,
596 Yim, J., Eisenach, H. E., Ahern, W., Borst, A. J., Ragotte,
597 R. J., Milles, L. F., et al. De novo design of protein struc-
598 ture and function with rdiffusion. *Nature*, 620(7976):
599 1089–1100, 2023.
- 601 Weininger, D. Smiles, a chemical language and information
602 system. 1. introduction to methodology and encoding
603 rules. *J. Chem. Inf. Comput. Sci.*, 28(1):31–36, 1988.
- 604 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 opti-
mization, 2025. URL <https://arxiv.org/abs/2505.15093>.
- Yuan, H., Chen, Z., Ji, K., and Gu, Q. Self-play fine-tuning
of diffusion models for text-to-image generation. *Ad-
vances in Neural Information Processing Systems*, 37:
73366–73398, 2024.
- 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 categori-
cal sampling. *arXiv preprint arXiv:2409.02908*, 2024.

A. Related Work

Diffusion and flow model reward adaptation for generative optimization. Several works adapt pre-trained diffusion and flow models for reward maximization, either via fine-tuning (e.g., Uehara et al., 2024; De Santi et al., 2025b; Liu et al., 2025) or inference-time sampling (e.g., Uehara et al., 2025; Perez Jensen et al., 2026; Uehara et al., 2025). Without injecting further validity information, these methods remain constrained by the pre-trained model coverage, leading to *over-optimization*, i.e., invalid samples, when steering the model excessively beyond the training distribution (Gutjahr et al., 2025). In this work, we formalize this limitation through the notion of *generable set*, and introduce ACTFLOW, a task-agnostic continued pre-training method that expands coverage to new valid regions of the design space.

Synthetic data generation for diffusion model self-adaptation Generative priors can provide useful synthetic data for downstream learning (e.g., Azizi et al., 2023; Nguyen et al., 2025), but recursively training on generated data can induce model collapse (Shumailov et al., 2024; Alemohammad et al., 2023). For generative model adaptation, prior work uses synthetic data for negative guidance (Alemohammad et al., 2024) or conservative verifier-free self-play (Yuan et al., 2024), with the primary goal of improving in-distribution sample quality (e.g., Alemohammad et al., 2024). In contrast, we study flow adaptation for out-of-distribution modeling: rather than further refining high-density modes, the goal is to reallocate mass toward newly verified valid regions beyond the training distribution. To our knowledge, ACTFLOW is the first theory-backed synthetic continued-pretraining method for out-of-distribution generative modeling.

Diffusion and flow based design space exploration Recent works introduce scalable methods for diffusion- and flow-based design space exploration via entropy maximization (e.g., De Santi et al., 2025a;b) or approximations (e.g., Celik et al., 2025). A more recent line assumes access to a known differentiable verifier and uses it to rebalance model density toward valid regions via verifier-constrained entropy maximization (De Santi et al., 2026). In contrast, we assume access only to black-box, unknown verifiers, covering highly-relevant settings in science, where feedback is experimental and/or non-differentiable. Moreover, ACTFLOW expands valid coverage through directed synthetic data generation rather than verifier-constrained optimization.

Safe (active) exploration theory Safe exploration studies learning under unknown safety constraints. Early work uses Gaussian-process confidence bounds to restrict exploration to points certified above a safety threshold (Sui et al., 2015; 2018). Subsequent safe-RL methods extend this principle to MDPs and control (e.g., Turchetta et al., 2016; Berkenkamp et al., 2017; Koller et al., 2018). We bridge *safety* in safe exploration with *validity* in generative modeling: expanding a generable set amounts to discovering new valid (i.e., safe) regions. In scientific discovery, invalid samples lead to rejected verifier queries, not unsafe actions in safety-critical systems; thus, ACTFLOW does not perform safe-set constrained planning. Rather, we employ this viewpoint to provide a first-of-its-kind reachability-based theory of out-of-distribution flow modeling.

B. Central Limitation of Standard Pre-training with an Imperfect Model

In the main text, Eq. (7) is stated for clarity under the idealization that the pre-trained model generates only valid designs. In practice, this need not hold: the generable set may contain invalid designs, so that

$$\Omega_{\theta_{\text{pre}}}^{\tau} \cap (\mathcal{X} \setminus \Omega^*) \neq \emptyset.$$

This does not change the central limitation. One simply applies the coverage statement to the *valid part* of the generable set,

$$\Omega_{\theta_{\text{pre}}}^{\tau, \text{val}} := \Omega_{\theta_{\text{pre}}}^{\tau} \cap \Omega^* \subseteq \Omega^*.$$

The realistic limitation is therefore

$$\Omega_{\theta_{\text{pre}}}^{\tau, \text{val}} \subset \Omega^*, \quad \text{Vol}(\Omega_{\theta_{\text{pre}}}^{\tau, \text{val}}) \ll \text{Vol}(\Omega^*),$$

i.e., the pre-trained model covers only a small fraction of the valid design space with non-negligible probability, even if it also assigns mass to invalid regions.

This is precisely the setting addressed by ACTFLOW. The algorithm does not require the initial generable set to be valid; it only requires verifier feedback to identify which generated samples are valid and to adapt the model toward newly verified valid regions. Thus, generable set expansion should be understood as expansion of the valid, verifier-certified portion of the model’s coverage. Moreover, when the initial model has partial validity, reallocating mass toward newly verified valid regions can also increase the model’s overall validity, as we observe in our experiments, where the pre-trained models are imperfect and ACTFLOW systematically improves not only coverage and diversity, but typically increases substantially validity as well.

Volume over the valid design space. Throughout, $\text{Vol}(\cdot)$ denotes volume with respect to a domain-appropriate reference measure on the valid design space Ω^* . Formally, let ν^* be a reference measure supported on Ω^* : Lebesgue measure when Ω^* is full-dimensional, intrinsic Hausdorff measure when Ω^* is a lower-dimensional manifold, and counting measure in discrete design spaces. For any $A \subseteq \mathcal{X}$, we write

$$\text{Vol}(A) := \nu^*(A \cap \Omega^*).$$

Thus, the central limitation in Eq. (7) measures how much of the valid design space is covered by the model’s τ -generable set, rather than ambient Lebesgue volume in \mathcal{X} . This avoids degeneracies in settings where valid designs lie on lower-dimensional manifolds or discrete spaces.

C. A Warm-Up Analysis of Failure Modes of Expansion via Synthetic Data

C.1. Compact Analysis

Warm Up Gaussian Setting. A widely adopted approach to synthetic data generation is recursive sampling with closed-loop verifier filtering. The purpose of this section is *not* to establish a general impossibility result for this paradigm, but to use a minimal analytical model to isolate possible failure modes relevant to generable-set expansion. To this end, we introduce a simple probabilistic notion of the *generable valid frontier*: valid samples that remain reliably generable under the current model, yet lie away from its dominant modes. Such points form a natural abstraction of expansion-enabling samples. The analysis then shows that standard recursive self-generation schemes may be unlikely to produce them even in a simple low-dimensional setting, thereby shedding light on possible failure modes of data generation schemes and motivating algorithmic desiderata for expansion. We report derivations in Apx. C.2.

A Gaussian abstraction of pre-training. We consider a design space $\mathcal{X} = \mathbb{R}^d$. Let $U \subset \mathbb{R}^d$ be a k -dimensional linear subspace, let $m := d - k$, and write $\mathbb{R}^d = U \oplus U^\perp$. We interpret U as a low-dimensional region well captured by pre-training, and U^\perp as orthogonal directions along which expansion beyond the dominant pre-trained modes could occur. For $x \in \mathbb{R}^d$, write $x = x_U + x_\perp$ with $x_U = \Pi_U x$ and $x_\perp = \Pi_{U^\perp} x$. We model the pre-trained generator p^{θ_0} by the anisotropic Gaussian

$$X \sim p^{\theta_0} =: p_0 = \mathcal{N}(0, I_U \oplus \sigma^2 I_{U^\perp}), \quad 0 < \sigma \ll 1.$$

Thus, the model is spread along U , while it is sharply concentrated around U in the orthogonal directions. For a density threshold $\tau > 0$, the corresponding generable set is

$$\Omega_{\theta_0}^\tau = \{x \in \mathbb{R}^d : \|x_U\|_2^2 + \sigma^{-2} \|x_\perp\|_2^2 \leq r_\tau^2\}, \quad r_\tau^2 := 2 \log \left(\frac{(2\pi)^{-d/2} \sigma^{-m}}{\tau} \right),$$

namely an ellipsoid elongated along U and very thin along U^\perp .

A probabilistic notion of generable valid frontier. To reason about samples that may support expansion, we isolate a shell inside $\Omega_{\theta_0}^\tau$ that is close to its orthogonal boundary. Fix $R_U > 0$ and radii $0 < \rho_- < \rho_+$ such that $R_U < r_\tau$ and $\rho_+ < \sigma \sqrt{r_\tau^2 - R_U^2}$. We denote the *generable frontier* by

$$\mathcal{F}_{\text{frontier}} := \{x \in \mathbb{R}^d : \|x_U\|_2 \leq R_U, \rho_- \leq \|x_\perp\|_2 \leq \rho_+\}.$$

By construction, $\mathcal{F}_{\text{frontier}} \subseteq \Omega_{\theta_0}^\tau$. Thus, points in $\mathcal{F}_{\text{frontier}}$ remain reliably generable under the current model, since they lie inside its generable set, while already exhibiting a substantial orthogonal deviation from the model dominant modes. Next, fix a unit vector $u_\star \in U^\perp$ and an opening angle $\phi \in (0, \pi/2)$. We define the cone of *valid directions of expansion*

$$C_\phi(u_\star) := \left\{ y \in U^\perp \setminus \{0\} : \left\langle \frac{y}{\|y\|_2}, u_\star \right\rangle \geq \cos \phi \right\}.$$

We call *generable valid frontier* the set

$$\mathcal{V}_{\text{frontier}} := \{x \in \mathcal{F}_{\text{frontier}} : x_\perp \in C_\phi(u_\star)\}.$$

Thus $\mathcal{V}_{\text{frontier}}$ consists of samples that lie in the model generable frontier, and are aligned with a valid direction of expansion. While $\mathcal{V}_{\text{frontier}}$ does not provide a general first-principles characterization of all useful samples, it formally captures in a minimal way the qualitative tension we want to study: useful self-generated samples should plausibly be (i) *generable* by the current model, yet (ii) *not* lying within its dominant modes, and (iii) aligned with valid directions of expansion. The next result shows that standard self-generation scheme might sample data within this natural class of frontier-valid samples with extremely low probability even in the introduced illustrative Gaussian generator setting.

Proposition 1 (Standard self-generation is unlikely to find frontier-valid samples). *Let $X \sim p_0$, and let $\mathcal{V}_{\text{frontier}}$ be defined above. For any $\eta \in (0, 1)$, define $\rho_\eta := \sigma \sqrt{m + 2\sqrt{m \log(1/\eta)} + 2 \log(1/\eta)}$. If $\rho_- = \rho_\eta$, then*

$$\Pr(X \in \mathcal{V}_{\text{frontier}}) \leq \eta \exp \left(-\frac{m-1}{2} \cos^2 \phi \right).$$

Consequently, for N i.i.d. passive samples $X_1, \dots, X_N \sim p_0$,

$$\Pr(\exists i \in [N] : X_i \in \mathcal{V}_{\text{frontier}}) \leq N\eta \exp \left(-\frac{m-1}{2} \cos^2 \phi \right).$$

What the proposition reveals. The proposition isolates two sources of difficulty in recursive self-generation. The η -term reflects the rarity of reaching the generable frontier at all, namely of producing a sample with substantial orthogonal deviation from the dominant modes. Beyond this, an additional exponential penalty governs the chance of landing in the correct valid cone. Hence, when valid expansion directions are sparse, the sample complexity of passively obtaining a frontier-valid point becomes exponential in the orthogonal dimension m .

The algorithmic implication is clear: closed-loop verifier filtering alone is not enough for reliable expansion. Successful expansion requires actively steering generation toward the generable valid frontier.

C.2. Extensive Analysis with Proofs

Gaussian warm-up: self-generation requires rare valid frontier events We present a simple Gaussian model that isolates the two difficulties behind standard self-generation for efficient out-of-distribution discovery: sampling a *large deviation* away from the low-dimensional region captured by pre-training, and doing so along a *valid direction* of expansion.

We work directly on the design space $\mathcal{X} = \mathbb{R}^d$, so in this warm-up there is no separate representation map and $\mathcal{X} = \mathcal{Z}$. Let $U \subset \mathbb{R}^d$ be a k -dimensional linear subspace, and write

$$m := d - k, \quad \mathbb{R}^d = U \oplus U^\perp.$$

For every $x \in \mathbb{R}^d$, denote by

$$x_U := \Pi_U x, \quad x_\perp := \Pi_{U^\perp} x,$$

so that $x = x_U + x_\perp$. Throughout this subsection we assume $m \geq 2$, since the directional effect of interest only appears when the orthogonal complement has dimension at least two.

