# Gumbel-Softmax Flow Matching with Straight-Through Guidance for Controllable Biological Sequence Generation

**Sophia Tang**[1], Yinuo Zhang[2], Alexander Tong[3,4], Pranam Chatterjee[1,5] †

[1]Department of Computer and Information Science, University of Pennsylvania
[2]Center of Computational Biology, Duke-NUS Medical School
[3]Mila, Quebec AI Institute, [4]Université de Montréal
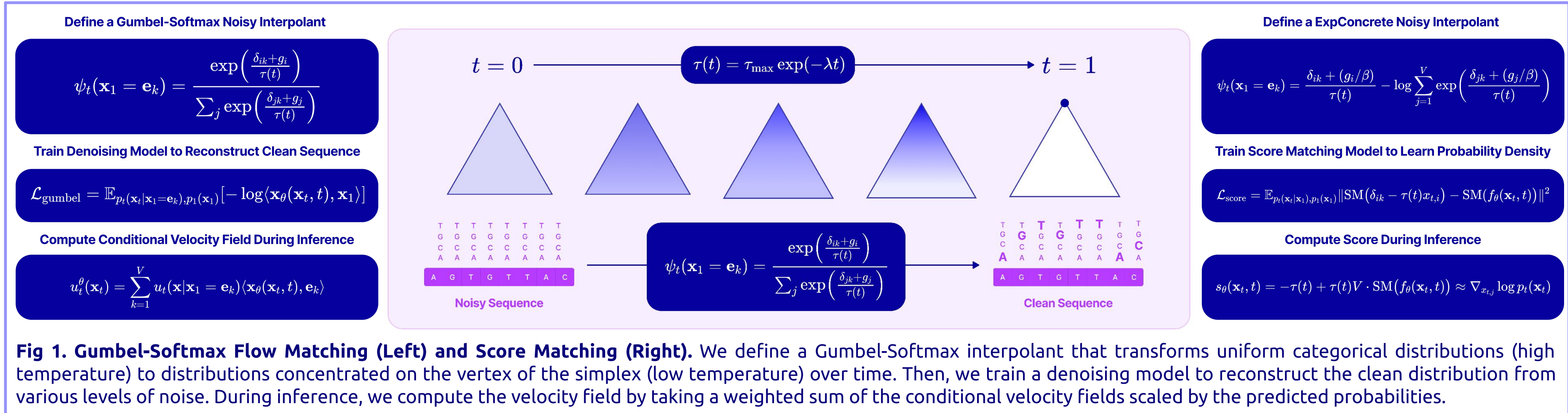[5]Department of Bioengineering, University of Pennsylvania
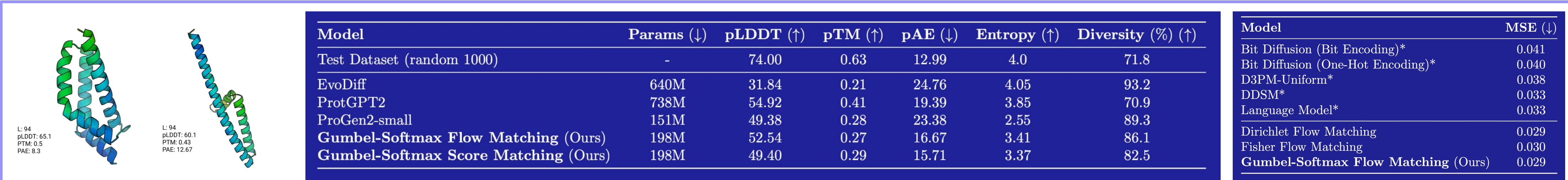†Correspondence to: pranam@seas.upenn.edu

