Active Flow Matching

Yashvir S. Grewal^{1,2}, Daniel M. Steinberg², Edwin V. Bonilla², and Thang D. Bui¹

¹ Australian National University, Australia¹
² Data61, CSIRO, Australia²

Abstract. Discrete diffusion and flow matching excel at capturing epistatic structure in protein fitness landscapes through parallel, iterative refinement. However, their implicit nature—sampling via learned dynamics without tractable densities—prevents direct use with principled variational frameworks like VSD and CbAS for budget-constrained design. We introduce $Active\ Flow\ Matching\ (AFM)$, which reformulates variational objectives to operate on conditional endpoint distributions along the flow rather than requiring $\log q_{\phi}(x)$. This enables gradient-based steering of flow models toward high-fitness regions while preserving the rigor of VSD and CbAS. We derive forward-KL and reverse-KL variants using self-normalised importance sampling. Across four protein design tasks forward-KL AFM consistently achieves lower regret and higher optimization performance than VSD and diffusion-based LaMBO-2, demonstrating effective exploration-exploitation under tight experimental budgets.

1 Introduction

Autoregressive (AR) decoders are commonly used across domains for discrete generation tasks, but their left-to-right factorisation cannot revise early tokens. This is a fundamental mismatch for *epistatic* systems where changing position i alters the effect of changing j. Protein fitness landscapes exhibit such coupling where distant residues interact through 3D folding and binding [Starr and Thornton, 2016, Phillips, 2008]. The *fitness square* formalizes this: independence requires $F_{11} = F_{10} + F_{01} - F_{00}$, but evolution frequently yields epistasis $\varepsilon = F_{11} - F_{10} - F_{01} + F_{00} \neq 0$. Capturing ε demands joint updates across sites.

Non-autoregressive iterative refinement models such as discrete diffusion and flow matching, generate all positions in parallel, enabling global coupling [Austin et al., 2021, Gat et al., 2024]. These models match or exceed AR/masked baselines across protein and RNA design (EvoDiff, DiMA, RFdiffusion, Chroma, RNAdiffusion, DNA-Diffusion), and can also enable structure-conditioned generation (RFdiffusion, FoldFlow, motif-scaffolding) [Alamdari et al., 2023, Meshchaninov et al., 2024, Watson et al., 2023, Ingraham et al., 2023, Huang et al., 2024, DaSilva et al., 2024, Bose et al., 2024, Trippe et al., 2022].

Translating these generative capabilities into practical discoveries requires navigating finite experimental budgets. Discovery loops face combinatorial search $(20^{20} \approx 10^{26} \text{ for } 20\text{-residue peptides})$ and expensive experiments ($\sim \$500\text{-}2000$

per assay) [Biophysics and Core, 2024, Core, 2024, for Macromolecular Interactions, 2024, GenScript, 2025]. We adopt active generation view to solve this problem, where we learn $q(x \mid y > \tau)$, the conditional distribution of high-fitness designs under fixed budgets [Steinberg et al., 2025a]. Practical requirements include (i) diverse batches for parallel screening [Jain et al., 2023, Steinberg et al., 2025b], (ii) multi-objective flexibility [Stanton et al., 2022, Jain et al., 2023], and (iii) interpretable structure discovery via co-occurrence patterns in the batch [Marks et al., 2011, Hopf et al., 2014].

Two principled approaches cast active generation as variational inference over rare events. **VSD** (reverse **KL**) minimizes $\mathrm{KL}(q_{\phi}(x) \parallel p(x \mid y \geq \tau, D_t))$, yielding an ELBO with prior, likelihood, and entropy terms; discrete sequences require score-function estimators, demanding access to $\nabla_{\phi} \log q_{\phi}(x)$ [Steinberg et al., 2025a]. **CbAS** (forward **KL**) minimizes $\mathrm{KL}(p(x \mid y \geq \tau) \parallel q_{\phi}(x))$, yielding weighted MLE $\mathbb{E}_{p(x)}[w(x) \log q_{\phi}(x)]$ where $w(x) \propto \Pr(y \geq \tau \mid x)$ [Brookes et al., 2019]. Both support informative priors and batch-sequential updates, but both require tractable $q_{\phi}(x)$.