Pre-trained model. We model a pre-trained generator p^{θ_0} by the anisotropic Gaussian

$$X \sim p^{\theta_0} =: p_0 = \mathcal{N}(0, I_U \oplus \sigma^2 I_{U^\perp}), \quad 0 < \sigma \ll 1.$$

Equivalently,

$$X_U \sim \mathcal{N}(0, I_k), \quad X_\perp \sim \mathcal{N}(0, \sigma^2 I_m),$$

independently. Its density is

$$p_0(x) = (2\pi)^{-d/2} \sigma^{-m} \exp\left(-\frac{1}{2}\left(\|x_U\|_2^2 + \sigma^{-2}\|x_\perp\|_2^2\right)\right).$$

Generable set. Fix a level $\epsilon \in (0, (2\pi)^{-d/2} \sigma^{-m})$. In the sense of Definition 1, the ϵ -level generable set of the pre-trained model is

$$\Omega_{\theta_0}^\epsilon := \{x \in \mathbb{R}^d : p_0(x) \geq \epsilon\}.$$

For the Gaussian above this is the ellipsoid

$$\Omega_{\theta_0}^\epsilon = \{x \in \mathbb{R}^d : \|x_U\|_2^2 + \sigma^{-2}\|x_\perp\|_2^2 \leq r_\epsilon^2\},$$

where

$$r_\epsilon^2 := 2 \log\left(\frac{(2\pi)^{-d/2} \sigma^{-m}}{\epsilon}\right).$$

Valid direction and frontier seed set. Fix a unit vector $u_\star \in U^\perp$ and an opening angle $\phi \in (0, \pi/2)$. We define the one-sided cone of valid orthogonal directions by

$$C_\phi(u_\star) := \left\{y \in U^\perp \setminus \{0\} : \left\langle \frac{y}{\|y\|_2}, u_\star \right\rangle \geq \cos \phi \right\}.$$

We also fix a radius $R_U > 0$ and shell radii $0 < \rho_- < \rho_+$ satisfying

$$R_U < r_\epsilon, \quad \rho_+ < \sigma \sqrt{r_\epsilon^2 - R_U^2}.$$

The corresponding *valid frontier seed set* is

$$\mathcal{V}_{\text{seed}} := \{x \in \mathbb{R}^d : \|x_U\|_2 \leq R_U, \rho_- \leq \|x_\perp\|_2 \leq \rho_+, x_\perp \in C_\phi(u_\star)\}.$$

These are rare points that are still inside the current generable set, but already lie on a specific valid orthogonal frontier along which future continued pre-training may expand.

Lemma C.1 (Frontier seeds lie inside the current generable set). *Under the above choice of R_U, ρ_+, ϵ , one has*

$$\mathcal{V}_{\text{seed}} \subseteq \Omega_{\theta_0}^\epsilon.$$

Proof. Take any $x \in \mathcal{V}_{\text{seed}}$. Then $\|x_U\|_2 \leq R_U$ and $\|x_\perp\|_2 \leq \rho_+$, hence

$$\|x_U\|_2^2 + \sigma^{-2} \|x_\perp\|_2^2 \leq R_U^2 + \sigma^{-2} \rho_+^2 < r_\epsilon^2.$$

By the explicit description of $\Omega_{\theta_0}^\epsilon$, this implies $x \in \Omega_{\theta_0}^\epsilon$. \square

A larger valid design space. We consider a valid design space that extends beyond the current generable set in the same direction u_\star . Fix any $R > \sigma \sqrt{r_\epsilon^2 - R_U^2}$, and define

$$\mathcal{V}_{\text{out}} := \left\{ x \in \mathbb{R}^d : \|x_U\|_2 \leq R_U, \sigma \sqrt{r_\epsilon^2 - R_U^2} < \|x_\perp\|_2 \leq R, x_\perp \in C_\phi(u_\star) \right\}.$$

We then set

$$S_0 := \{x \in \mathbb{R}^d : \|x_U\|_2 \leq R_U, \|x_\perp\|_2 \leq \rho_-\}, \quad \Omega^\star := S_0 \cup \mathcal{V}_{\text{seed}} \cup \mathcal{V}_{\text{out}}.$$

Thus Ω^\star contains an already-valid core S_0 , a rare but still generable frontier $\mathcal{V}_{\text{seed}}$, and a genuinely out-of-distribution valid region $\mathcal{V}_{\text{out}} \subset (\Omega_{\theta_0}^\epsilon)^c$. The role of the negative result below is to show that, even before trying to reach \mathcal{V}_{out} , passive self-generation is already unlikely to find the expansion-enabling seeds $\mathcal{V}_{\text{seed}}$.

A convenient large-deviation scale. Since $X_\perp \sim \mathcal{N}(0, \sigma^2 I_m)$, one has

$$\frac{\|X_\perp\|_2^2}{\sigma^2} \sim \chi_m^2.$$

Therefore, for any $\eta \in (0, 1)$, the choice

$$\rho_\eta := \sigma \sqrt{m + 2\sqrt{m \log(1/\eta)} + 2 \log(1/\eta)}$$

satisfies

$$\Pr(\|X_\perp\|_2 \geq \rho_\eta) \leq \eta$$

by the Laurent-Massart inequality (Laurent & Massart, 2000; Boucheron et al., 2013).

Proposition 1 (Standard self-generation is unlikely to find frontier-valid samples). *Let $X \sim p_0$, and let $\mathcal{V}_{\text{frontier}}$ be defined above. For any $\eta \in (0, 1)$, define $\rho_\eta := \sigma \sqrt{m + 2\sqrt{m \log(1/\eta)} + 2 \log(1/\eta)}$. If $\rho_- = \rho_\eta$, then*

$$\Pr(X \in \mathcal{V}_{\text{frontier}}) \leq \eta \exp\left(-\frac{m-1}{2} \cos^2 \phi\right).$$

Consequently, for N i.i.d. passive samples $X_1, \dots, X_N \sim p_0$,

$$\Pr(\exists i \in [N] : X_i \in \mathcal{V}_{\text{frontier}}) \leq N\eta \exp\left(-\frac{m-1}{2} \cos^2 \phi\right).$$

Proof. By definition of $\mathcal{V}_{\text{seed}}$,

$$\{X \in \mathcal{V}_{\text{seed}}\} \subseteq \{\|X_\perp\|_2 \geq \rho_-\} \cap \left\{ \left\langle \frac{X_\perp}{\|X_\perp\|_2}, u_\star \right\rangle \geq \cos \phi \right\}.$$

Therefore

$$\Pr(X \in \mathcal{V}_{\text{seed}}) \leq \Pr\left(\|X_{\perp}\|_2 \geq \rho_-, \left\langle \frac{X_{\perp}}{\|X_{\perp}\|_2}, u_{\star} \right\rangle \geq \cos \phi\right).$$

Now write

$$X_{\perp} = \sigma G, \quad G \sim \mathcal{N}(0, I_m).$$

Let

$$R := \|G\|_2, \quad S := \frac{G}{\|G\|_2} \in \mathbb{S}^{m-1}.$$

For an isotropic Gaussian, R and S are independent, and S is uniform on the unit sphere \mathbb{S}^{m-1} . Since $X_{\perp} = \sigma G$, the previous probability equals

$$\Pr\left(R \geq \frac{\rho_-}{\sigma}, \langle S, u_{\star} \rangle \geq \cos \phi\right) = \Pr\left(R \geq \frac{\rho_-}{\sigma}\right) \cdot \Pr(\langle S, u_{\star} \rangle \geq \cos \phi),$$

which proves

$$\Pr(X \in \mathcal{V}_{\text{seed}}) \leq \Pr(\|X_{\perp}\|_2 \geq \rho_-) \cdot \Pr\left(\left\langle \frac{X_{\perp}}{\|X_{\perp}\|_2}, u_{\star} \right\rangle \geq \cos \phi\right).$$

For the directional term, a standard spherical-cap bound for $S \sim \text{Unif}(\mathbb{S}^{m-1})$ gives, for every $t \in [0, 1]$,

$$\Pr(\langle S, u_{\star} \rangle \geq t) \leq \exp\left(-\frac{m-1}{2}t^2\right).$$

Applying this with $t = \cos \phi$ yields

$$\Pr\left(\left\langle \frac{X_{\perp}}{\|X_{\perp}\|_2}, u_{\star} \right\rangle \geq \cos \phi\right) = \Pr(\langle S, u_{\star} \rangle \geq \cos \phi) \leq \exp\left(-\frac{m-1}{2} \cos^2 \phi\right).$$

Combining the two bounds proves

$$\Pr(X \in \mathcal{V}_{\text{seed}}) \leq \Pr(\|X_{\perp}\|_2 \geq \rho_-) \exp\left(-\frac{m-1}{2} \cos^2 \phi\right).$$

If $\rho_- = \rho_{\eta}$, then by the above χ^2 -tail bound,

$$\Pr(\|X_{\perp}\|_2 \geq \rho_{\eta}) \leq \eta,$$

hence

$$\Pr(X \in \mathcal{V}_{\text{seed}}) \leq \eta \exp\left(-\frac{m-1}{2} \cos^2 \phi\right).$$

Finally, for N i.i.d. passive samples X_1, \dots, X_N , the union bound gives

$$\Pr(\exists i \in [N] : X_i \in \mathcal{V}_{\text{seed}}) \leq \sum_{i=1}^N \Pr(X_i \in \mathcal{V}_{\text{seed}}) = N \Pr(X \in \mathcal{V}_{\text{seed}}) \leq N\eta \exp\left(-\frac{m-1}{2} \cos^2 \phi\right).$$

□

D. Gradient Descent and Ascent on flow matching loss

D.1. Signed replay-based continued pre-training and connection to data unlearning

This appendix formalizes the signed continued-pretraining update used in Sec. 4 and clarifies its connection to data unlearning objectives for diffusion models (Alberti et al., 2025). The key point is that, after each active-query round, verifier feedback naturally partitions the queried synthetic data into an *accepted* set and a *rejected* set. This induces a retain/forget decomposition analogous to data unlearning, except that here the partition is generated online by the verifier and used to improve out-of-distribution validity coverage.

Per-round replay buffers. At round t , let

$$\mathcal{D}_t^+ \subseteq \{x_i : y_i = 1, i \leq t\}, \quad \mathcal{D}_t^- \subseteq \{x_i : y_i = 0, i \leq t\}$$

denote the accepted and rejected replay buffers accumulated up to round t . Thus the queried synthetic data available for continued pre-training are partitioned into

$$X_t := \mathcal{D}_t^+ \uplus \mathcal{D}_t^-, \quad A_t := \mathcal{D}_t^-,$$

where \uplus denotes disjoint union of indexed samples. By construction,

$$A_t \subseteq X_t, \quad X_t \setminus A_t = \mathcal{D}_t^+.$$

This is exactly the retain/forget decomposition considered in data unlearning: X_t is the current empirical dataset, A_t is the subset to be unlearned, and $X_t \setminus A_t$ is the retained subset.

Underlying flow-matching loss. Let $\ell(\theta; x)$ denote the standard flow-matching training loss on a sample $x \in \mathcal{X}$, namely the same loss used in the original pre-training objective. For any finite empirical collection S , define

$$L_S(\theta) := \frac{1}{|S|} \sum_{x \in S} \ell(\theta; x). \quad (23)$$

In particular,

$$L_{X_t \setminus A_t}(\theta) = L_{\mathcal{D}_t^+}(\theta), \quad L_{A_t}(\theta) = L_{\mathcal{D}_t^-}(\theta).$$

Reduction to the deletion objective. A first natural objective is to train the model as if the rejected samples had been removed from the replay set, namely

$$L_t^{\text{del}}(\theta) := L_{X_t \setminus A_t}(\theta) = L_{\mathcal{D}_t^+}(\theta). \quad (24)$$

Writing

$$m_t := |\mathcal{D}_t^+|, \quad k_t := |\mathcal{D}_t^-|, \quad n_t := m_t + k_t,$$

the same algebra used in (Alberti et al., 2025) gives the exact decomposition

$$L_t^{\text{del}}(\theta) = \frac{n_t}{m_t} L_{X_t}(\theta) - \frac{k_t}{m_t} L_{A_t}(\theta) = \frac{n_t}{m_t} L_{X_t}(\theta) - \frac{k_t}{m_t} L_{\mathcal{D}_t^-}(\theta). \quad (25)$$

Thus, even the valid-only objective can be rewritten using *both* accepted and rejected samples. This is precisely the same deletion identity underlying SISS in (Alberti et al., 2025), specialized here to the verifier-induced partition of queried synthetic data. In the terminology of (Alberti et al., 2025), our practical implementation uses the corresponding *non-importance-sampled* variant, i.e., the analogue of SISS (No IS), rather than the one-pass importance-sampled estimator.

Why we use a signed objective. For out-of-distribution generable-set expansion, Eq. (24) is often too conservative: it treats rejected samples only as points to be discarded. In contrast, rejected verifier queries contain useful information about directions in design space toward which the current model should allocate *less* probability mass. This motivates the signed replay objective used in the main text:

$$L_t^{\text{signed}}(\theta) := L_{\mathcal{D}_t^+}(\theta) - \alpha_t L_{\mathcal{D}_t^-}(\theta). \quad (26)$$

Its gradient is exactly the signed update:

$$\nabla_{\theta} L_t^{\text{signed}}(\theta_t) = \nabla_{\theta} L_{\mathcal{D}_t^+}(\theta_t) - \alpha_t \nabla_{\theta} L_{\mathcal{D}_t^-}(\theta_t). \quad (27)$$

At the minibatch level, our implementation follows the *non-importance-sampled* variant: we draw separate minibatches

$$U_t^+ \subseteq \mathcal{D}_t^+, \quad U_t^- \subseteq \mathcal{D}_t^-,$$

and form the empirical losses

$$\widehat{L}_t^+(\theta) = \frac{1}{|U_t^+|} \sum_{x \in U_t^+} \ell(\theta; x), \quad \widehat{L}_t^-(\theta) = \frac{1}{|U_t^-|} \sum_{x \in U_t^-} \ell(\theta; x).$$

The practical signed update is then

$$g_t = \nabla \widehat{L}_t^+(\theta_t) - \alpha_t \nabla \widehat{L}_t^-(\theta_t).$$

Hence accepted samples provide an attractive update, while rejected samples provide a repulsive update of controlled magnitude. This corresponds to the analogue of SISS (No IS) in (Alberti et al., 2025): two separate replay-buffer minibatches are used, rather than a single importance-sampled mixture minibatch.