## Motivation

1. **Discretization Errors as a Result of Discrete Iterative Steps.** Discrete diffusion and flow matching models operate in the fully discrete state space, which means that the noisy sequence at each time step is a fully discrete sequence of one-hot vectors sampled from continuous categorical distributions. This can result in discretization errors during sampling when abruptly restricting continuous distributions to a single token.
2. **Deterministic vs Stochastic Flows for *De Novo* Design Tasks.** Many flow matching and optimal transport strategies learn strictly deterministic paths with minimal stochasticity, which is optimal for tasks like matching trajectories, but lacks expressivity and diversity for de novo design tasks like protein or peptide-binder design.
3. **Lack of Training-Free Guidance Methods for Discrete Flow Matching.** Due to the non-differentiability of discrete sequences sampled from relaxed categorical distributions, guidance strategies often involve training classifiers on noisy distributions (classifier-based) or training a separate guided flow model (classifier-free).
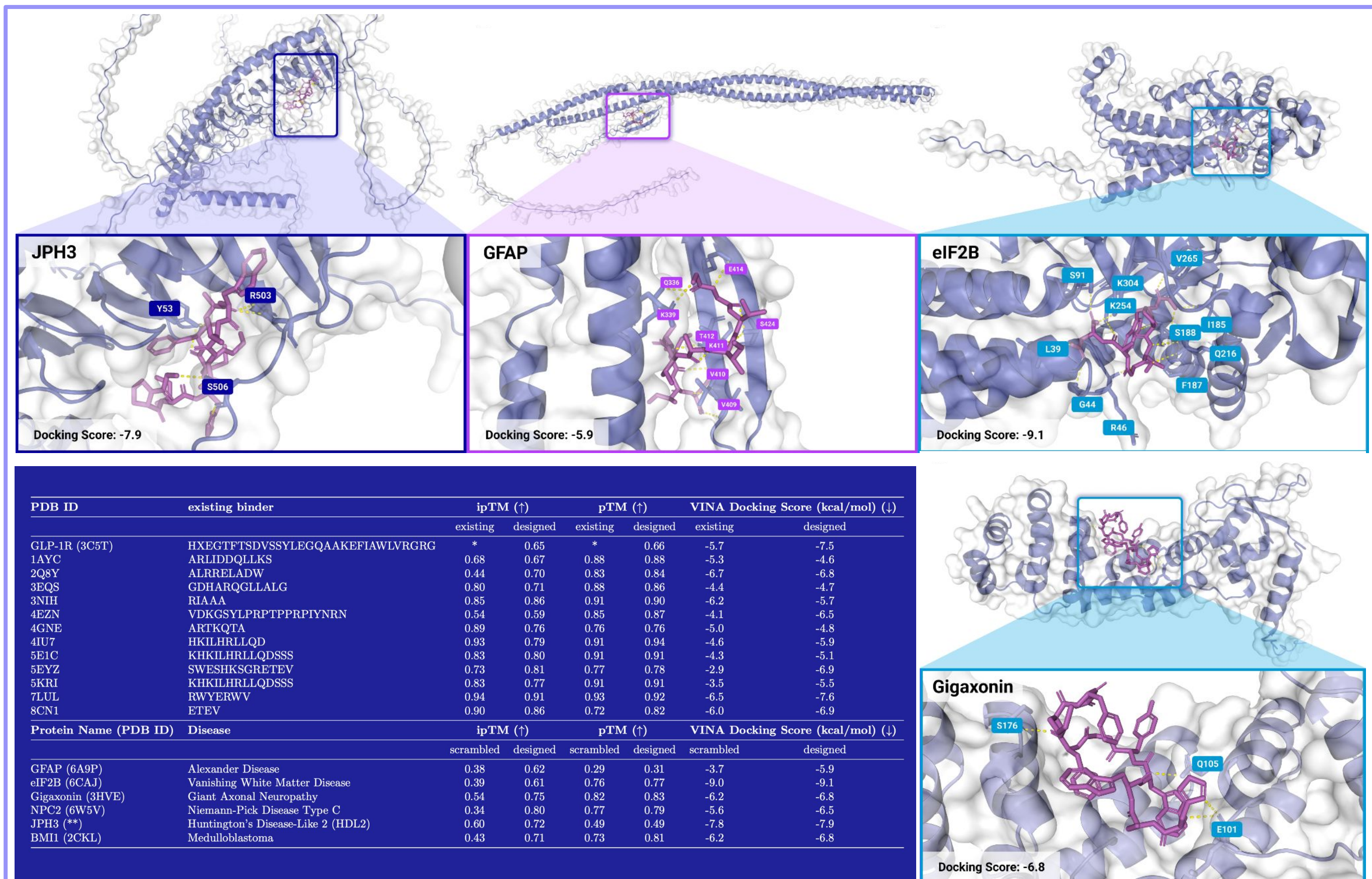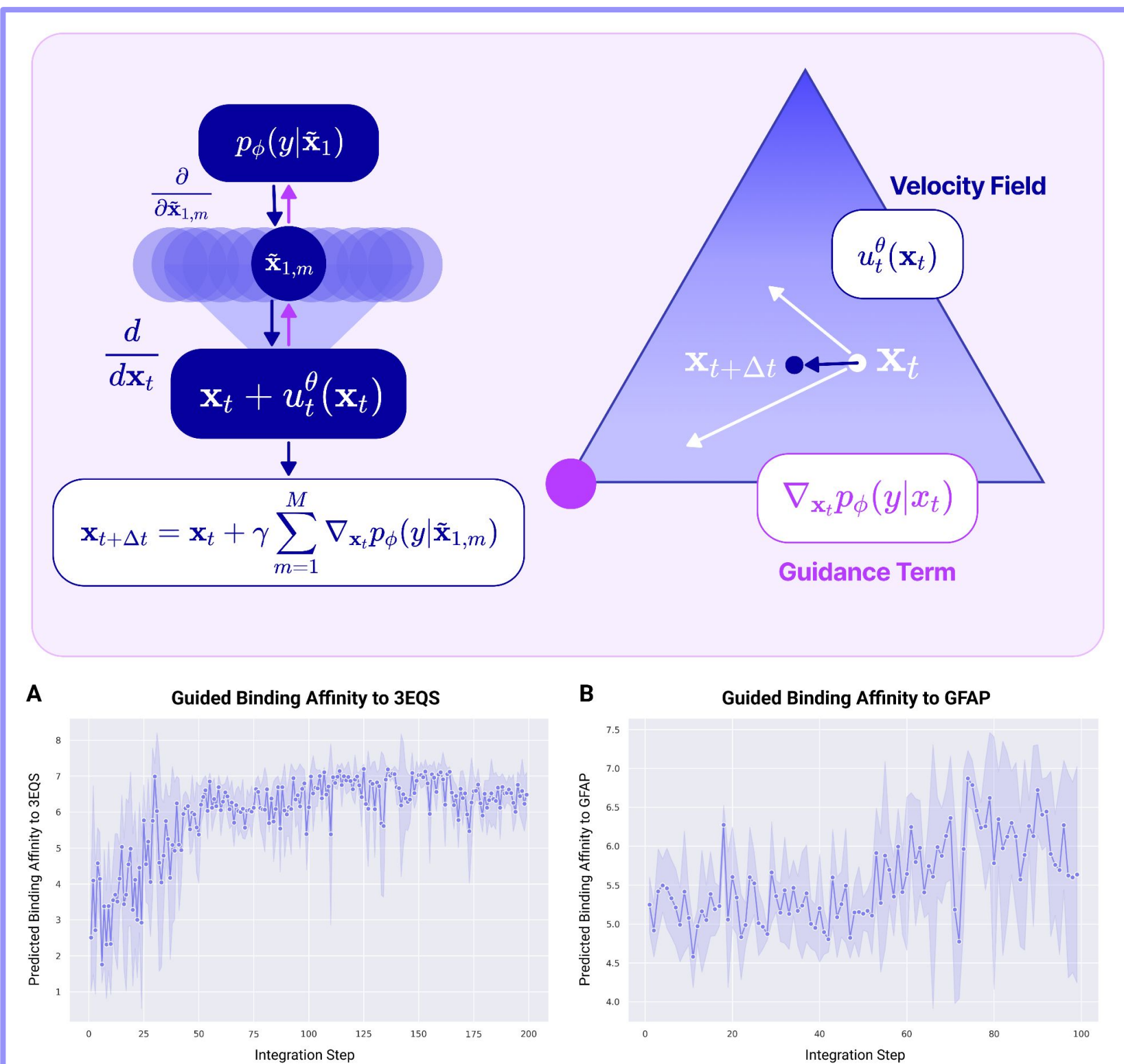
## Gumbel-Softmax Flow and Score Matching for Discrete Generation on the Multi-Dimensional Simplex