State-of-the-art discrete diffusion and flow-matching models are *implicit* generators: they optimise score/denoising or flow-regression objectives rather than a tractable likelihood, and thus do not yield normalised densities over discrete sequences. Consequently, evaluating or differentiating $\log q_{\phi}(x)$ is generally intractable. These models sample via learned dynamics but lack a usable mass function $q_{\phi}(x)$: for discrete diffusion, exact $\log q_{\phi}(x)$ requires summing over exponentially many corruption paths [Austin et al., 2021]; for discrete flow matching, current formulations provide no simple closed-form mass function [Lipman et al., 2022, Gat et al., 2024]. Objectives that require $\log q_{\phi}(x)$ or its score $\nabla_{\phi} \log q_{\phi}(x)$ are therefore incompatible with these generative models.

Active Flow Matching (AFM). We resolve this by reformulating variational objectives to operate on conditional endpoint distributions along the flow rather than on $q_{\phi}(x)$ itself. Active Flow Matching preserves the principled foundations of VSD and CBAS while leveraging implicit generators for principled, budget-efficient design.

2 Active Flow Matching (AFM)

Setup. Let $\mathcal{X} = \Sigma^L$ denote the sequence space. We train a discrete-state flow that induces, for each $t \in [0,1]$, the conditional endpoint distribution $q_t^{\theta}(\mathbf{x}_1 \mid \mathbf{x}_t)$. The flow starts from uniform $u(\mathbf{x}) = |\Sigma|^{-L}$ at t = 0. A class probability estimator provides scores $p(y=1 \mid \mathbf{x}, \mathcal{D})$ for desirable sequences, where y denotes the property label.

Forward-KL AFM If we could sample from $p(\mathbf{x}_1|y)$, we would learn the flow by simply minimising

$$\mathcal{L}_{\text{gVFM}}(\theta) = \mathbb{E}_{t,\mathbf{x}_t|y} \left[\text{KL} \left[p_t(\mathbf{x}_1|\mathbf{x}_t, y) \| q_t^{\theta}(\mathbf{x}_1|\mathbf{x}_t) \right] \right] = -\mathbb{E}_{t,\mathbf{x}_1|y,\mathbf{x}_t} \left[\log q_t^{\theta}(\mathbf{x}_1|\mathbf{x}_t) \right] + \text{const.}$$
(1)

Since sampling from $p(\mathbf{x}_1|y)$ is intractable, we use self-normalized importance sampling (SNIS) with a proposal distribution $q(\mathbf{x}_1)$:

$$\mathcal{L}_{\text{gVFM}}(\theta) = -\mathbb{E}_{t, \mathbf{x}_1 \sim q(\mathbf{x}_1), \mathbf{x}_t} \left[\frac{p(\mathbf{x}_1 | y)}{q(\mathbf{x}_1)} \log q_t^{\theta}(\mathbf{x}_1 | \mathbf{x}_t) \right]$$
(2)

$$\approx -\mathbb{E}_{t,\mathbf{x}_t} \left[\frac{\sum_{k=1}^K w_k \log q_t^{\theta}(\mathbf{x}_{1,k}|\mathbf{x}_t)}{\sum_{k=1}^K w_k} \right], \tag{3}$$

where $\{\mathbf{x}_{1,k}\}_{k=1}^K \sim q(\mathbf{x}_1)$, $w_k = \frac{p(\mathbf{x}_{1,k},y)}{q(\mathbf{x}_{1,k})}$, $t \sim \text{Unif}[0,1]$, and \mathbf{x}_t is sampled from the model's CTMC.

 $Reverse-KL\ AFM$ At each round, we steer the base flow (from the previous round) toward high-property regions by minimizing

$$\mathcal{L}_{\text{srVFM}}(\phi) = \mathbb{E}_{t,\mathbf{x}_t} \left[\text{KL} \left[q_t^{\phi}(\mathbf{x}_1|\mathbf{x}_t) || p_t(\mathbf{x}_1|\mathbf{x}_t) \right] \right], \tag{4}$$

where $p_t(\mathbf{x}_1|\mathbf{x}_t) \propto q_t^{\theta}(\mathbf{x}_1|\mathbf{x}_t)p(y|\mathbf{x}_1,\mathcal{D})$ and θ denotes the base flow parameters. Using SNIS with proposal $q(\mathbf{x}_1)$ yields:

$$\mathcal{L}_{\text{srVFM}}(\phi) \approx \mathbb{E}_{t,\tilde{\mathbf{x}}_1 \sim q(\tilde{\mathbf{x}}_1),\mathbf{x}_t} \left[\frac{p(\tilde{\mathbf{x}}_1|y)}{q(\tilde{\mathbf{x}}_1)} \mathbb{E}_{q_t^{\phi}(\mathbf{x}_1|\mathbf{x}_t)} \left[\log q_t^{\phi}(\mathbf{x}_1|\mathbf{x}_t) - \log q_t^{\theta}(\mathbf{x}_1|\mathbf{x}_t) - \log p(y|\mathbf{x}_1, \mathcal{D}) \right] \right]$$
(5)