Relation to weighted data-unlearning objectives. The weighted SISS objective of (Alberti et al., 2025) can be written as

$$L^{\text{wSISS}}(\theta) = L_{X \setminus A}(\theta) - s \frac{|A|}{|X \setminus A|} L_A(\theta). \quad (28)$$

Applying this to the verifier-induced partition

$$X = X_t = \mathcal{D}_t^+ \uplus \mathcal{D}_t^-, \quad A = A_t = \mathcal{D}_t^-,$$

yields

$$L_t^{\text{wSISS}}(\theta) = L_{\mathcal{D}_t^+}(\theta) - s_t \frac{k_t}{m_t} L_{\mathcal{D}_t^-}(\theta). \quad (29)$$

Therefore our signed replay objective in Eq. (26) is exactly the same class of weighted retain-minus-forget objective, with the identification

$$\alpha_t = s_t \frac{k_t}{m_t}. \quad (30)$$

In this sense, the update used by ACTFLOW is a direct reduction of the weighted unlearning objective of (Alberti et al., 2025) to the online verifier-guided expansion setting considered here.

Possible importance sampling (IS) extension. As mentioned, our implementation is the direct analogue of SISS (No IS) in (Alberti et al., 2025), rather than the one-pass importance-sampled estimator introduced there. An importance-sampled extension could also be used in our setting by sampling from a suitable mixture over accepted and rejected replay samples and reweighting by the corresponding importance ratios, but we avoid this additional estimator complexity here for simplicity.

Interpretation in our setting. The reduction above is exact at the level of the empirical replay buffers held fixed at round t . Conditionally on the current buffers $(\mathcal{D}_t^+, \mathcal{D}_t^-)$, the signed objective

$$L_t^{\text{signed}}(\theta) = L_{\mathcal{D}_t^+}(\theta) - \alpha_t L_{\mathcal{D}_t^-}(\theta)$$

should be read as follows:

1. the first term increases likelihood of synthetic samples that have been certified as valid by the verifier and therefore support expansion into new valid regions;
2. the second term decreases likelihood of queried samples that were explicitly rejected by the verifier and therefore encode evidence against those directions of expansion.

Thus, unlike standard valid-only pre-training, the signed replay update uses *both* outcomes of the verifier query and converts the online generation-and-verification loop into a principled retain/forget training signal.

1045 **Practical choice of α_t .** In the practical implementation, the loss-level weight α_t is calibrated online as a gradient-norm
1046 fraction. Concretely, for a target ratio $\alpha \in [0, 1]$, we set

$$\alpha_t := \alpha \frac{\|\nabla \widehat{L}_t^+(\theta_t)\|_2}{\|\nabla \widehat{L}_t^-(\theta_t)\|_2},$$

1050 This yields

$$\|\alpha_t \nabla \widehat{L}_t^-(\theta_t)\|_2 \approx \alpha \|\nabla \widehat{L}_t^+(\theta_t)\|_2,$$

1051 so the rejected-sample contribution is scaled to have a prescribed norm fraction relative to the accepted-sample contribution.

1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

E. Intermezzo: Theory of Active Safe Logistic Regression

Suppose that given any query point x , the verifier provides binary feedback $y \in \{0, 1\}$ generated according to

$$\Pr[y = 1 \mid x] = s(f(x)), \quad (31)$$

where $s : \mathbb{R} \rightarrow [0, 1]$ is the sigmoid function, and $f : \mathcal{X} \rightarrow \mathbb{R}$ is a validity function. We assume that f lies in the RKHS \mathcal{H}_k associated with a kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, with bounded norm $\|f\|_k \leq B$. We also assume that f is L_f -smooth with respect to a metric d , namely

$$|f(x) - f(x')| \leq L_f d(x, x') \quad \forall x, x' \in \mathcal{X}.$$

Without loss of generality, we further assume that $k(x, x') \leq 1$ for all $x, x' \in \mathcal{X}$.

We adopt a *pre-query* indexing convention. At the beginning of round t , the learner has access to the history

$$H_{t-1} \triangleq \{(x_\tau, y_\tau)\}_{\tau=1}^{t-1},$$

fits a probabilistic model using H_{t-1} , constructs a safe set S_t , and then selects the next query x_t . Given H_{t-1} , the learner estimates f by minimizing the regularized negative log-likelihood

$$\begin{aligned} \mu_t &\triangleq \arg \min_{g \in \mathcal{H}_k, \|g\|_k \leq B} \mathcal{L}(g, H_{t-1}), \\ \mathcal{L}(g, H_{t-1}) &\triangleq \sum_{\tau=1}^{t-1} \left[-y_\tau \log(s(g(x_\tau))) - (1 - y_\tau) \log(1 - s(g(x_\tau))) \right] + \frac{\lambda}{2} \|g\|_k^2, \end{aligned} \quad (32)$$

where $\lambda > 0$ is a regularization coefficient. By the Representer Theorem (Schölkopf et al., 2001), the solution lies in the span of the kernel sections at the previously queried points.

Given μ_t , the learner predicts the verifier output via $s(\mu_t(x))$. Previous work gives anytime-valid confidence sets of the form

$$[s(\mu_t(x)) \pm \beta_t(\delta) \sigma_t(x)]$$

for kernelized logistic regression (Pásztor et al., 2024). Under our indexing convention, the uncertainty scores are

$$\sigma_t^2(x) = k(x, x) - k_{t-1}^\top(x) (K_{t-1} + \lambda \kappa I_{t-1})^{-1} k_{t-1}(x),$$

where

$$\kappa \triangleq \sup_{a \leq B} \frac{1}{\dot{s}(a)},$$

$k_{t-1}(x)$ is the kernel vector with entries $[k_{t-1}(x)]_i = k(x_i, x)$, and K_{t-1} is the kernel matrix with entries $[K_{t-1}]_{i,j} = k(x_i, x_j)$.

We will use the following pre-query reindexing of the confidence-sequence guarantee.

Theorem E.1 (Kernelized Logistic Confidence Sequences). *Assume $f \in \mathcal{H}_k$ and $\|f\|_k \leq B$. Assume that the data x_1, \dots, x_t used to fit the model μ_t lie in a compact subset $A \subseteq \mathcal{X}$. Let $0 < \delta < 1$ and define*

$$\beta_t(\delta) \triangleq 4L_s B + 2L_s \sqrt{\frac{2\kappa}{\lambda} \left(\gamma_{t-1}^A + \log(1/\delta) \right)}, \quad (33)$$

where

$$\gamma_t^A \triangleq \max_{x_1, \dots, x_t \in A} \frac{1}{2} \log \det \left(I_t + (\lambda \kappa)^{-1} K_t \right), \quad L_s \triangleq \sup_{a \leq B} \dot{s}(a).$$

Then

$$\Pr(\forall t \geq 1, \forall x \in A : |s(\mu_t(x)) - s(f(x))| \leq \beta_t(\delta) \sigma_t(x)) \geq 1 - \delta.$$

Using these confidence sets, the learner may perform active safe logistic regression, querying only points that are certified to satisfy $s(f(x)) \geq h$ for some threshold $h \in [0, 1]$. This parallels safe Bayesian optimization in the regression setting (e.g., Sui et al., 2015; 2018).

In this setting, safe designs are those for which the verifier returns label 1 with sufficiently high probability. Suppose the learner is given a safe seed set $S_0 \subseteq \mathcal{X}$. At round t , the learner updates its certified safe set as

$$S_t = S_{t-1} \cup \{x \in \mathcal{X} \mid \exists x' \in S_{t-1} : s(\mu_t(x')) - \beta_t(\delta)\sigma_t(x') - L_s L_f d(x, x') \geq h\}. \quad (34)$$

The next lemma shows that, conditioned on the confidence event of Theorem E.1, every point ever added to S_t is indeed safe.

Lemma E.2 (Safety). *Condition on the event that the confidence intervals of Theorem E.1 are valid. Then for every $t \geq 0$, every $x \in S_t$ satisfies $s(f(x)) \geq h$.*

Proof. We argue by induction on t . The claim holds at $t = 0$ by assumption on the seed set S_0 . Now fix $t \geq 1$, and let $x \in S_t$. Either $x \in S_{t-1}$, in which case the claim follows by induction, or

$$x \in \{x \in \mathcal{X} \mid \exists x' \in S_{t-1} : s(\mu_t(x')) - \beta_t(\delta)\sigma_t(x') - L_s L_f d(x, x') \geq h\}.$$

In the latter case, there exists $x' \in S_{t-1}$ such that

$$s(\mu_t(x')) - \beta_t(\delta)\sigma_t(x') - L_s L_f d(x, x') \geq h.$$

By the confidence event,

$$s(f(x')) \geq s(\mu_t(x')) - \beta_t(\delta)\sigma_t(x').$$

Moreover, since s is L_s -Lipschitz and f is L_f -Lipschitz,

$$s(f(x)) \geq s(f(x')) - L_s L_f d(x, x').$$

Combining the two inequalities gives

$$s(f(x)) \geq s(\mu_t(x')) - \beta_t(\delta)\sigma_t(x') - L_s L_f d(x, x') \geq h.$$

Thus every $x \in S_t$ is safe. \square

Notice that by construction, the sequence $\{S_t\}_{t \geq 0}$ is monotone:

$$S_{t-1} \subseteq S_t \quad \forall t \geq 1.$$

Using this set of safe decisions, the learner may follow an active learning procedure which samples the next point at which to query the verifier as the one in the set of safe decisions that has the highest uncertainty score, often referred to as *safe uncertainty sampling*:

$$x_t = \arg \max_{x \in S_t} \sigma_t(x). \quad (35)$$

After querying the verifier at x_t and observing y_t , the history updates as

$$H_t = H_{t-1} \cup \{(x_t, y_t)\}.$$

We ultimately want to show some notion of suitable expansion of the safe set by following the above sampling procedure. As our safety set at any time is defined as an expansion of the safe set at the previous time, we cannot guarantee convergence to the entire subset of \mathcal{X} which is safe. Instead, we must define some notion of *reachable safe set* of decisions. To this end, we consider the tightened one step reachability operator defined as

$$R_\varepsilon(S) \triangleq \{x \in \mathcal{X} \mid \exists x' \in S : s(f(x')) - L_s L_f d(x, x') - \varepsilon \geq h\}.$$

The H -fold recursive application of this operator is denoted by $R_\varepsilon^H(S)$, and its closure as $H \rightarrow \infty$ is denoted $\bar{R}_\varepsilon(S)$. Our goal is to show that safe uncertainty sampling expands the certified safe set toward $\bar{R}_0(S_0)$. We first prove an uncertainty-reduction bound over the current safe set. Let

$$\Omega^* \triangleq \{x \in \mathcal{X} : s(f(x)) \geq h\}$$

denote the subset of \mathcal{X} corresponding to the *valid design space*, and assume that Ω^* is compact.

Lemma E.3 (Uncertainty Reduction). Fix $t \geq 0$ and $T \geq 1$. Under safe uncertainty sampling (35), the epistemic uncertainty at time $t + T$ decays as

$$\max_{x \in S_t} \sigma_{t+T}(x) \leq \sqrt{\frac{2\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\frac{\gamma_T^{\Omega^*}}{T}},$$

where $\bar{\sigma} \triangleq \max_{x \in \mathcal{X}} \sigma_1(x)$.

Proof. By monotonicity of the uncertainty score, it holds that

$$\begin{aligned} \max_{x \in S_t} \sigma_{t+T}(x) &\leq \frac{1}{T} \sum_{n=0}^{T-1} \max_{x \in S_t} \sigma_{t+n}(x) && \text{(uncertainty is monotone)} \\ &\leq \frac{1}{T} \sum_{n=0}^{T-1} \max_{x \in S_{t+n}} \sigma_{t+n}(x) && (S_t \text{ are monotone}) \\ &= \frac{1}{T} \sum_{n=0}^{T-1} \sigma_{t+n}(x_{t+n}) && \text{(safe uncertainty sampling (35))} \\ &\leq \frac{1}{\sqrt{T}} \sqrt{\sum_{n=0}^{T-1} \sigma_{t+n}^2(x_{t+n})} && \text{(Cauchy-Schwarz).} \end{aligned}$$

It holds that $x \leq \frac{\bar{x}}{\log(1+\bar{x})} \log(1+x)$ for $0 \leq x \leq \bar{x}$. Let $K_{t:t+T-1}$ be the Gram matrix of the queried points x_t, \dots, x_{t+T-1} , let $K_{1:t-1}$ be the Gram matrix of the previously queried points x_1, \dots, x_{t-1} , let $K_{t:t+T-1, 1:t-1}$ be the corresponding cross-Gram matrix, and define

$$\Sigma_{t:t+T-1|1:t-1} \triangleq K_{t:t+T-1} - K_{t:t+T-1, 1:t-1} (K_{1:t-1} + \lambda\kappa I_{t-1})^{-1} K_{1:t-1, t:t+T-1}.$$

Then the above inequality implies that

$$\begin{aligned} \max_{x \in S_t} \sigma_{t+T}(x) &\leq \frac{1}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\sum_{n=0}^{T-1} \log\left(1 + \frac{\sigma_{t+n}^2(x_{t+n})}{\lambda\kappa}\right)} \\ &= \frac{1}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\log \det(I_T + (\lambda\kappa)^{-1} \Sigma_{t:t+T-1|1:t-1})} \\ &\leq \frac{1}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\log \det(I_T + (\lambda\kappa)^{-1} K_{t:t+T-1})} \\ &\leq \frac{1}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\max_{x_1, \dots, x_T \in \Omega^*} \log \det(I_T + (\lambda\kappa)^{-1} K_T)} \\ &= \sqrt{2} \frac{\sqrt{\gamma_T^{\Omega^*}}}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}}. \end{aligned}$$

□

The next lemma shows that once the uncertainty over the current safe set is sufficiently small, additional safe uncertainty sampling expands the safe set according to the one-step reachability operator.