**Define a Gumbel-Softmax Noisy Interpolant**
$$\psi_t(\mathbf{x}_1 = \mathbf{e}_k) = \frac{\exp\left(\frac{\delta_{ik}+g_i}{\tau(t)}\right)}{\sum_j \exp\left(\frac{\delta_{jk}+g_j}{\tau(t)}\right)}$$

**Train Denoising Model to Reconstruct Clean Sequence**
$$\mathcal{L}_{\text{gumbel}} = \mathbb{E}_{p_t(\mathbf{x}_t|\mathbf{x}_1=\mathbf{e}_k),p_1(\mathbf{x}_1)}[-\log\langle\mathbf{x}_\theta(\mathbf{x}_t,t),\mathbf{x}_1\rangle]$$

**Compute Conditional Velocity Field During Inference**
$$u_t^\theta(\mathbf{x}_t) = \sum_{k=1}^V u_t(\mathbf{x}|\mathbf{x}_1 = \mathbf{e}_k)\langle\mathbf{x}_\theta(\mathbf{x}_t,t),\mathbf{e}_k\rangle$$

$$\tau(t) = \tau_{\max}\exp(-\lambda t)$$

$$\psi_t(\mathbf{x}_1 = \mathbf{e}_k) = \frac{\exp\left(\frac{\delta_{ik}+g_i}{\tau(t)}\right)}{\sum_j \exp\left(\frac{\delta_{jk}+g_j}{\tau(t)}\right)}$$

**Define a ExpConcrete Noisy Interpolant**
$$\psi_t(\mathbf{x}_1 = \mathbf{e}_k) = \frac{\delta_{ik}+(g_i/\beta)}{\tau(t)} - \log\sum_{j=1}^V \exp\left(\frac{\delta_{jk}+(g_j/\beta)}{\tau(t)}\right)$$

**Train Score Matching Model to Learn Probability Density**
$$\mathcal{L}_{\text{score}} = \mathbb{E}_{p_t(\mathbf{x}_t|\mathbf{x}_1),p_1(\mathbf{x}_1)}\|\text{SM}(\delta_{ik}-\tau(t)x_{t,i}) - \text{SM}(f_\theta(\mathbf{x}_t,t))\|^2$$

**Compute Score During Inference**
$$s_\theta(\mathbf{x}_t,t) = -\tau(t) + \tau(t)V\cdot\text{SM}(f_\theta(\mathbf{x}_t,t)) \approx \nabla_{x_{t,i}}\log p_t(\mathbf{x}_t)$$

**Fig 1. Gumbel-Softmax Flow Matching (Left) and Score Matching (Right).** We define a Gumbel-Softmax interpolant that transforms uniform categorical distributions (high temperature) to distributions concentrated on the vertex of the simplex (low temperature) over time. Then, we train a denoising model to reconstruct the clean distribution from various levels of noise. During inference, we compute the velocity field by taking a weighted sum of the conditional velocity fields scaled by the predicted probabilities.