$$\approx \sum_{k=1}^{K} \tilde{w}_{k} \mathbb{E}_{q_{t}^{\phi}(\mathbf{x}_{1}|\mathbf{x}_{t,k})} \left[\log q_{t}^{\phi}(\mathbf{x}_{1}|\mathbf{x}_{t,k}) - \log q_{t}^{\theta}(\mathbf{x}_{1}|\mathbf{x}_{t,k}) - \log p(y|\mathbf{x}_{1}, \mathcal{D}) \right],$$
(6)

where
$$\tilde{w}_k = w_k / \sum_{k'=1}^K w_{k'}$$
, $w_k = \frac{p(\mathbf{x}_{1,k}, y)}{q(\mathbf{x}_{1,k})}$, $t \sim \mathcal{U}(0,1)$, $\mathbf{x}_{1,k} \sim q(\mathbf{x}_1)$.

Our choice of proposal distribution (a mixture of unlabelled data, flow endpoints, and a replay buffer) provides good coverage as demonstrated in the strong experiment results.

Symmetric-KL baseline. For completeness, we also report a symmetric-KL variant that adds the forward- and reverse-KL objectives above, implemented with the same SNIS endpoint sampling scheme.

3 Experiments

We evaluate on four protein design tasks: Ehrlich synthetic objectives (lengths 32, 64) [Stanton et al., 2024], FoldX stability [Guerois et al., 2002, Schymkowitz et al., 2005, Delgado et al., 2019], SASA optimization [Lee and Richards, 1971, Shrake and Rupley, 1973, pol]. We compare Forward-KL AFM against VSD and diffusion-based LaMBO-2. For tasks with known optima (Ehrlich), we report simple regret $r_t = f(x^*) - \max_{1 \le s \le t} f(x_s)$. For unknown optima (FoldX, SASA),

Y. S. Grewal et al.

4

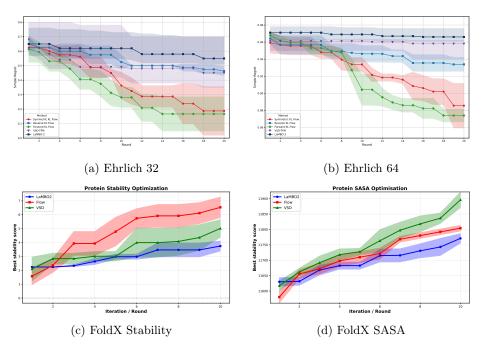


Fig. 1: Performance comparison across four protein design tasks. Forward-KL AFM achieves superior optimization compared to VSD and Lambo2 baselines.

we report highest value uptil each round $I_t = \max_{1 \le s \le t} f(x_s)$. We use a batch size of 128 in ehrlich sequences and batch size of 10 in FoldX experiments

Forward-KL AFM achieves lowest regret (Ehrlich) and highest scores in FoldX stability. (Figure 1;). Reverse-KL's performs relatively poorly. Symmetric-KL performs competitively but trails Forward-KL on Ehrlich-32/64, indicating Forward-KL's mass-covering better balances exploration-exploitation in these sequence spaces.

4 Conclusion

We introduced Active Flow Matching (AFM), which enables principled variational optimization with implicit discrete generators by reformulating objectives on conditional endpoint distributions. This resolves the incompatibility between state-of-the-art flow models and likelihood-based frameworks like VSD and CbAS, allowing gradient-based steering without tractable $q_{\phi}(x)$. Across protein design tasks, forward-KL AFM consistently outperforms existing methods under tight experimental budgets, demonstrating effective exploration-exploitation. Our framework opens the door to leveraging powerful pretrained flow and diffusion models (e.g., EvoDiff, ESM-2) for budget-constrained discovery in proteins.

