

Simulate co-transcriptional folding kinetics by genetic algorithm

Zhuoran Qiao

August 8, 2018

We developed a genetic algorithm based approach to simulate kinetics of co-transcriptional folding. Our method is built on two following assumptions:

- 1 Most of RNA secondary structures (SS) are linkage of locally optimal or sub-optimal structures (foldons) at various folding sites;
- 2 Global structural rearrangement of a partial RNA segment is permitted only if it's folding to the optimal SS on that segment.

Formally, we denote direct product... irreducible foldon representation... green's function...

1 Algorithm procedure

During every elongation step, an active species pool of strands with unique SS and different population is updated. New candidate strands with length of $L + \Delta L$ are generated by a recombination process: for every old strand, all indices in its IFR is identified as possible rearrangement site, then its child strands is generated by linking partial segments with a foldon that terminated at $L + \Delta L$. We assume that elongation will not change the initial population distribution of secondary structures: child strands with the exact parental SS on $[0, L]$ will also inherit the population of their parents.

After structural generation the rate matrix among all candidate strands within the new active species pool is calculated (see part 3). Then the population distribution of strands after elongation is computed by chemical master equation, and one iterative elongation step is finished.

Pseudocodes of the whole procedure are as follows:

Algorithm 1 Co-transcriptional folding elongation procedure

```
1: Initialize active pool
2: while sequence length > current length do
3:   old pool  $\leftarrow$  active pool
4:   renew active pool
5:   current length  $+= dL$ 
6:    $dt \leftarrow dL / \text{transcription rate}$ 
7:   for left boundary  $\in \{0, dL, 2dL, \dots, \text{current length} - dL\}$  do       $\triangleright$  Get all new foldons
8:      $D^{foldon}(\text{left boundary}, \text{current length}) \leftarrow \text{numpy.mfe}(\text{sequence}[\text{left boundary}, \text{current length}])$ 
9:   end for
10:  for strand  $\in$  old pool do       $\triangleright$  Recombination
11:    for rearrangement site  $\in$  strand.IFR do
12:      Candidate  $\leftarrow$ 
13:       $D^{strand}(0, \text{rearrangement site}) \oplus D^{foldon}(\text{rearrangement site}, \text{current length})$ 
14:    end for
15:  end for
16: end while
```

Algorithm 2 Euclid's algorithm

```
1: procedure EUCLID( $a, b$ )       $\triangleright$  The g.c.d. of  $a$  and  $b$ 
2:    $r \leftarrow a \bmod b$ 
3:   while  $r \neq 0$  do       $\triangleright$  We have the answer if  $r$  is 0
4:      $a \leftarrow b$ 
5:      $b \leftarrow r$ 
6:      $r \leftarrow a \bmod b$ 
7:   end while
8:   for <some condition> do
9:     <do stuff>
10:  end for
11:  return  $b$        $\triangleright$  The gcd is  $b$ 
12: end procedure
```

2 Secondary structure generation

Our goal is to figure out a general approach to calculate transition path statistics (average transit time and transit time distribution) for free energy landscape with discrete states. The Markov rate matrix \mathbf{W} element between state σ_i and state σ_j is given by

$$\mathbf{W}_{ij} = \langle \sigma_i | \mathbf{W} | \sigma_j \rangle = k_0 \exp\left(-\frac{F_j - F_i}{2}\right) \quad (1)$$

$$\begin{aligned} P(t) &= \sum_{\varphi} P(t|\varphi)P(\varphi) \\ &= \sum_{\varphi} P(t|\vec{H}[\varphi])P(\varphi) \\ &= \sum_{\vec{H}} P(t|\vec{H})P(\vec{H}) \\ &= \sum_{\vec{H}, l} P(t|\vec{H})P(\vec{H}|l)P(l) \end{aligned} \quad (2)$$

3 Folding pathway identification & Rate calculation

- a. Using a genetic-algorithm based method and NUPACK to predict populated RNA configurations during cotranscriptional folding (in progress);
- b. Using master equation method to simulate evolution of folding configurations and SD sequence accessibility.