

DiffusionSELEX

1 Introduction

Systematic Evolution of Ligands by Exponential Enrichment (SELEX) is a powerful method for discovering novel ligands with high affinity and specificity for target molecules. However, the experimental process is time-consuming, resource-intensive, and involves numerous parameters that can significantly impact the outcome. This thesis aims to develop an in-silico SELEX simulator using diffusion models, which are a class of generative models that learn to denoise data by iteratively refining a signal through a series of noise-removal steps. By leveraging the power of diffusion models, we can streamline ligand discovery and optimize experimental conditions.

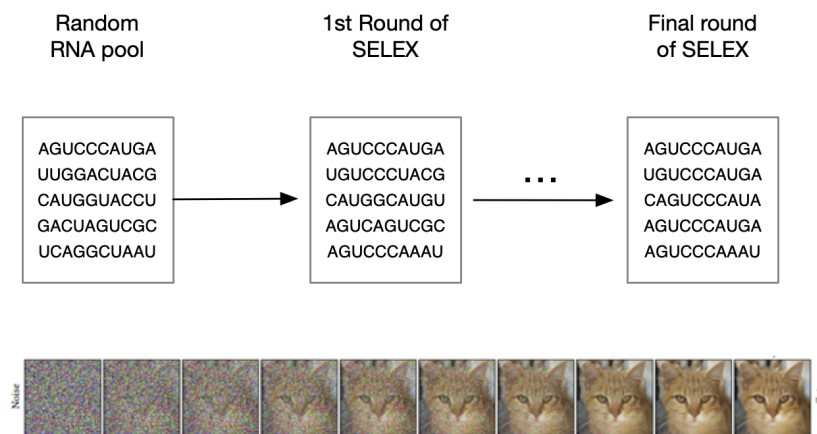


Figure 1: Comparison of a SELEX diffusion approach (top) and a standard image diffusion process (bottom). RNA pools evolve from random (left) to specific (right), with enrichment of specific target-affine sequences. Image diffusion starts with Gaussian noise (left) and progressively denoises to generate a realistic image.

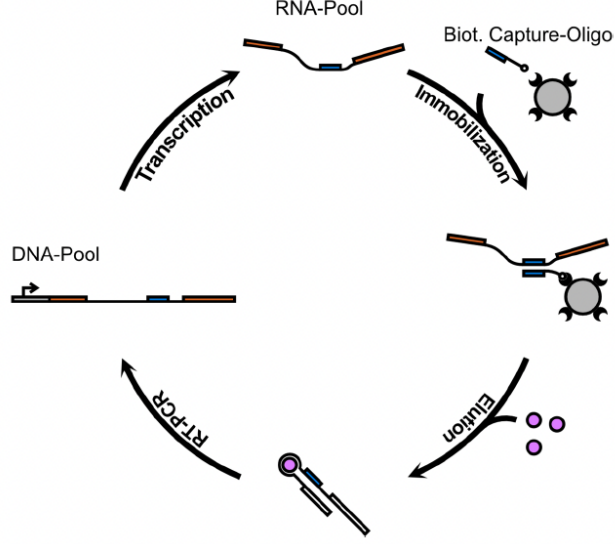


Figure 2: Schematic representation of the SELEX process for aptamer selection. Starting from a random RNA library, iterative rounds of selection, binding, and amplification are performed to enrich for high-affinity aptamers against a specific target molecule. Image from [1].

2 Approach

When transferring the process of the SELEX to a diffusion process, the idea is to consider each sequence in the pool of the final round as a sample, analogously to the image of the cat in fig. 1. This sequence should be diffused, i.e. its entries should randomly change from one round to the next. However, one sequence in one round of the SELEX will also appear unchanged in the previous round. In fact, the additional sequences that exist in previous rounds compared to later rounds, increase the randomness, while the entries of the sequences do not change. This is a problem and we will address this in the following.

To consider the SELEX as a diffusion process, we need to connect a chosen sample from the last round to one sequence of the previous rounds, respectively. The most straightforward way would be to choose the sequence of the previous round that is closest to the current sample. However, this would usually result in identical sequences from the final round to the random pool. In order to increase the randomness, the idea for the construction of the forward pass is the following: Choose a sample s_0^i from the final round, compare it with sequence s_1^j , $j \in \{1, \dots, N\}$ from the previous pool and make a similarity score for each pair. This similarity score should be transferred to probabilities for the transition of the current sample to the corresponding sequence in the previous pool $p(x_1 = s_1^j | x_0 = s_0^i)$. Then sample x_1 from this distribution and do the same for the

next round until the SELEX has reached its first round. With this approach, we have a well-defined forward process for the diffusion model which can be inserted into the diffusion model framework.

3 Use sampling without replacement to match the pools

The following procedure is done by image diffusion:

Algorithm 1 High-level Diffusion Model Training

Require: \mathcal{D} : dataset of clean images, T : number of diffusion steps, f_θ : neural network with parameters θ

Ensure: f_θ : neural network

```

1: for  $i = 1, \dots, \text{num\_epochs}$  do
2:   for  $x_0 \in \mathcal{D}$  do ▷ Iterate over clean images
3:      $t \sim \text{Uniform}(\{1, \dots, T\})$  ▷ Sample a random time step
4:      $x_t \leftarrow \text{FORWARDDIFFUSION}(x_0, t)$  ▷ Apply forward diffusion to
       obtain a noisy image
5:      $\epsilon_\theta \leftarrow f_\theta(x_t, t)$  ▷ Predict noise using the neural network
6:      $\mathcal{L} \leftarrow \text{COMPUTELOSS}(\epsilon_\theta, \epsilon)$  ▷ Compute the loss between predicted
       and actual noise
7:      $\theta \leftarrow \theta - \eta \cdot \nabla_\theta \mathcal{L}$  ▷ Update neural network parameters using gradient
       descent
8:   end for
9: end for
10: return  $f_\theta$ 

```

Noise needs to be added in several steps to each pixel to obtain noised images which can be used for training to reconstruct original image estimation from noised images. After training, one can sample pure noise and apply the neural network to create a random image. Experiments show that this is very effective to create new realistic images.