## Gumbel-Softmax FM and SM for *De Novo* Protein and DNA Promoter Design Tasks



L: 94
pLDDT: 65.1
PTM: 0.5
PAE: 8.3

L: 94
pLDDT: 60.1
PTM: 0.43
PAE: 12.67

| Model | Params (↓) | pLDDT (↑) | pTM (↑) | pAE (↓) | Entropy (↑) | Diversity (%) (↑) |
|---|---|---|---|---|---|---|
| Test Dataset (random 1000) | - | 74.00 | 0.63 | 12.99 | 4.0 | 71.8 |
| EvoDiff | 640M | 31.84 | 0.21 | 24.76 | 4.05 | 93.2 |
| ProtGPT2 | 738M | 54.92 | 0.41 | 19.39 | 3.85 | 70.9 |
| ProGen2-small | 151M | 49.38 | 0.28 | 23.38 | 2.55 | 89.3 |
| **Gumbel-Softmax Flow Matching** (Ours) | 198M | 52.54 | 0.27 | 16.67 | 3.41 | 86.1 |
| **Gumbel-Softmax Score Matching** (Ours) | 198M | 49.40 | 0.29 | 15.71 | 3.37 | 82.5 |

| Model | MSE (↓) |
|---|---|
| Bit Diffusion (Bit Encoding)* | 0.041 |
| Bit Diffusion (One-Hot Encoding)* | 0.040 |
| D3PM-Uniform* | 0.038 |
| DDSM* | 0.033 |
| Language Model* | 0.033 |
| Dirichlet Flow Matching | 0.029 |
| Fisher Flow Matching | 0.030 |
| **Gumbel-Softmax Flow Matching** (Ours) | 0.029 |

**Fig 2. Gumbel-Softmax FM and SM for protein and DNA promoter generation. Left:** Evaluation metrics for the generative quality of protein sequences. Metrics were calculated on 100 unconditionally generated sequences from each model. **Right:** Evaluation of promoter DNA generation conditioned on transcription profile. MSE was evaluated across all validation batches between the predicted signal of a conditionally generated sequence and the true sequence. Regulatory signals were predicted with a pre-trained Sei model.

## Straight-Through Guided Flow Matching (STGFlow) for Target-Binding Peptide Design



$$p_\phi(y|\bar{\mathbf{x}}_1)$$
$$\frac{\partial}{\partial\bar{\mathbf{x}}_{1,m}}$$
$$\bar{\mathbf{x}}_{1,m}$$
$$\frac{d}{d\mathbf{x}_t}$$
$$\mathbf{x}_t + u_t^\theta(\mathbf{x}_t)$$
$$\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + \gamma\sum_{m=1}^M \nabla_{\mathbf{x}_t}p_\phi(y|\bar{\mathbf{x}}_{1,m})$$

**Velocity Field** $u_t^\theta(\mathbf{x}_t)$
$\mathbf{x}_{t+\Delta t}$ — $\mathbf{x}_t$
**Guidance Term** $\nabla_{\mathbf{x}_t}p_\phi(y|x_t)$



JPH3 — Docking Score: -7.9
GFAP — Docking Score: -5.9
eIF2B — Docking Score: -9.1
Gigaxonin — Docking Score: -6.8