Lemma E.4 (One-step Reachable Expansion). Consider $\varepsilon > 0$, $t \geq 0$, and $T \geq 1$. Suppose that

$$T \geq \frac{8\beta_{t+T}^2 \gamma_T^{\Omega^*} (\lambda\kappa + \bar{\sigma}^2)}{\varepsilon^2}.$$

Then, conditioned on the event that the intervals of Theorem E.1 are valid, it holds that

$$R_\varepsilon(S_t) \subseteq S_{t+T}.$$

1265 *Proof.* By Lemma E.3 and the condition on T ,

$$1266 \max_{x \in S_t} \sigma_{t+T}(x) \leq \frac{\varepsilon}{2\beta_{t+T}}. \quad (36)$$

1267 Consider any point $x \in R_\varepsilon(S_t)$. By definition, there exists $x' \in S_t$ such that

$$1271 h \leq s(f(x')) - \varepsilon - L_s L_f d(x, x').$$

1272 By the confidence event at round $t + T$,

$$1273 s(f(x')) \leq s(\mu_{t+T}(x')) + \beta_{t+T} \sigma_{t+T}(x').$$

1274 Hence

$$\begin{aligned} 1275 h &\leq s(f(x')) - \varepsilon - L_s L_f d(x, x') \\ 1276 &\leq s(\mu_{t+T}(x')) + \beta_{t+T} \sigma_{t+T}(x') - \varepsilon - L_s L_f d(x, x') \\ 1277 &\leq s(\mu_{t+T}(x')) - \beta_{t+T} \sigma_{t+T}(x') - L_s L_f d(x, x'), \end{aligned}$$

1278 where the last step uses (36), since $x' \in S_t$ implies

$$1282 \beta_{t+T} \sigma_{t+T}(x') \leq \frac{\varepsilon}{2}.$$

1283 By monotonicity of the safe sets, $x' \in S_t \subseteq S_{t+T-1}$. Therefore there exists $x' \in S_{t+T-1}$ such that

$$1284 s(\mu_{t+T}(x')) - \beta_{t+T} \sigma_{t+T}(x') - L_s L_f d(x, x') \geq h.$$

1285 By the definition of S_{t+T} in (34), this implies $x \in S_{t+T}$. □

1286 Applying the above result recursively leads to H -step expansion of S_0 at time T^* . In particular, suppose $T^* = HT$ for some T satisfying

$$1287 T \geq \frac{8\beta_{HT}^2 \gamma_T^{\Omega^*} (\lambda\kappa + \bar{\sigma}^2)}{\varepsilon^2}.$$

1288 This implies that the condition of the above lemma holds for each interval of length T , since β_t is monotonically increasing. Consequently, it holds that

$$1289 R_\varepsilon^H(S_0) \subseteq S_{T^*}.$$

1290 We can find T^* large enough to satisfy this condition as long as the complexity term $\beta_t^2 \gamma_t^{\Omega^*}$ grows sublinearly.

F. Theoretical Analysis of Diffusion Models Active Diffusion Expansion

In the following, we first (i) report the general analysis for safe logistic regression that we introduced and presented in Apx. E, where we now re-interpret safety as validity, then we (ii) introduce a probabilistic modeling framework of the generative density p_t via energy-based models, and (iii) derive coverage guarantees with sample complexity for the proposed active expansion algorithm.

F.1. Probabilistic Modeling of Binary Verifier over Generative Model Learned Representation

We denote by $\phi : \mathcal{X} \rightarrow \mathcal{Z}$ the representation map learned by the pre-trained generative model. We now describe the active learning process of the verifier over the learned representation space \mathcal{Z} . Suppose that given any query point x with latent representation $z := \phi(x)$, the verifier provides binary feedback $y \in \{0, 1\}$ generated according to

$$\Pr[y = 1 \mid z] = s(g(z)), \quad (37)$$

where $s : \mathbb{R} \rightarrow [0, 1]$ is the sigmoid function, and $g : \mathcal{Z} \rightarrow \mathbb{R}$ is an unknown validity function. We assume that g lies in the RKHS \mathcal{H}_k associated with a kernel $k : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$, with bounded norm $\|g\|_k \leq B$ and that g is L_g -smooth with respect to a metric d , namely

$$|g(z) - g(z')| \leq L_g d(z, z') \quad \forall z, z' \in \mathcal{Z},$$

for B and L_g known. We additionally assume there exists a known constant \bar{Z} such that $\int_{\mathcal{Z}} \exp(g(z)) dz \leq \bar{Z}$.

On unbounded domains with respect to Lebesgue measure, the energy condition is incompatible with globally bounded kernels such as the squared-exponential kernel. Indeed, if k is bounded on the diagonal and $\|g\|_k \leq B$, then

$$|g(z)| \leq \|g\|_k \sqrt{k(z, z)} \text{ for all } z \in \mathcal{Z}$$

implies that g is uniformly bounded. Hence $\exp(g)$ is bounded below by a positive constant, and therefore cannot be integrable over an infinite-measure domain. Thus, in this setting, the energy condition should be understood as requiring a kernel class capable of representing functions with sufficiently negative tails. For example, sums of bounded kernels with even-degree polynomial kernels yield RKHS that contain coercive negative polynomials, such as $g(z) = c - z^\top A z$ with $A \succ 0$, which satisfy the energy condition.

Instead, for ease of analysis, we restrict to compact domains \mathcal{Z} , and assume without loss of generality that $k(z, z') \leq 1$ for all $z, z' \in \mathcal{Z}$.

We adopt a *pre-query* indexing convention. At the beginning of round t , the learner has access to the history

$$H_{t-1} \triangleq \{(z_\tau, y_\tau)\}_{\tau=1}^{t-1}.$$

Given H_{t-1} , the learner estimates g by minimizing the regularized negative log-likelihood

$$\begin{aligned} \mu_t &\triangleq \arg \min_{u \in \mathcal{H}_k, \|u\|_k \leq B} \mathcal{L}(u, H_{t-1}), \\ \mathcal{L}(u, H_{t-1}) &\triangleq \sum_{\tau=1}^{t-1} \left[-y_\tau \log(s(u(z_\tau))) - (1 - y_\tau) \log(1 - s(u(z_\tau))) \right] + \frac{\lambda}{2} \|u\|_k^2, \end{aligned} \quad (38)$$

where $\lambda > 0$ is a regularization coefficient. By the Representer Theorem (Schölkopf et al., 2001), the solution lies in the span of the kernel sections at the previously queried points.

Given μ_t , the learner predicts the verifier output via $s(\mu_t(z))$. Previous work gives anytime-valid confidence sets of the form

$$[s(\mu_t(z)) \pm \beta_t \sigma_t(z)]$$

for kernelized logistic regression (Pásztor et al., 2024). Under our indexing convention, the uncertainty scores are

$$\sigma_t^2(z) = k(z, z) - k_{t-1}^\top(z) (K_{t-1} + \lambda \kappa I_{t-1})^{-1} k_{t-1}(z),$$

where

$$\kappa \triangleq \sup_{a \leq B} \frac{1}{\dot{s}(a)},$$

1375 $k_{t-1}(z)$ is the kernel vector with entries $[k_{t-1}(z)]_i = k(z_i, z)$, and K_{t-1} is the kernel matrix with entries $[K_{t-1}]_{i,j} =$
 1376 $k(z_i, z_j)$.

1377 We will use the following pre-query reindexing of the confidence-sequence guarantee.

1378 **Theorem F.1** (Kernelized Logistic Confidence Sequences (Pásztor et al., 2024)). *Assume $g \in \mathcal{H}_k$ and $\|g\|_k \leq B$. Assume*
 1379 *that the queried points lie in a compact subset $A \subseteq \mathcal{Z}$. Let $0 < \delta < 1$ and define*

$$1380 \beta_t^A(\delta) \triangleq 4L_s B + 2L_s \sqrt{\frac{2\kappa}{\lambda} \left(\gamma_{t-1}^A + \log(1/\delta) \right)}, \quad (39)$$

1381 where

$$1382 \gamma_t^A \triangleq \max_{z_1, \dots, z_t \in A} \frac{1}{2} \log \det \left(I_t + (\lambda\kappa)^{-1} K_t \right), \quad L_s \triangleq \sup_{a \leq B} \dot{s}(a).$$

1383 Then

$$1384 \Pr \left(\forall t \geq 1, \forall z \in A : |s(\mu_t(z)) - s(g(z))| \leq \beta_t^A(\delta) \sigma_t(z) \right) \geq 1 - \delta.$$

1385 Using these confidence sets, the learner may perform guarded logistic regression, querying only points that are certified to
 1386 satisfy $s(g(z)) \geq h$ for some threshold $h \in [0, 1]$. This parallels safe Bayesian optimization in the regression setting (e.g.,
 1387 Sui et al., 2015; 2018).

1388 In this setting, safe designs are those for which the probabilistic verifier returns label 1 with sufficiently high probability.
 1389 Suppose the learner is given a safe seed set $S_0 \subseteq \mathcal{Z}$. Given a monotonically increasing sequence of calibration coefficients
 1390 at level δ , $\{\hat{\beta}_t(\delta)\}_{t \geq 0}$, the set of points that the learner can verify as safe may be expanded at round t as

$$1391 S_t = S_{t-1} \cup \left\{ z \in \mathcal{Z} \mid \exists z' \in S_{t-1} : s(\mu_t(z')) - \hat{\beta}_t(\delta) \sigma_t(z') - L_s L_g d(z, z') \geq h \right\}. \quad (40)$$

1392 We differentiate between the sequence $\{\hat{\beta}_t(\delta)\}$ and $\{\hat{\beta}_t^A(\delta)\}$ as computing the latter requires knowledge of the subset A ,
 1393 whereas our downstream analysis shows that our sampling remains restricted to an unknown set. Despite this, the learner
 1394 may reasonably have an upper bound on the information capacity of the kernel restricted to the unknown subset, and thus
 1395 may reasonably have access to such upper bounds.

1396 The next lemma shows that, conditioned on the confidence event of Theorem F.1, every point ever added to S_t satisfies the
 1397 validity condition.

1398 **Lemma F.2** (Safety). *Condition on the event that the confidence intervals of Theorem F.1 are valid at level $\delta \in (0, 1)$ under*
 1399 *a set A satisfying $\Omega^* \subset A$ and suppose that the sequence $\{\hat{\beta}_t(\delta)\}_{t \geq 0}$ satisfies $\hat{\beta}_t(\delta) \geq \beta_t^A(\delta)$. Then for every $t \geq 0$, every*
 1400 *$z \in S_t$ satisfies $s(g(z)) \geq h$.*

1401 *Proof.* We argue by induction on t . The claim holds at $t = 0$ by assumption on the seed set S_0 .

1402 Now fix $t \geq 1$, and let $z \in S_t$. Either $z \in S_{t-1}$, in which case the claim follows by induction, or

$$1403 z \in \left\{ z \in \mathcal{Z} \mid \exists z' \in S_{t-1} : s(\mu_t(z')) - \hat{\beta}_t(\delta) \sigma_t(z') - L_s L_g d(z, z') \geq h \right\}.$$

1404 In the latter case, there exists $z' \in S_{t-1}$ such that

$$1405 s(\mu_t(z')) - \hat{\beta}_t(\delta) \sigma_t(z') - L_s L_g d(z, z') \geq h.$$

1406 By the confidence event,

$$1407 s(g(z')) \geq s(\mu_t(z')) - \hat{\beta}_t(\delta) \sigma_t(z').$$

1408 Moreover, since s is L_s -Lipschitz and g is L_g -Lipschitz,

$$1409 s(g(z)) \geq s(g(z')) - L_s L_g d(z, z').$$

1410 Combining the two inequalities gives

$$1411 s(g(z)) \geq s(\mu_t(z')) - \hat{\beta}_t(\delta) \sigma_t(z') - L_s L_g d(z, z') \geq h.$$

1412 Thus every $z \in S_t$ exceeds the validity threshold. □

Notice that by construction, the sequence $\{S_t\}_{t \geq 0}$ is monotone:

$$S_{t-1} \subseteq S_t \quad \forall t \geq 1.$$

F.2. Abstraction of the Generative Model as an EBM

Our primary abstraction of the generative model is that it approximately tracks the sequence of certified safe sets constructed with a sequence of calibration parameters $\{\hat{\beta}_t(\delta)\}$. We achieve this abstraction using an energy based model. Specifically, for a fixed $\gamma > 0$ and $0 < \ell \leq h$, we consider an energy function which satisfies

$$\mu'_t \in \left\{ \begin{aligned} f : \int_{\mathcal{Z}} e^{f(z)} dz &\leq \bar{Z}, \\ f(z) &\geq \text{logit}(h) \mathbf{1}(S_t) \forall z \in \mathcal{Z}, \\ f(z) &\leq \text{logit}(\ell) \text{ for all } z \text{ such that } d(z, S_t) \geq \gamma \end{aligned} \right\}. \quad (41)$$

The lower bound on f captures the behavior where the updates to our generative model maintain high density on regions that our verifier certifies as valid according to (34). The upper bound on f instead captures the behavior where the model maintains low density on points far from this validated set. The explicit bound on the partition function \bar{Z} ensures that the normalization constant does not cause the density placed on the valid region to vanish. These capture the desiderata that our generative model updates maintain large density in regions which have been verified as valid, and maintain small density in regions far away from this verified region. The generative model updates using both positive and negative samples, defined in Apx. D.1, are designed to satisfy these desiderata. For our theoretical analysis, we assume that this abstraction is valid.