Bibliography

- Protein solvent accessibility (using foldx) poli objective. https://machinelearninglifescience.github.io/polidocs/using_poli/objective_repository/foldx_sasa.html. AccessedOct29, 2025.
- Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Neil Tenenholtz, Robert Strome, Alan M. Moses, Alex X. Lu, Nicolo Fusi, Ava P. Amini, and Kevin K. Yang. Protein generation with evolutionary diffusion: sequence is all you need. bioRxiv, 2023. https://doi.org/10.1101/2023.09.11.556673. URL https://www.biorxiv.org/content/10.1101/2023.09.11.556673v1. Preprint; Microsoft Research page lists Nov 2024.
- Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state spaces. arXiv, 2021. URL https://arxiv.org/abs/2107.03006.
- Scripps Research Biophysics and Biochemistry Core. Surface plasmon resonance (spr) fees, 2024. URL https://www.scripps.edu/science-and-medicine/cores-and-services/biophysics-and-biochemistry-core/spr/index.html. Biacore S200 full-service \$500 setup + \$260/sample (includes 3 replicates).
- Avishek Joey Bose, Alexander Tong, Guillaume Huguet, James Vuckovic, Kilian Fatras, Eric Thibodeau-Laufer, Pablo Lemos, Riashat Islam, Cheng-Hao Liu, Jarrid Rector-Brooks, Tara Akhound-Sadegh, and Michael Bronstein. Se(3)-stochastic flow matching for protein backbone generation (foldflow). In *Proc. International Conference on Learning Representations (ICLR)*, 2024. URL https://arxiv.org/abs/2310.02391. ICLR 2024; arXiv:2310.02391.
- David H. Brookes, Hahnbeom Park, and Jennifer Listgarten. Conditioning by adaptive sampling for robust design. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 773–782. PMLR, 2019. URL https://proceedings.mlr.press/v97/brookes19a.html.
- Duke University Biomolecular Interaction Analysis Core. Reservations, policies, and rates, 2024. URL https://dhvi.duke.edu/programs-and-centers/shared-resources/cores/biomolecular-interaction-analysis-bia/services-and-2. SPR/BLI hourly rates; sample prep/data analysis fees.
- Lucas Ferreira DaSilva, Simon Senan, Zain Munir Patel, Aniketh Janardhan Reddy, Sameer Gabbita, Zach Nussbaum, César Miguel Valdez Córdova, Aaron Wenteler, Noah Weber, Tin M. Tunjic, et al. Dna-diffusion: Leveraging generative models for controlling chromatin accessibility and gene expression via synthetic regulatory elements. bioRxiv, 2024. https://doi.org/10.1101/2024.02.01.578352. URL https://www.biorxiv.org/content/10.1101/2024.02.01.578352v1.
- Jesús Delgado, René Radusky, Daniel Cianferoni, and Luis Serrano. Foldx 5.0: Working with rna, small molecules and a new graphical interface. *Bioinformatics*, 35(20):4168–4169, 2019. https://doi.org/10.1093/bioinformatics/btz185.

- Harvard Medical School Center for Macromolecular Interactions. Access fees, 2024.
 URL https://cmi.hms.harvard.edu/harvard-access-fees. SPR/BLI day/hour rates.
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky T. Q. Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. arXiv, 2024. https://doi.org/10.48550/arXiv.2407.15595. URL https://arxiv.org/abs/2407.15595.
- GenScript. Bli & spr real-time affinity measurement services, 2025. URL https://www.genscript.com/bli-spr-real-time-affinity-measurement-services.html. High-throughput BLI from 120-136/4 design; 8120-136/4 design;
- Raphaël Guerois, Jens Erik Nielsen, and Luis Serrano. Predicting changes in the stability of proteins and protein complexes: A study of more than 1000 mutations. *Journal of Molecular Biology*, 320(2):369–387, 2002. https://doi.org/10.1016/S0022-2836(02)00442-4.
- Thomas A. Hopf et al. Sequence co-evolution gives 3d contacts and structures of protein complexes. *eLife*, 3:e03430, 2014. https://doi.org/10.7554/eLife.03430.
- Kaixuan Huang, Yukang Yang, Kaidi Fu, Yanyi Chu, Le Cong, and Mengdi Wang. Latent diffusion models for controllable rna sequence generation. arXiv, 2024. URL https://arxiv.org/abs/2409.09828.
- John B. Ingraham, Max Baranov, Zak Costello, Karl W. Barber, Wujie Wang, Ahmed Ismail, Vincent Frappier, Dana M. Lord, Christopher Ng-Thow-Hing, Erik R. Van Vlack, Shan Tie, Vincent Xue, Sarah C. Cowles, Alan Leung, João V. Rodrigues, Claudio L. Morales-Perez, Alex M. Ayoub, Robin Green, Katherine Puentes, Frank Oplinger, Nishant V. Panwar, Fritz Obermeyer, Adam R. Root, Andrew L. Beam, Frank J. Poelwijk, and Gevorg Grigoryan. Illuminating protein space with a programmable generative model. *Nature*, 623 (7988):1070–1078, 2023. https://doi.org/10.1038/s41586-023-06728-8. URL https://doi.org/10.1038/s41586-023-06728-8.
- Moksh Jain. Emmanuel Bengio, Alex Hernandez-Garcia, al. et Multi-objective ICML2023. URL gflownets. In2023. https://proceedings.mlr.press/v202/jain23a.html.
- B. Lee and F. M. Richards. The interpretation of protein structures: Estimation of static accessibility. *Journal of Molecular Biology*, 55:379–400, 1971. https://doi.org/10.1016/0022-2836(71)90324-X.
- Т. Q. Chen, Yaron Lipman, Ricky Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. 2022. https://doi.org/10.48550/arXiv.2210.02747. arXiv, URL https://arxiv.org/abs/2210.02747.
- Debora S. Marks, Thomas A. Hopf, and Chris Sander. Protein 3d structure computed from evolutionary sequence variation. *PLoS ONE*, 6(12):e28766, 2011. https://doi.org/10.1371/journal.pone.0028766.
- Viacheslav Meshchaninov, Pavel Strashnov, Andrey Shevtsov, Fedor Nikolaev, Nikita Ivanisenko, Olga Kardymon, and Dmitry Vetrov. Diffusion on language model encodings for protein sequence genera-