For images, there is a 1:1 relationship from data to noise. I.e. we exactly know to which original image pixel a noise-added pixel corresponds. In our case, having different pools this is not so easy. But an one-to-one relationship can be enforced by removing matched sequences from the pool. This is similar to flow-matching.

4 Simple Markov Model Estimation

Algorithm 2 Forward Process

Require: $\mathcal{X} = \{\mathbb{X}_0, \dots, \mathbb{X}_R\}$: dataset of pools of R rounds, round 0 is random pool, round R is specific pool. Pools containing sequences \mathbf{s} together with frequencies $n_{\mathbf{s}}$ for each sequence. We denote $n_{\mathbf{s}} = n_{\mathbf{s}, \mathbb{X}_i} \in \mathbb{R}$ the frequency of sequence \mathbf{s} (in pool \mathbb{X}_i), and $\mathbf{n}_{\mathbb{X}_i} \in \mathbb{R}^{N_{\text{unique}}}$ the vector containing all frequencies for all unique sequences in pool i . T : number of diffusion steps, f_{θ} : neural network with parameters θ .

Require: All pools have the same number of sequences, i.e. $\mathbf{1}^T \mathbf{n}_{\mathbb{X}_i} =: N$ is the same for all i . Since this is in general not the case, sample N sequences of every pool, weighting each probability of a sequence by its relative frequency in the pool.

Ensure: f_{θ} : trained neural network

```
1: for  $t = 1, \dots, T$  do
2:    $i \leftarrow \text{UNIFORM}(\{1, \dots, R\})$  Pick random round on which to add noise.
3:   for sequence  $\mathbf{s}$ , randomly selected, in pool  $\mathbb{X}_i$  do  $\triangleright$  While  $\mathbf{1}^T \mathbf{n}_{i, \mathbb{X}_i} > 0$ 
4:      $\mathbf{d}_{\mathbf{s}, \mathbb{X}_{i-1}} = (d(\mathbf{s}, \tilde{\mathbf{s}}))_{\tilde{\mathbf{s}} \in \mathbb{X}_{i-1}}$   $\triangleright$  calculate (Levenstein) distances from  $\mathbf{s}$ 
      to all sequences of pool  $i-1$ 
5:      $\tilde{\mathbf{p}}_{\mathbf{s}, \mathbb{X}_{i-1}} = \text{COMPUTESCORE}(\mathbf{d}_{\mathbf{s}, \mathbb{X}_{i-1}}) \odot \mathbf{n}_{\mathbb{X}_{i-1}}$   $\triangleright$  Calculate
      probabilities proportional to the frequencies of Sequences in pool  $i-1$ 
6:      $\mathbf{p}_{\mathbf{s}, \mathbb{X}_{i-1}} = \frac{\tilde{\mathbf{p}}_{\mathbf{s}, \mathbb{X}_{i-1}}}{\mathbf{1}^T \tilde{\mathbf{p}}_{\mathbf{s}, \mathbb{X}_{i-1}}}$   $\triangleright$  Normalize to obtain probabilities
7:      $\mathbf{s}_{i-1} \leftarrow \text{SAMPLE}(\mathbb{X}_{i-1}, \mathbf{p}_{\mathbf{s}, \mathbb{X}_{i-1}})$   $\triangleright$  Apply forward diffusion by
      discrete sampling to obtain a noisy sequence from the previous pool
8:      $n_{\mathbf{s}_i, \mathbb{X}_i} \leftarrow n_{\mathbf{s}_{i-1}, \mathbb{X}_{i-1}} - 1$ 
9:      $n_{\mathbf{s}_{i-1}, \mathbb{X}_{i-1}} \leftarrow n_{\mathbf{s}_{i-1}, \mathbb{X}_{i-1}} - 1$   $\triangleright$  Remove matched sequences from the
      pools to match all sequences from pool  $i$  with a sequence from pool  $i-1$ 
10:     $\hat{\mathbf{s}}_i \leftarrow f_{\theta}(\mathbf{s}_{i-1}, i)$   $\triangleright$  Predict original sequence
11:     $\mathcal{L} \leftarrow \text{COMPUTELOSS}(\hat{\mathbf{s}}_i, \mathbf{s}_i)$   $\triangleright$  Compute the loss between predicted
      and actual noise
12:     $\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} \mathcal{L}$   $\triangleright$  Update neural network parameters using gradient
      descent
13:   end for
14: end for
15: return  $f_{\theta}$ 
```

Algorithm 3 Pseudocode for SELEX diffusion sampling

Ensure: f_{θ} : trained neural network

```
1:  $s_T \sim \text{UNIFORM}(\{A, C, T, G\})$  Sample random sequence from uniform distri-
   bution
2: for  $t = T-1, \dots, 1$  do
3:    $p_t \leftarrow f_{\theta}(s_{t+1})$  Compute probabilities for each time step given the state
     at time  $t+1$ 
4:    $s_t \sim p_t$  Sample the next state for time  $t$ 
5: end for return  $s_0$ 
```