Assumption F.1. Consider $\gamma > 0$ and let ℓ be such that

$$\frac{\exp(\text{logit}(\ell))}{\underline{Z}} \leq \frac{\exp(\text{logit}(h))}{\bar{Z}}. \quad (42)$$

where $\underline{Z} := \text{vol}(S_0) \exp(\text{logit}(h))$ and $\text{logit}(u) := \log\left(\frac{u}{1-u}\right)$. We assume that the generative model can be abstracted as an EBM by considering the density

$$p_t(z) = \frac{\exp(\mu'_t(z))}{Z_t}, \quad Z_t := \int_{z \in \mathcal{Z}} \exp(\mu'_t(z)) dz. \quad (43)$$

for the energy function μ'_t defined in (41).

We define the corresponding *generable set* at level τ as the high-probability set

$$\Omega_t^\tau := \{x \in \mathcal{X} : p_t(x) \geq \tau\}. \quad (44)$$

The EBM-abstraction of the generative model allows us to draw a relationship between the generable set of the EBM and the certified safe sets in which the generable set at some level τ and round t contains the set S_t .

Lemma F.3 (Generable and Verifier Set Inequality for the Recursive Safe Set). *Suppose that Assumption F.1 holds with some $\gamma > 0$. Then, for any*

$$\tau \leq \frac{\exp(\text{logit}(h))}{\bar{Z}}, \quad (45)$$

the high-probability set Ω_t^τ is a superset of S_t for the EBM model defined by Assumption F.1, namely $S_t \subseteq \Omega_t^\tau$.

Proof. By definition of the high-probability set,

$$\Omega_t^\tau = \{x \in \mathcal{X} : p_t^\pi(x) \geq \tau\} = \left\{ x \in \mathcal{X} : \frac{\exp(\mu'_t(x))}{Z_t} \geq \tau \right\} = \{x \in \mathcal{X} : \mu'_t(x) \geq \log \tau + \log Z_t\}. \quad (46)$$

Now let $z \in S_t$. By definition of $\mu'_t, \mu'_t(z) \geq \text{logit}(h)$. Therefore, a sufficient condition for $x \in \Omega_t^\tau$ is $\text{logit}(h) \geq \log \tau + \log Z_t$, or equivalently, $\tau \leq \exp(\text{logit}(h))/Z_t$.

It remains to upper bound Z_t . This holds immediately by definition of μ'_t such that it satisfies

$$Z_t = \int_{\mathcal{Z}} \exp(\mu'_t(z)) dz \leq \bar{Z}.$$

Plugging this bound in above, we find that a sufficient condition is,

$$\tau \leq \frac{\exp(\text{logit}(h))}{\bar{Z}}. \quad (47)$$

Under this condition, every $z \in S_t$ satisfies $\mu_t(z) \geq \log \tau + \log Z_t$, hence by (46) we have $z \in \Omega_t^\tau$. Therefore, $S_t \subseteq \Omega_t^\tau$, concluding the proof. \square

At the same time, we can show that the generable set lies within an extension of the valid design space. To this end, we define the inflation of the valid design space.

Definition 3. We define the inflation of the valid set by amount $\zeta \in \mathbb{R}$ in metric d as

$$\Gamma(\zeta) = \{z \in \mathcal{Z} \setminus \Omega^* : d(z, \Omega^*) \leq \zeta\}.$$

Assumption F.2. Let $\gamma > 0$ be the same as Assumption F.1. We suppose that the set $\Omega^* \cup \Gamma(\gamma)$ is compact.

Lemma F.4. Fix some $\gamma > 0$. Suppose that the sequence of calibration coefficients $\{\hat{\beta}_t(\delta)\}$ used to construct the sets S_t satisfies $\hat{\beta}_t(\delta) \geq \beta_t^{\Omega^* \cup \Gamma(\gamma)}(\delta)$. Let Assumptions F.1-F.2 hold at level γ and condition on the calibration event of Theorem F.1 with $A \leftarrow \Omega^* \cup \Gamma(\gamma)$. For $\tau \geq \frac{\exp(\text{logit}(\ell))}{\bar{Z}}$ it holds that $\Omega_t^\tau \subseteq \Omega^* \cup \Gamma(\gamma)$.

Proof. By (46) and the fact that $Z_t \geq \underline{Z}$ it holds that

$$\begin{aligned} \Omega_t^\tau &\subseteq \{z \in \mathcal{Z} : \mu'_t(z) \geq \log(\tau) + \log(\underline{Z})\} \\ &\subseteq \{z \in \mathcal{Z} : \mu'_t(z) \geq \text{logit}(\ell)\} \\ &\subseteq S_t \cup \{z \in \mathcal{Z} \setminus S_t : d(z, S_t) \leq \gamma\} \\ &\subseteq \Omega^* \cup \{z \in \mathcal{Z} \setminus \Omega^* : d(z, \Omega^*) \leq \gamma\} \\ &= \Omega^* \cup \Gamma(\gamma), \end{aligned}$$

where the final set inequality follows from Lemma F.2. \square

F.3. Active Diffusion Expansion Core Analysis

By building on the two previous subsections, we can finally analyze the active expansion process of the set Ω_t^τ . Intuitively, we proceed in two steps: first, we analyze the expansion of the verifier valid set S_t under samples obtained via generative sampling; second, we use Lemma F.3 to transfer this expansion guarantee to the generable set Ω_t^τ , which is the object we ultimately wish to study.

We adopt the same pre-query indexing convention as in Section F.1. Accordingly, at round t the learner first constructs Ω_t^τ and the verifier model (μ_t, σ_t) , and then selects the next query point. The sampling scheme used by Algorithm 1 is given by as follows.

Assumption F.3. For some fixed $\alpha \geq 1$, we suppose that the sample points queried by the learner satisfy

$$z_t \in \Omega_t^\tau \quad \text{s.t.} \quad \sigma_t(z_t) \geq \frac{1}{\alpha} \max_{z \in \Omega_t^\tau} \sigma_t(z). \quad (48)$$

Crucially, the sampling oracle in (48) requires the generative model (e.g., via inference-time techniques) to sample approximate maximizers of $\sigma_t(\cdot)$ over the current generable set Ω_t^τ .

We ultimately want to show a suitable notion of expansion of the valid set by following the above procedure. As our valid set at any time is defined as an expansion of the valid set at the previous time, we cannot guarantee convergence to the entire subset of \mathcal{Z} which is valid. Instead, we must define some notion of *reachable valid set*. To this end, we consider the tightened one-step reachability operator

$$R_\varepsilon(S) \triangleq \{z \in \mathcal{Z} \mid \exists z' \in S : s(g(z')) - L_s L_g d(z, z') - \varepsilon \geq h\}.$$

The H -fold recursive application of this operator is denoted $R_\varepsilon^H(S)$, and its closure as $H \rightarrow \infty$ is denoted $\bar{R}_\varepsilon(S)$. Ultimately we show that the valid set discovered by the learner converges to $\bar{R}_0(S_0)$, leading to the following theorem.

Theorem F.5. Fix some $\gamma > 0, \delta > 0, \varepsilon > 0$. Let H be a positive integer. Suppose that the sequence of calibration coefficients $\{\hat{\beta}_t(\delta)\}$ used to construct the sets S_t satisfies $\hat{\beta}_t(\delta) \geq \beta_t^{\Omega^* \cup \Gamma(\gamma)}(\delta)$. Let Assumptions F.1-F.2 hold at level γ and condition on the calibration event of Theorem F.1 with $A \leftarrow \Omega^* \cup \Gamma(\gamma)$. Let $\frac{\exp(\text{logit}(\ell))}{\bar{Z}} \leq \tau \leq \frac{\exp(\text{logit}(h))}{\bar{Z}}$. Consider sampling with the oracle defined by Assumption F.3 for $T^* = TH$ steps, with T satisfying

$$T \geq \frac{8\alpha^2 \hat{\beta}_{HT}(\delta)^2 \gamma_T^{\Omega^* \cup \Gamma(\gamma)} (\lambda\kappa + \bar{\sigma}^2)}{\varepsilon^2}.$$

Then it holds that $R_\varepsilon^H(S_0) \subseteq \Omega_{T^*}^\tau$.

To show the above result, we first prove a few basic inequalities.

Our first inequality bounds the uncertainty over the set of valid decisions after T additional algorithm steps.

Lemma F.6 (Uncertainty Reduction via Local Generative Sampling). Consider the setting of Lemma F.4. Fix $t \geq 0$ and $T \geq 1$. Under local generative sampling (48), it holds that the epistemic uncertainty at time $t + T$ satisfies

$$\max_{z \in S_t} \sigma_{t+T}(z) \leq \alpha \sqrt{\frac{2\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\frac{\gamma_T^{\Omega^* \cup \Gamma(\gamma)}}{T}},$$

where $\bar{\sigma} \triangleq \max_{z \in \mathcal{Z}} \sigma_1(z)$.

Proof. By monotonicity of the uncertainty score, it holds that

$$\begin{aligned} \max_{z \in S_t} \sigma_{t+T}(z) &\leq \frac{1}{T} \sum_{n=0}^{T-1} \max_{z \in S_t} \sigma_{t+n}(z) && \text{(uncertainty is monotone)} \\ &\leq \frac{1}{T} \sum_{n=0}^{T-1} \max_{z \in S_{t+n}} \sigma_{t+n}(z) && (S_t \text{ are monotone}) \\ &\leq \frac{1}{T} \sum_{n=0}^{T-1} \max_{z \in \Omega_{t+n}^\tau} \sigma_{t+n}(z) && \text{(Lemma F.3)} \\ &\leq \frac{\alpha}{T} \sum_{n=0}^{T-1} \sigma_{t+n}(z_{t+n}) && \text{(local generative sampling (48))} \\ &\leq \frac{\alpha}{\sqrt{T}} \sqrt{\sum_{n=0}^{T-1} \sigma_{t+n}^2(z_{t+n})} && \text{(Cauchy-Schwarz).} \end{aligned}$$

It holds that $x \leq \frac{\bar{x}}{\log(1+\bar{x})} \log(1+x)$ for $0 \leq x \leq \bar{x}$. Let $K_{t:t+T-1}$ be the Gram matrix of the queried points z_t, \dots, z_{t+T-1} , let $K_{1:t-1}$ be the Gram matrix of the previously queried points z_1, \dots, z_{t-1} , let $K_{t:t+T-1, 1:t-1}$ be the corresponding cross-Gram matrix, and define

$$\Sigma_{t:t+T-1|1:t-1} \triangleq K_{t:t+T-1} - K_{t:t+T-1, 1:t-1} (K_{1:t-1} + \lambda\kappa I_{t-1})^{-1} K_{1:t-1, t:t+T-1}.$$

Then the above inequality implies that

$$\begin{aligned}
 \max_{z \in S_t} \sigma_{t+T}(z) &\leq \frac{\alpha}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\sum_{n=0}^{T-1} \log\left(1 + \frac{\sigma_{t+n}^2(z_{t+n})}{\lambda\kappa}\right)} \\
 &= \frac{\alpha}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\log \det(I_T + (\lambda\kappa)^{-1} \Sigma_{t:t+T-1|1:t-1})} \\
 &\leq \frac{\alpha}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\log \det(I_T + (\lambda\kappa)^{-1} K_{t:t+T-1})} \\
 &\leq \frac{\alpha}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}} \sqrt{\max_{z_1, \dots, z_T \in \Omega^* \cup \Gamma(\gamma)} \log \det(I_T + (\lambda\kappa)^{-1} K_T)} \\
 &= \alpha \sqrt{2} \frac{\sqrt{\gamma_T^{\Omega^* \cup \Gamma(\gamma)}}}{\sqrt{T}} \sqrt{\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)}}.
 \end{aligned}$$

Here the last inequality uses Lemma F.4, which implies that all queried points z_t, \dots, z_{t+T-1} lie in $\Omega^* \cup \Gamma(\gamma)$. \square

The next lemma shows that once the uncertainty over the current valid set is sufficiently small, additional local generative sampling expands the valid set according to the one-step reachability operator.