- tion. *arXiv*, 2024. https://doi.org/10.48550/arXiv.2403.03726. URL https://arxiv.org/abs/2403.03726. Accepted to ICML 2025 (poster).
- Patrick C. Phillips. Epistasis the essential role of gene interactions in the structure and evolution of genetic systems. *Nature Reviews Genetics*, 9(11):855–867, 2008. https://doi.org/10.1038/nrg2452. URL https://doi.org/10.1038/nrg2452.
- Joost Schymkowitz, Jesper Borg, François Stricher, Robby Nys, Frédéric Rousseau, and Luis Serrano. The foldx web server: An online force field. *Nucleic Acids Research*, 33(Web Server issue):W382–W388, 2005. https://doi.org/10.1093/nar/gki387.
- A. Shrake and J. A. Rupley. Environment and exposure to solvent of protein atoms: Lysozyme and insulin. *Journal of Molecular Biology*, 79(2):351–371, 1973. https://doi.org/10.1016/0022-2836(73)90011-9.
- Samuel Stanton, Wesley J. Maddox, Nate Gruver, Phillip Maffettone, Emily Delaney, Peyton Greenside, and Andrew Gordon Wilson. Accelerating bayesian optimization for biological sequence design with denoising autoencoders. In *ICML* 2022, 2022. URL https://proceedings.mlr.press/v162/stanton22a.html.
- Samuel Stanton, Robert Alberstein, Nathan Frey, Andrew Watkins, and Kyunghyun Cho. Closed-form test functions for biophysical sequence optimization algorithms, 2024.
- Tyler N. Starr and Joseph W. Thornton. Epistasis in protein evolution. *Protein Science*, 25(7):1204–1218, 2016. https://doi.org/10.1002/pro.2897. URL https://doi.org/10.1002/pro.2897.
- Daniel M. Steinberg, Rafael Oliveira, Cheng Soon Ong, and Edwin V. Bonilla. Variational search distributions. ICLR 2025, 2025a. URL https://arxiv.org/abs/2409.06142.
- Daniel M. Steinberg, Asiri Wijesinghe, Rafael Oliveira, Piotr Koniusz, Cheng Soon Ong, and Edwin V. Bonilla. Amortized active generation of pareto sets. arXiv:2510.21052, 2025b. URL https://arxiv.org/abs/2510.21052.
- Brian L. Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and Tommi Jaakkola. Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. *arXiv*, 2022. URL https://arxiv.org/abs/2206.04119. ICLR 2023.
- Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, Helen E. Eisenach, Woody Ahern, Andrew J. Borst, Robert J. Ragotte, Lukas F. Milles, Basile I. M. Wicky, Nikita Hanikel, Samuel J. Pellock, Alexis Courbet, William Sheffler, Jue Wang, Preetham Venkatesh, Isaac Sappington, Susana Vázquez Torres, Anna Lauko, Valentin De Bortoli, Emile Mathieu, Sergey Ovchinnikov, Regina Barzilay, Tommi S. Jaakkola, Frank DiMaio, Minkyung Baek, and David Baker. De novo design of protein structure and function with rfdiffusion. Nature, 620 (7976):1089–1100, 2023. https://doi.org/10.1038/s41586-023-06415-8. URL https://doi.org/10.1038/s41586-023-06415-8.