**A** Guided Binding Affinity to 3EQS
**B** Guided Binding Affinity to GFAP

| PDB ID | existing binder | ipTM (↑) | | pTM (↑) | | VINA Docking Score (kcal/mol) (↓) | |
|---|---|---|---|---|---|---|---|
| | | existing | designed | existing | designed | existing | designed |
| GLP-1R (3C5T) | HXEGTFTSDVSSYLEGQAAKEFIAWLVRGRG | * | 0.65 | * | 0.66 | -5.7 | -7.5 |
| 1AYC | ARLDDDQLLKS | 0.68 | 0.67 | 0.88 | 0.88 | -5.3 | -4.6 |
| 2Q8Y | ALRRELADW | 0.44 | 0.70 | 0.83 | 0.84 | -6.7 | -6.8 |
| 3EQS | GDHARQGLLALG | 0.80 | 0.71 | 0.88 | 0.86 | -4.4 | -4.7 |
| 3NIH | RIAAA | 0.85 | 0.86 | 0.91 | 0.90 | -6.2 | -5.7 |
| 4EZN | VDKGSYLPRPTPPRPIYNRN | 0.54 | 0.59 | 0.85 | 0.87 | -4.1 | -6.5 |
| 4GNE | ARTKQTA | 0.89 | 0.76 | 0.76 | 0.76 | -5.0 | -4.8 |
| 4IU7 | HKILHRLLQD | 0.93 | 0.79 | 0.91 | 0.94 | -4.6 | -5.9 |
| 5E1C | KHKILHRLLQDSSS | 0.83 | 0.80 | 0.91 | 0.91 | -4.3 | -5.1 |
| 5EYZ | SWEISHKSGRETEV | 0.73 | 0.81 | 0.77 | 0.78 | -2.9 | -6.9 |
| 5KRI | KHKILHRLLQDSSS | 0.83 | 0.77 | 0.91 | 0.91 | -3.5 | -5.5 |
| 7LUL | RWYERWV | 0.83 | 0.91 | 0.93 | 0.92 | -6.5 | -7.7 |
| 8CN1 | ETEV | 0.90 | 0.86 | 0.72 | 0.82 | -6.0 | -6.9 |

| Protein Name (PDB ID) | Disease | ipTM (↑) | | pTM (↑) | | VINA Docking Score (kcal/mol) (↓) | |
|---|---|---|---|---|---|---|---|
| | | scrambled | designed | scrambled | designed | scrambled | designed |
| GFAP (6A9P) | Alzheimer Disease | 0.38 | 0.62 | 0.29 | 0.31 | -3.7 | -5.9 |
| eIF2B (6CAJ) | Vanishing White Matter Disease | 0.39 | 0.61 | 0.76 | 0.77 | -9.0 | -9.1 |
| Gigaxonin (3HVE) | Giant Axonal Neuropathy | 0.54 | 0.75 | 0.82 | 0.83 | -6.2 | -6.8 |
| NPC2 (6W5V) | Niemann-Pick Disease Type C | 0.34 | 0.80 | 0.77 | 0.79 | -5.6 | -6.5 |
| JPH3 (**) | Huntington's Disease-Like 2 (HDL2) | 0.60 | 0.72 | 0.49 | 0.49 | -7.8 | -7.9 |
| BMI1 (2CKL) | Medulloblastoma | 0.43 | 0.71 | 0.73 | 0.81 | -6.2 | -6.8 |

**Fig 3. Straight-Through Guided Flow Matching (STGFlow).** To fill the gap in inference-time guidance algorithms for discrete flow matching, we introduce STGFlow, a classifier-based guidance algorithm that leverages straight-through gradient estimators to guide the flow trajectory towards high-scoring sequences sampled from the Gumbel-Softmax distribution.

**Fig 4. Comparison of ipTM and VINA docking scores for peptide binders generated from Gumbel-Softmax FM with STGFlow to experimentally-validated binders and scrambled controls** (for targets without known binders). Generated binders show consistently stronger binding affinity across both tasks.

## Conclusions

1. We define a temperature-controlled Gumbel-Softmax interpolation and derive a velocity field that enables smooth transport from noisy to clean distributions on the interior of the simplex.
2. By applying Gumbel noise during training, Gumbel-Softmax FM avoids overfitting the training data, increasing the exploration of diverse flow trajectories.
3. To address the lack of training-free guidance methods for discrete flow matching, we propose STGFlow, a classifier-based guidance method that leverages straight-through gradient estimators to steer the velocity field toward optimal sequences on the data manifold.

## Paper