Lemma F.7 (One-step Reachable Expansion via Local Generative Sampling). *Consider the setting of Lemma F.6. Let $\varepsilon > 0$, $\delta > 0$, $t \geq 0$, and $T \geq 1$. Suppose*

$$T \geq \frac{8\alpha^2 \hat{\beta}_{t+T}(\delta)^2 \gamma_T^{\Omega^* \cup \Gamma(\gamma)} (\lambda\kappa + \bar{\sigma}^2)}{\varepsilon^2}.$$

Then it holds that

$$R_\varepsilon(S_t) \subseteq S_{t+T}.$$

Proof. By Lemma F.6 and the elementary inequality

$$\frac{\bar{\sigma}^2}{\log(1 + (\lambda\kappa)^{-1}\bar{\sigma}^2)} \leq \lambda\kappa + \bar{\sigma}^2,$$

the condition on T implies that

$$\max_{z \in S_t} \sigma_{t+T}(z) \leq \frac{\varepsilon}{2\hat{\beta}_{t+T}}, \quad (49)$$

where we have adopted the shorthand $\hat{\beta}_t$ for $\hat{\beta}_t(\delta)$. Consider any point $z \in R_\varepsilon(S_t)$. By definition, there exists $z' \in S_t$ such that

$$h \leq s(g(z')) - \varepsilon - L_s L_g d(z, z').$$

By the confidence event at round $t + T$,

$$s(g(z')) \leq s(\mu_{t+T}(z')) + \hat{\beta}_{t+T} \sigma_{t+T}(z').$$

Hence

$$\begin{aligned}
 h &\leq s(g(z')) - \varepsilon - L_s L_g d(z, z') \\
 &\leq s(\mu_{t+T}(z')) + \hat{\beta}_{t+T} \sigma_{t+T}(z') - \varepsilon - L_s L_g d(z, z') \\
 &\leq s(\mu_{t+T}(z')) - \hat{\beta}_{t+T} \sigma_{t+T}(z') - L_s L_g d(z, z') \quad (z' \in S_t, \text{ substitute (49)}).
 \end{aligned}$$

Since the sets S_t are monotone, $z' \in S_t$ implies $z' \in S_{t+T-1}$. Therefore there exists $z' \in S_{t+T-1}$ such that

$$h \leq s(\mu_{t+T}(z')) - \hat{\beta}_{t+T} \sigma_{t+T}(z') - L_s L_g d(z, z').$$

By the definition of S_{t+T} in (40), this implies $z \in S_{t+T}$. \square

Applying the above result recursively leads to H -step expansion of S_0 at some timestep T^* . In particular, suppose $T^* = HT$ for some T satisfying

$$T \geq \frac{8\alpha^2 \hat{\beta}_{HT}(\delta)^2 \gamma_T^{\Omega^* \cup \Gamma(\gamma)} (\lambda\kappa + \bar{\sigma}^2)}{\varepsilon^2}.$$

This implies that the condition of the above lemma holds for each interval of length T , since γ_t and $\hat{\beta}_t$ are monotonically increasing. Consequently, it holds that

$$R_\varepsilon^H(S_0) \subseteq S_{T^*}.$$

We can find T^* large enough to satisfy this condition as long as the complexity term $\hat{\beta}_t^2 \gamma_t^{\Omega^* \cup \Gamma(\gamma)}$ grows sublinearly. Ultimately, by Lemma F.3, we have

$$S_{T^*} \subseteq \Omega_{T^*}^\tau,$$

which implies that the generable set $\Omega_{T^*}^\tau$ obtained after running the algorithm is a superset of the valid reachable set, namely

$$R_\varepsilon^H(S_0) \subseteq S_{T^*} \subseteq \Omega_{T^*}^\tau. \quad (50)$$

While (50) states that the generable set of the extended diffusion model is a superset of the reachable valid set, it does not clarify whether enough model density has been placed within the valid reachable set $R_\varepsilon^H(S_0) \subseteq S_{T^*}$ rather than within $\Omega_{T^*}^\tau \setminus R_\varepsilon^H(S_0)$, which might contain invalid points. To shed light on this question, we next derive a pointwise lower bound on the final density over the generable set.

F.4. Deriving a Lower Bound on Validity

We now show that samples from the generative model trained after T^* are valid with high probability.

Corollary F.8. *Under the setting of Theorem F.5 and by the construction of the EBM model, it holds that*

$$\mathbb{P}_{z \sim p_{T^*}}[z \in \Omega^*] \geq \text{vol}(R_\varepsilon^H(S_0)) \frac{\exp(\text{logit}(h))}{\bar{Z}}.$$

Proof. It holds by definition of p_{T^*} that

$$\begin{aligned} \mathbb{P}_{z \sim p_{T^*}}[z \in \Omega^*] &= \int_{z \in \mathcal{Z}} \mathbf{1}(\Omega^*)(z) p_{T^*}(z) dz \\ &= \int_{z \in \mathcal{Z}} \mathbf{1}(\Omega^*)(z) \frac{\exp(\mu'_{T^*}(z))}{\int_{\mathcal{Z}} \exp(\mu'_{T^*}(z)) dz} dz \\ &\geq \int_{z \in R_\varepsilon^H(S_0)} \mathbf{1}(\Omega^*)(z) \frac{\exp(\mu'_{T^*}(z))}{\bar{Z}} dz \\ &= \int_{z \in R_\varepsilon^H(S_0)} \frac{\exp(\mu'_{T^*}(z))}{\bar{Z}} dz \\ &\geq \int_{z \in R_\varepsilon^H(S_0)} \frac{\exp(\text{logit}(h))}{\bar{Z}} dz, \end{aligned}$$

where the last equality follows from the fact that $R_\varepsilon^H(S_0) \subseteq \Omega^*$, and the last inequality follows from the fact that $R_\varepsilon^H(S_0) \subseteq S_{T^*}$ and $\mu'_{T^*}(z) \geq \text{logit}(h)$ for $z \in S_{T^*}$. \square

F.5. Proof of design-space coverage corollary

Assumption F.4 (Fixed representation with nondegenerate Jacobian). *The representation map $\phi : \mathcal{X} \rightarrow \mathcal{Z}$ is fixed throughout ADE, is a C^1 -diffeomorphism onto $\mathcal{Z} = \phi(\mathcal{X})$, and satisfies $|\det J\phi(x)| \geq j_{\min} > 0$ for all $x \in \mathcal{X}$.*

Corollary 4.2 (Design-space coverage of the induced reachable valid set). *Assume the conditions of Theorem 4.1 and Assumption F.4, with $j_{\min} := \inf_{x \in \mathcal{X}} |\det J\phi(x)| > 0$. Let $S_0^X := \phi^{-1}(S_0)$ and $\tau_X := j_{\min} \tau$. Then, after the same number T^* of verified samples as in Theorem 4.1, with probability at least $1 - \delta$,*

$$(R_\varepsilon^{X,\phi})^H(S_0^X) \subseteq \Omega_{T^*}^{X,\tau_X}. \quad (22)$$

Proof. Condition on the event of Theorem 4.1, which holds with probability at least $1 - \delta$. On this event,

$$R_\epsilon^H(S_0) \subseteq \Omega_{T^*}^\tau = \{z \in \mathcal{Z} : p'_{T^*}(z) \geq \tau\}. \quad (51)$$

We first show that for every $A \subseteq \mathcal{X}$,

$$\phi(R_\epsilon^{X,\phi}(A)) = R_\epsilon(\phi(A)). \quad (52)$$

Indeed, if $x \in R_\epsilon^{X,\phi}(A)$, then for some $x' \in A$,

$$s(g(\phi(x'))) - L_s L_g d(\phi(x), \phi(x')) - \epsilon \geq h.$$

Writing $z = \phi(x)$ and $z' = \phi(x')$, we get $z' \in \phi(A)$ and

$$s(g(z')) - L_s L_g d(z, z') - \epsilon \geq h,$$

hence $z \in R_\epsilon(\phi(A))$. This proves $\phi(R_\epsilon^{X,\phi}(A)) \subseteq R_\epsilon(\phi(A))$.

Conversely, if $z \in R_\epsilon(\phi(A))$, then for some $z' \in \phi(A)$,

$$s(g(z')) - L_s L_g d(z, z') - \epsilon \geq h.$$

Since ϕ is bijective, there exist unique $x = \phi^{-1}(z)$ and $x' = \phi^{-1}(z')$ with $x' \in A$. Substituting $z = \phi(x)$ and $z' = \phi(x')$ gives

$$s(g(\phi(x'))) - L_s L_g d(\phi(x), \phi(x')) - \epsilon \geq h,$$

so $x \in R_\epsilon^{X,\phi}(A)$ and thus $z = \phi(x) \in \phi(R_\epsilon^{X,\phi}(A))$. Therefore (52) holds.

Iterating (52) yields

$$\phi\left((R_\epsilon^{X,\phi})^H(S_0^X)\right) = R_\epsilon^H(\phi(S_0^X)) = R_\epsilon^H(S_0), \quad (53)$$

where we used $S_0^X = \phi^{-1}(S_0)$.

Now let $x \in (R_\epsilon^{X,\phi})^H(S_0^X)$. By (53), its image $z := \phi(x)$ belongs to $R_\epsilon^H(S_0)$, hence by (51),

$$p'_{T^*}(z) \geq \tau. \quad (54)$$

Since ϕ is a C^1 -diffeomorphism, the change-of-variables formula implies that the design-space and representation-space densities satisfy

$$p_1^{\pi T^*}(x) = p'_{T^*}(\phi(x)) |\det J\phi(x)|.$$

Combining this with (54) and the Jacobian lower bound from Assumption F.4 gives

$$p_1^{\pi T^*}(x) \geq \tau j_{\min} = \tau_X.$$

Therefore $x \in \Omega_{T^*}^{X,\tau_X}$ by (20). Since x was arbitrary, (22) follows. \square

F.6. Full design-space coverage under global reachability

Corollary F.9 (Full design-space coverage under global reachability). *Assume the conditions of Theorem 5.1 and Assumption 5.1. Let*

$$S_0^X := \phi^{-1}(S_0), \quad \Omega_t^{X,\tau_X} := \{x \in X : p_1^{\pi t}(x) \geq \tau_X\}, \quad \tau_X := j_{\min} \tau.$$

Assume in addition that the valid design space is reachable from the seed set in at most H steps, namely

$$\Omega^* \subseteq (R_\epsilon^{X,\phi})^H(S_0^X).$$

Then, after the same number T^ of verified samples as in Theorem 5.1, with probability at least $1 - \delta$,*

$$\Omega^* \subseteq \Omega_{T^*}^{X,\tau_X}.$$

Equivalently, every valid design belongs to the τ_X -level generable set of the final model.

1760 *Proof.* By Corollary 5.2, on an event of probability at least $1 - \delta$,

$$1761 \quad (R_\varepsilon^{X,\phi})^H(S_0^X) \subseteq \Omega_{T^*}^{X,\tau X}.$$

1764 By the additional reachability assumption,

$$1765 \quad \Omega^* \subseteq (R_\varepsilon^{X,\phi})^H(S_0^X).$$

1766 Combining the two inclusions yields

$$1767 \quad \Omega^* \subseteq \Omega_{T^*}^{X,\tau X},$$

1769 as claimed. □

1771 **Corollary F.10** (Full coverage under finite-chain global reachability). *Assume the conditions of Theorem 5.1 and Assumption 5.1. Suppose that for every $x \in \Omega^*$ there exists a finite chain*

$$1774 \quad x_0, \dots, x_m \in \Omega^*, \quad x_0 \in S_0^X, \quad x_m = x, \quad m \leq H,$$

1775 *such that for every $i = 1, \dots, m$, $L_s L_g d(\phi(x_i), \phi(x_{i-1})) + \varepsilon \leq s(g(\phi(x_{i-1}))) - h$. Then*

$$1777 \quad \Omega^* \subseteq \Omega_{T^*}^{X,\tau X}$$

1779 *with probability at least $1 - \delta$, after the same number T^* of verified samples as in Theorem 5.1.*

1781 *Proof.* Fix $x \in \Omega^*$, and let x_0, \dots, x_m be the chain given by the assumption. We prove by induction that $x_i \in (R_\varepsilon^{X,\phi})^i(S_0^X)$ for all $i \leq m$.

1784 The base case $i = 0$ is immediate since $x_0 \in S_0^X$. Assume $x_{i-1} \in (R_\varepsilon^{X,\phi})^{i-1}(S_0^X)$. Since $x_{i-1} \in \Omega^*$, we have $s(g(\phi(x_{i-1}))) \geq h$. By the chain condition,

$$1787 \quad s(g(\phi(x_{i-1}))) - L_s L_g d(\phi(x_i), \phi(x_{i-1})) - \varepsilon \geq h,$$

1789 so by definition of $R_\varepsilon^{X,\phi}$,

$$1790 \quad x_i \in R_\varepsilon^{X,\phi} \left((R_\varepsilon^{X,\phi})^{i-1}(S_0^X) \right) = (R_\varepsilon^{X,\phi})^i(S_0^X).$$

1792 Thus $x = x_m \in (R_\varepsilon^{X,\phi})^H(S_0^X)$. Since $x \in \Omega^*$ was arbitrary,

$$1795 \quad \Omega^* \subseteq (R_\varepsilon^{X,\phi})^H(S_0^X).$$

1796 The conclusion now follows from Corollary F.9. □

G. ActFlow for Discrete Diffusion

To adapt ACTFLOW for *discrete diffusion models*, we leverage the algorithm introduced in prior work for adapting a pre-trained discrete diffusion model to an intractable reward-tilted distribution (Tang et al., 2025). Given a reward function defined over clean sequences, this method provably tilts a pre-trained discrete diffusion model to the reward-tilted distribution using off-policy reinforcement learning. The full algorithm is given in Alg 2.

G.1. Discrete Diffusion as Continuous-Time Markov Chains

In contrast to the continuous state space, where generative models are defined by stochastic differential equations (SDEs) or ordinary differential equations (ODEs), generative models in the discrete state space $\mathcal{V} := \{1, \dots, V\}$ are defined by **continuous-time Markov chains** (CTMCs). A CTMC is a stochastic process $X_{0:1} := (X_s)_{s \in [0,1]}$ whose probability law is defined by a *generator* $(\mathbf{Q}_s)_{s \in [0,1]} \in \mathbb{R}^{\mathcal{V} \times \mathcal{V}}$ defined as:

$$\mathbf{Q}_s(x, y) = \lim_{\Delta s \rightarrow 0} \frac{1}{\Delta s} (\Pr(X_{s+\Delta s} = y | X_s = x) - \mathbf{1}_{x=y}) \quad (55)$$

which defines the probability of transitioning from state $x \in \mathcal{V}$ to state $y \in \mathcal{V}$ at time s .

