
Algorithm 1 SRI Generator

Given: We define $\mathbb{B} = (A, T, C, G) = (\mathbb{B}_i)_{i=1}^4$ to be the sequence of canonical bases. Furthermore, let $\mathbb{O}^K = (\mathbb{O}_i^K)_{i=1}^{4^K}$ be the sequence of *oligomers* of length K , e.g., $\mathbb{O}^3 = (AAA, AAT, \dots, CGG, GGG)$. (Notice that \mathbb{O}^3 has $4^3 = 64$ elements.) Similarly, let $\mathbb{F}^K = (\mathbb{F}_i^K)_{i=1}^{4^K}$ be the sequence of oligomer *frequencies* of length K , such that $\sum_{i=1}^{4^K} \mathbb{F}_i^K = 1$ and $\mathbb{F}_i^K \in [0, 1] \forall i$. For each fixed i , the oligomer \mathbb{O}_i^K has one corresponding frequency value \mathbb{F}_i^K . These frequencies are computed from a sample input sequence using SRI Analyzer. Let $\mathcal{N} \geq 1$ be the user-chosen oligomer length for which the frequency composition is being approximated. Finally, we assume the input sequence is composed purely of A, T, C, and G. For a modified version of this algorithm which handles impure samples, please see our source code.

Ensure: The output sequence(s) *approximates* the short-range inhomogeneity of the input sequence. In other words, the input and output sequence(s) will share a similar \mathcal{N} -mer frequency composition. The output sequence(s) will be *randomly* constructed to satisfy this constraint.

```
1:  $\langle \mathbb{O}^{\mathcal{N}}, \mathbb{F}^{\mathcal{N}} \rangle = \text{ReadCompositionTable}()$  #Example for  $\mathcal{N} = 1$ :  $\langle (A, T, C, G), (0.25, 0.40, 0.20, 0.15) \rangle$ 
   #Calculate the demarcation table for the random selection of the first oligomer.
2:  $sum = 0$ 
3: for  $i = 1$  to  $4^{\mathcal{N}}$  do
4:    $sum \leftarrow sum + \mathbb{F}_i^{\mathcal{N}}$ 
   #Let partialSum be an array such that  $partialSum^{(i)} \in [0, 1]$ , for  $1 \leq i \leq 4^{\mathcal{N}}$ .
5:    $partialSum^{(i)} \leftarrow sum$ 
6: end for
7: for each FASTA sequence, “seq” do
8:    $length = \text{GetSequenceLength}(seq)$ 
9:   if  $length \geq \mathcal{N}$  then
10:     $R_1 \leftarrow \text{GenerateRandomNumber}(0, 1)$  #Generate some random  $R_1 \in [0, 1]$ .
11:    for  $i = 1$  to  $4^{\mathcal{N}}$  do
12:      if  $R_1 < partialSum^{(i)}$  then
13:         $randomSeq \leftarrow \mathbb{O}_i^{\mathcal{N}}$ 
14:        Exit Loop.
15:      end if
16:    end for
17:    for  $i = 1$  to  $length - \mathcal{N}$  do
18:       $sum \leftarrow 0$ 
      #The next base is chosen randomly using the frequencies of the 4 possible overlapping  $\mathcal{N}$ -mer sequence tails.
19:       $R_2 \leftarrow \text{GenerateRandomNumber}(0, 1)$ 
20:       $tail \leftarrow \text{Suffix}(seq, \mathcal{N} - 1)$  #Example:  $\text{Suffix}(\text{“hotdog”}, 3) = \text{“dog”}$ .
21:      for  $i = 1$  to 4 do
22:         $oligo \leftarrow \text{Concatenate}(tail, \mathbb{B}_i)$  #Example:  $\text{Concatenate}(\text{“hot”}, \text{“dog”}) = \text{“hotdog”}$ .
23:         $f \leftarrow \text{GetOligoFrequency}(oligo)$  #GetOligoFrequency(oligo) =  $\mathbb{F}_j^{\mathcal{N}}$  for one  $j$  such that  $\mathbb{O}_j^{\mathcal{N}} = oligo$ .
24:         $sum \leftarrow sum + f$ 
        #Let demarcation be an array such that  $demarcation^{(i)} \in [0, 1]$ , for  $1 \leq i \leq 4$ .
25:         $demarcation^{(i)} \leftarrow sum$ 
26:        if  $R_2 < demarcation^{(i)}$  then
27:           $randSeq \leftarrow \text{Concatenate}(randSeq, \mathbb{B}_i)$ 
28:          Exit Loop.
29:        end if
30:      end for
31:      if  $R_2 \geq sum$  then
32:         $randomBase \leftarrow \text{PickRandomBase}()$  #Randomly choose a base from {A,T,C,G}.
33:         $randSeq \leftarrow \text{Concatenate}(randSeq, randomBase)$ 
34:      end if
35:    end for
36:     $\text{WriteOutputFile}(randSeq)$ 
37:  end if
38: end for
```