Masked discrete diffusion models (MDMs) (Shi et al., 2024; Sahoo et al., 2024; Ou et al., 2024; Zheng et al., 2024) are a class of discrete diffusion models that aim to generate sequences $x_1 \sim p_{\text{data}} \in \mathcal{V}^L$ of length L from a sequence of absorbing *mask tokens* M . The generative process is defined as a CTMC that evolves from a prior distribution p_0 defined as the Dirac delta of fully masked sequences to the data distribution $p_1 \equiv p_{\text{data}}$. Since each position along the sequence $\ell \in \{1, \dots, L\}$ transitions from a mask token $X^\ell = M$ to a clean token $X^\ell = x^\ell$, the generator can be parameterized as:

$$\mathbf{Q}_s(x, x^{\ell \leftarrow v}) = \gamma(s) p^{\theta_t}(\cdot | x)_{\ell, v} \quad (56)$$

where $x^{\ell \leftarrow v}$ denotes the sequence where the ℓ th token is replaced with state $v \in \mathcal{V}$ and $\gamma(t)$ is the forward noising schedule. To train $p^{\theta_t}(\cdot | x)$ to reconstruct sequences from the data distribution, we can optimize a **denoising cross-entropy (DCE)** loss (Shi et al., 2024; Sahoo et al., 2024; Ou et al., 2024) defined as:

$$\mathcal{L}_{\text{DCE}}(\theta; x_1) := \mathbb{E}_{s \sim \mathcal{U}(0,1)} \left[\frac{1}{s} \mathbb{E}_{p_s(\tilde{x}_s | x_1)} \sum_{\ell: \tilde{x}_s^\ell = M} -\log p^{\theta_t}(\cdot | \tilde{x}_s^\ell)_\ell \right], \quad x_1 \sim p_{\text{data}} \quad (57)$$

where $p_s(\tilde{x}_s | x_1)$ is the distribution of partially masked sequences obtained by masking each position with probability s given a clean sequence $x_1 \sim p_{\text{data}}$.

G.2. Entropy-Regularized Uncertainty Optimization for Discrete Diffusion

Given a function $\sigma(\cdot) : \mathcal{V}^L \rightarrow \mathbb{R}$ that returns the epistemic uncertainty of a sequence $x_1 \in \mathcal{V}^L$ and a pretrained discrete diffusion model that generates the path measure \mathbb{P}^{θ_0} , we define the **uncertainty-tilted path measure** as:

$$\mathbb{P}^\sigma(X_{0:1}) := \frac{1}{Z} \mathbb{P}^{\theta_0}(X_{0:1}) \exp\left(\frac{\sigma(X_1)}{\beta}\right), \quad p_1^\sigma(X_1) \propto p_{\text{data}}(X_1) \exp\left(\frac{\sigma(X_1)}{\beta}\right) \quad (58)$$

This uncertainty-tilted path measure coincides with the solution to the **entropy-regularized reward optimization** problem with reward defined as the uncertainty $r(\cdot) := \sigma(\cdot)$ given by:

$$\arg \max_{\theta_{t+1}} \mathbb{E}_{X_{0:1} \sim \mathbb{P}^{\theta_t}} [\sigma(X_1)] - \beta D_{\text{KL}}(\mathbb{P}^{\theta_t} \| \mathbb{P}^{\theta_0}) \quad (59)$$

where \mathbb{P}^{θ_t} is the CTMC path measure generated from the adapted model with parameters θ and \mathbb{P}^{θ_0} is the frozen pre-trained model. β is the weight of the KL regularization with the pre-trained model.

G.3. Uncertainty-Aware Fine-Tuning of Discrete Diffusion

To adapt a pre-trained discrete diffusion model to align with the *uncertainty-tilted distribution* defined in (58), we leverage the off-policy reinforcement learning algorithm introduced in prior work (Tang et al., 2025), which is minimized when

Algorithm 2 ACTFLOW for Discrete Diffusion

```

1: Input: pre-trained model  $p^{\theta_0}(\cdot|X_s)$ , uncertainty function  $\sigma : \mathcal{V}^L \rightarrow \mathbb{R}$ , number of iterations  $T$ , number of WDCE
   repeats  $R$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:    $\{x_1^i, w^\sigma\}_{i=1}^B \leftarrow \text{Generate}(p^{\theta_0}, p^{\theta_t})$ 
4:    $\mathcal{B} \leftarrow \{x_1^i, w^\sigma\}_{i=1}^B$  {replay buffer}
5:   for step in  $1, \dots, N_{\text{step}}$  do
6:      $\{\tilde{x}_s^i, w^\sigma\}_{i=1}^{B \times R} \leftarrow \text{ResampleWithMask}(\mathcal{B}; R)$ 
7:     Compute  $\mathcal{L}_{\text{WDCE}}$  from (60) with  $\{\tilde{x}_s^i, w^\sigma\}_{i=1}^{B \times R}$ 
8:     Update  $\theta_{t+1}$  with  $\nabla_\theta \mathcal{L}_{\text{WDCE}}$ 
9:   end for
10: end for

```

$\mathbb{P}^{u_\theta} = \mathbb{P}^\sigma$. The training objective is defined as the **weighted denoising cross-entropy** loss $\mathcal{L}_{\text{WDCE}}$ given by:

$$\mathcal{L}_{\text{WDCE}} := \mathbb{E}_{X_{0:1} \sim \mathbb{P}^\sigma} [\mathcal{L}_{\text{DCE}}(\theta; x_1)] = \mathbb{E}_{X_{0:1} \sim \mathbb{P}^{\bar{u}}} \left[\underbrace{\frac{d\mathbb{P}^\sigma}{d\mathbb{P}^{\bar{u}}}}_{\exp(w^\sigma(X_1))} \mathcal{L}_{\text{DCE}}(\theta; x_1) \right] \quad (60)$$

where $\mathbb{P}^{\bar{\theta}_t} = \text{stopgrad}(\mathbb{P}^{\theta_t})$ is the CTMC path measure generated from the non-gradient-tracking model. We define the importance weight $w^\sigma(X_1) := \log \frac{d\mathbb{P}^\sigma}{d\mathbb{P}^{\bar{\theta}_t}}$ as the Radon-Nikodym derivative between the uncertainty-tilted path measure \mathbb{P}^σ and the non-gradient-tracking adapted model $\mathbb{P}^{\bar{\theta}_t}$:

$$\log \frac{d\mathbb{P}^\sigma}{d\mathbb{P}^{\bar{\theta}_t}}(X_{0:1}) = \frac{\sigma(X_1)}{\beta} + \underbrace{\sum_{s: X_{s+\Delta s} \neq X_s} \sum_{\ell: X_{s+\Delta s}^\ell \neq X_s^\ell} \log \frac{p^{\theta_0}(X_{s+\Delta s}^\ell | X_s)}{p^{\bar{u}}(X_{s+\Delta s}^\ell | X_s)}}_{w^\sigma(X_{0:1})} - \log Z \quad (61)$$

which reweights trajectories $X_{0:1}$ from the current frozen model by their likelihood under the uncertainty-tilted path measure. To optimize this objective, we iterate through the following steps:

- (i) Sample B trajectories $X_{0:1}^i$ from the adapted model $\mathbb{P}^{\bar{\theta}_t}$ without gradient tracking while tracking the log-likelihoods of each step $\log p^{\theta_0}(X_{s+\Delta s} | X_s) - \log p^{\bar{\theta}_t}(X_{s+\Delta s} | X_s)$.
- (ii) Compute the importance weights $w^\sigma(X_{0:1}^i)$ using the log-likelihoods and the uncertainties evaluated on the clean sequences $\sigma(X_1^i)$.
- (iii) Store the clean sequences and their weights in a replay buffer $\mathcal{B} \leftarrow \{x_1^i, w^\sigma(x_1^i)\}_{i=1}^B$.
- (iv) Resample R partially masked versions of each clean sequence in the buffer x_1^i at different time steps $s \in [0, 1]$ to obtain $\{\tilde{x}_s^i, w^\sigma(x_1^i)\}_{i=1}^{B \times R}$.
- (v) Compute $\mathcal{L}_{\text{WDCE}}(\theta)$ from Eq (60) using $\{\tilde{x}_s^i, w^\sigma(x_1^i)\}_{i=1}^{B \times R}$ and update θ with $\nabla_\theta \mathcal{L}_{\text{WDCE}}$ for N_{iter} iterations.

H. Experimental Details

H.1. Domain-agnostic evaluation metrics for real-world OOD generative modeling.

Standard evaluation criteria for deep generative models, such as Frechet Inception Distance (FID) (Heusel et al., 2017), aim to assess whether the generative model well-approximates the data distribution p_{data} . This is misaligned with OOD generative modeling: successful expansion should increase valid coverage beyond the initial generable region, and may therefore *increase* distributional distance from the pre-trained model. In low-dimensional illustrative experiments, valid coverage can be estimated directly by discretizing the design space, together with validity. In molecular and protein spaces, however, such direct coverage computation is infeasible. We therefore employ several domain-agnostic evaluation metrics for OOD generative modeling, namely: (i) number of fixed-threshold clusters covered by valid samples, to measure model coverage; (ii) Vendi diversity (Friedman & Dieng, 2022), to measure distributional diversity; (iii) FID to quantify divergence from the pre-trained model distribution, which we wish to increase; and validity percentage, to assess whether model expansion is happening over valid regions. In particular, the number of clusters corresponds to the number of hyper-spheres with a fixed data-specific radius, forming a packing over a finite draw of samples of fixed size across methods. This allows to approximately describe the volume covered by the union of hyper-spheres centered at the generated data points, and is an efficient-to-compute proxy metric to assess coverage of the model generable set.

H.2. Illustrative 2D Experiments

Overview. We evaluate ACTFLOW in a two-dimensional illustrative design space. The base model is a continuous flow model over \mathbb{R}^2 . The initial model is deliberately misspecified by centering the pretraining data at $(-1.1, 0)$, using 512 Gaussian samples with standard deviation 0.1. The valid region is a 3×3 checkerboard over $[-3.5, 3.5]^2$: a point is valid if it lies in-bounds and falls in a checkerboard cell with even parity. The base model is pretrained for 2500 steps with Adam, learning rate 10^{-3} , and batch size 256. Results are acquired over 20 seeds and 95% CIs are shown.

Algorithm configuration. We run ACTFLOW for 500 iterations and 20 random seeds. At each iteration, ACTFLOW self-generates 64 samples. Fine-tuning uses batch size 256 and 250 gradient steps per iteration. Evaluation is performed every 50 iterations using 3000 samples for evaluation curves. We use $\beta = 1/13$ and GP with RBF kernel with lengthscale 0.08, and negative-gradient scale $\alpha_t = 0.005$.

Uncertainty estimation. We use a Gaussian process with RBF kernel lengthscale set to 0.08 and employ flow representation timestep $s = 0.9$

Validity estimation (i.e., verifier). The verifier is the deterministic checkerboard validity function. The domain $[-3.5, 3.5]^2$ is partitioned into 3×3 equal cells. A sample is valid if it lies inside the domain and its cell index (i, j) satisfies $(i + j) \bmod 2 = 0$.

Coverage metric. Coverage is measured as generable coverage over the valid region. We draw 3,000 samples from the current model and estimate its generable set on a 100×100 histogram over $[-3.5, 3.5]^2$. A bin is considered generable if its estimated density is at least $\tau = 0.01$.

Ablations: no negative gradient in flow matching loss. We report in Fig. 7 visual results for the illustrative experiments, run with same parameters as in 2, except for α_t , which is now set to 0.0. As one can notice from Fig. 7, the change in this parameters leads to seemingly minimal changes in the results according to the metrics assessed.

Baselines. Both baselines use the same number of iterations, samples per iteration, fine-tuning steps, fine-tuning batch size, evaluation schedule, initial invalid model setting, and seeds as ACTFLOW.

Hardware and Compute. Each 2D experiment job requests one RTX 2080 GPU, 16GB memory per CPU, and a four-hour wall-clock limit.

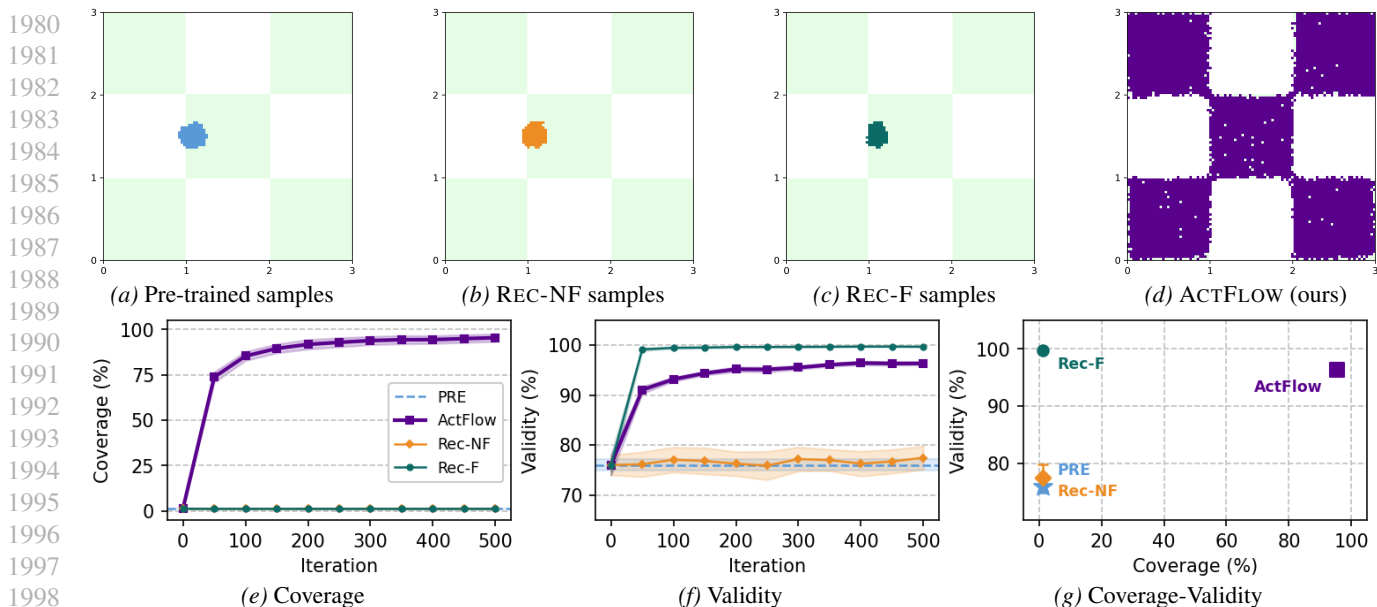


Figure 7. Results for illustrative experiments with same parameters as in 2, except for $\alpha_t = 0$.

H.3. Molecular Design: QM9 Experiments

Overview. We apply our method to small-molecule generation using FlowMol Gaussian (Dunn & Koes, 2024) pre-trained on QM9 (Ramakrishnan et al., 2014). Results are acquired over 5 seeds and 95% CIs are shown.

Algorithm configuration. We run ACTFLOW for 1,066 iterations, with 66 initial warm-up iterations during which the model is not fine-tuned. Each iteration consists of 64 samples, followed by 500 fine-tuning gradient steps. Fine-tuning uses AdamW with learning rate 10^{-4} , and batch size 64. Fine-tuning is deferred until 4096 valid samples have been collected, after which the accumulated buffer is used jointly with new guided samples at each iteration. We employ $\beta = 1/10 = 0.1$, and $s = 0.9$.

Uncertainty estimation. We use a deep bootstrapped ensemble of 5 MLPs, each with two hidden layers of 100 units, ReLU activations, and 10% dropout. Each ensemble member is trained independently on a 90% bootstrap subsample of the accumulated feature label pairs, using Adam with learning rate 10^{-3} , for up to 1000 steps. The ensemble standard deviation across members is used as the uncertainty signal.

Validity estimation (i.e., verifier). A generated molecule is deemed valid if its RDKit-sanitised representation passes valence and bond-order checks and consists of a single connected fragment. Sanitisation is preceded by an MMFF geometry relaxation (Halgren, 1996).

Coverage metric. Coverage is measured as the number of distinct molecular clusters obtained, computed via greedy sphere exclusion on Morgan fingerprints (radius 2, 2048 bits) using Tanimoto similarity with threshold $\tau = 0.85$: a candidate is added as a new cluster centre if its Tanimoto similarity to all previously selected centres is below τ . This is applied independently to the 500 valid molecules in each evaluation batch.

Diversity metric. Vendi score (Friedman & Dieng, 2022) is computed on the same 500 valid molecules per evaluation. The kernel matrix K is the pairwise Tanimoto similarity over 2048-bit Morgan fingerprints (radius 2).

Baselines. Both baselines use the same number of iterations, samples per iteration, fine-tuning steps, fine-tuning batch size, evaluation schedule, initial invalid model setting, and seeds as ACTFLOW.

FID We report in Fig. 8 FID results over iterates for QM9.

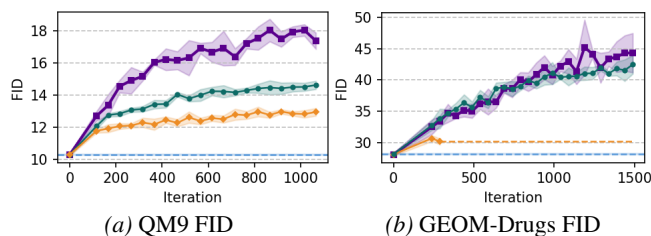


Figure 8. FID results for QM9 and GEOM-Drugs

Hardware and Compute. Each run used a single NVIDIA RTX 4090 GPU with 4 CPU cores and 32 GB of system memory, allocated for up to 120 h.

H.4. Molecular Design: GEOM-Drugs Experiments

Overview. We apply our method to small-molecule generation using FlowMol Gaussian (Dunn & Koes, 2024) pre-trained on GEOM-Drugs (Axelrod & Gomez-Bombarelli, 2022), a dataset of ~ 300 k drug-like organic molecules with energy-annotated conformers. Results are acquired over 3 seeds and 95% CIs are shown.

Algorithm configuration. We run ACTFLOW for 1,500 iterations, with 190 initial warm-up iterations during which the model is not fine-tuned. Each iteration consists of 64 samples, followed by 2000 fine-tuning gradient steps. Fine-tuning uses AdamW with learning rate 10^{-4} , and effective batch size 64. Fine-tuning is deferred until 4096 valid samples have been collected. We employ $\beta = 1/7 \approx 0.14$, and $s = 0.8$.

Uncertainty estimation. We use a deep bootstrapped ensemble of 5 MLPs, each with two hidden layers of 100 units, ReLU activations, and 10% dropout. Each ensemble member is trained independently on a 90% bootstrap subsample of the accumulated feature-label pairs, using Adam with learning rate 10^{-3} , for up to 1000 steps. The ensemble standard deviation across members is used as the uncertainty signal.

Validity estimation (i.e., verifier). A generated molecule is deemed valid if its RDKit-sanitised representation passes valence and bond-order checks and consists of a single connected fragment. Sanitisation is preceded by an MMFF geometry relaxation (Halgren, 1996).

Coverage metric. Coverage is measured as the number of distinct molecular clusters obtained, computed via greedy sphere exclusion on Morgan fingerprints (radius 2, 2048 bits) using Tanimoto similarity with threshold $\tau = 0.85$: a candidate is added as a new cluster centre if its Tanimoto similarity to all previously selected centres is below τ . This is applied independently to the 500 valid molecules in each evaluation batch.

Diversity metric. Vendi score (Friedman & Dieng, 2022) is computed on the same 500 valid molecules per evaluation. The kernel matrix K is the pairwise Tanimoto similarity over 2048-bit Morgan fingerprints (radius 2).

Baselines. Both baselines use the same number of iterations, samples per iteration, fine-tuning steps, fine-tuning batch size, evaluation schedule, initial invalid model setting, and seeds as ACTFLOW.

FID We report in Fig. 8 FID results over iterates for GEOM-Drugs.

Hardware and Compute. Each run used a single NVIDIA RTX 4090 GPU with 4 CPU cores and 128 GB of system memory, allocated for up to 120 h.

H.5. Therapeutic Peptide Design Experiments

Overview. We apply ACTFLOW for therapeutic peptide design by adapting the pre-trained discrete diffusion peptide SMILES generator from PepTune (Tang et al., 2025) trained on 11 million peptide SMILES. The Simplified Molecular-Input Line-Entry System (SMILES) (Weininger, 1988) representation enables the generation of non-natural amino acids containing diverse chemical backbone and side-chain modifications, and cyclic modifications, significantly expanding the design space

of therapeutic peptides over the standard 20 natural amino acids. The tokenization uses the SMILES Pair Encoding (SPE) tokenization scheme (Li & Fourches, 2021) from PeptideCLM (Feller & Wilke, 2025) with vocabulary size 586, including 5 special tokens. Results are acquired over 5 seeds and 95% CIs are shown.

Algorithm configuration. We run ACTFLOW for 5 rounds following Alg 2, with an initial pool of 1000 peptide sequences generated from the pretrained model and 100 iterations per round. At each round, we draw 100 sequences from the pool, score them with the verifier and Gaussian process uncertainty model, update the policy with the WDCE loss defined in Apx G.3 (16 replicates per sequence), and refresh the training pool by sampling from the updated policy. We train with Adam at learning rate 1×10^{-4} and batch size 100. The reward is scaled by a parameter $\beta = 0.005$.

Uncertainty estimation. Sequence-level uncertainty is obtained from a Gaussian process (GP) fit on top of the diffusion model’s RoFormer encoder. For each peptide, we extract attention-pooled, L_2 -normalised hidden states from the final transformer layer (768-dimensions) and treat the GP posterior variance at that point as the uncertainty signal. The GP uses an RBF kernel length-scale is initialised from the mean pairwise embedding distance of 50 samples from the pretrained model. The GP’s posterior is refit at each iteration on all sampled sequences.

Validity estimation (i.e., verifier). Generated SMILES are first parsed with RDKit. Those that yield a valid Mol object are then decoded with their SMILES2PEPTIDE verifier (Tang et al., 2025) which decodes the SMILES into a sequence of natural and non-natural amino acids split on their peptide bonds.

Coverage metric. Coverage is measured as the fraction of generated peptides that occupy distinct chemical neighbourhoods, using sphere-exclusion clustering on Morgan/ECFP4 fingerprints (radius 2, 2048 bits). We report the number of clusters at a Tanimoto-similarity threshold of 0.001.

Diversity metric. We compute the Vendi score (Friedman & Dieng, 2022) computed using the RBF kernel on the pretrained RoFormer’s L_2 -normalised hidden states.

Baselines. We compare against three baselines: (1) the pretrained model (Tang et al., 2025), (2) REC-NF where the policy is updated on its own samples with no uncertainty tilting (uniform weights on WDCE loss), (3) REC-F where the policy is updated on its own samples filtered to retain only valid peptides classified by the verifier. We hold hyperparameters fixed across all baselines.

Hardware and Compute. Each mode and seed run is trained on a single NVIDIA B200 GPU with 8 CPU cores and 80 GB of system RAM. A full 100-iteration run fits within a 48-hour wallclock budget per GPU.

H.6. Protein Sequence Design Experiments

Overview. We apply our method to protein sequence design using a continuous ESM diffusion model from SGPO (Yang et al., 2025). The base model is a continuous-space denoising network operating over ESM token-probability vectors of dimension equal to the vocabulary size (31 tokens, including the 20 standard amino acids and special tokens), pre-trained on the CreiLOV fluorescence dataset (Chen et al., 2023). CreiLOV is a 119-residue fluorescent protein; the dataset contains experimentally measured fluorescence fitness values for sequence variants. The diffusion process uses a cosine noise schedule, where $\alpha_t = \cos\left(\frac{(1-t)\pi}{2}\right)$ and $\beta_t = \sqrt{1 - \alpha_t^2}$, so that $t = 1$ corresponds to data and $t = 0$ to pure noise.

Algorithm configuration. We run ACTFLOW for 512 iterations, each consisting of 64 samples, followed by 1000 fine-tuning gradient steps on the accumulated valid samples. Fine-tuning uses AdamW with learning rate 10^{-4} , batch size 64, and no weight decay. Fine-tuning is deferred until 4096 valid samples have been collected (warm-up period), after which the accumulated buffer is used jointly with new guided samples at each iteration. We use $\beta = 1/50$. Features for the uncertainty estimator are extracted from the encoder of the ESM network at flow representation timestep $s = 0.8$, mean-pooled over the sequence dimension.

Uncertainty estimation. We use a deep ensemble of 5 MLPs, each with two hidden layers of 100 units, ReLU activations, and 10% dropout. Each ensemble member is trained independently on a 90% bootstrap subsample of the accumulated

feature-label pairs, using Adam with learning rate 10^{-3} , for up to 1000 steps. The ensemble standard deviation across members is used as the uncertainty signal.

Validity estimation (i.e., verifier). A generated sequence is deemed valid if its mean predicted local distance difference test (pLDDT), computed by ESMFold (Lin et al., 2022), exceeds a threshold of 65. ESMFold is run in batches of 32 sequences.

Coverage metric. Coverage is measured as the number of distinct sequence clusters accumulated across all 512 iterations, computed via greedy sphere exclusion on sequence identity with threshold $\tau = 0.35$: a candidate is added as a new cluster center if its sequence identity to all previously selected center is below τ .

Diversity metric. Vendi score (Friedman & Dieng, 2022) is computed on token-level ESM embeddings. For each generated sequence, we compute its embedding as the probability vector over the vocabulary projected through the ESM token embedding table (L2-normalised and scaled by $\sqrt{d_{\text{model}}}$), then mean-pooled over the sequence length. Vendi score is computed using an RBF kernel with lengthscale $\ell = 2.0$ applied to these mean-pooled embeddings.

Ablation: feature timestep. We ablated the fine-tuning flow representation time-steps $t \in \{0.5, 0.8, 0.9, 0.95\}$ in an alternative configuration with 1024 valid samples are initially queried. Among the feature timestep variants, $t = 0.95$ performed worst in terms of coverage (12 clusters at $\lambda = 50$, Vendi = 18.6), showing that extracting features close to the data level might be significantly sub-optimal.

Baselines. REC-F runs the same fine-tuning loop without uncertainty-guided sampling, training only on valid samples (pLDDT > 65). REC-NF additionally disables the validity filter during fine-tuning, training on all generated samples regardless of pLDDT. Both baselines use identical hyperparameters (fine-tuning steps, learning rate, batch size, warm-up threshold) to ACTFLOW, differing only in whether uncertainty guidance and validity filtering are applied.

Hardware and Compute. Each run used a single RTX 4090 GPU for 24 hours and 256 GB of memory.

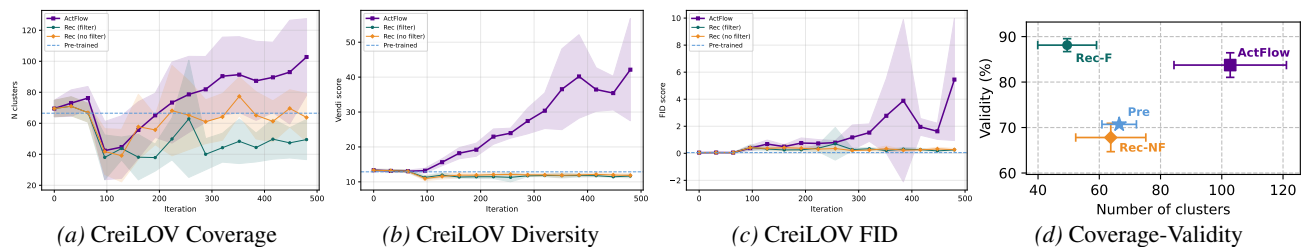


Figure 9. Results of protein sequence design experiments over iterations (Figs 9a - 9c), and diversity-validity tradeoff at final iteration (Fig 9d). ACTFLOW significantly outperforms REC-NF and REC-F in all diversity metrics (FID, Vendi, number of clusters), while maintaining a competitive validity.