

1 Basic definitions

My first attempt at providing a coherent mathematical definition of a repeat element family in the context of a spaced seed.

Definition 1 *An L -mer is a string of length l .*

Definition 2 *Given a spaced seed s of length l and an l -mer x , the seeded l -mer s' is created by removing all characters of x corresponding to 0's in s . Let $\sigma_s(x)$ denote the seeded l -mer for s .*

Example: for $s = 11011$ and $x = AACGG$, $\sigma_s(x) = AACC$.

Definition 3 *A repeat element family descriptor (refd) is a string over the alphabet $\{A, C, G, T, *\}$, describing the contents of any string in the family.*

Definition 4 *Given refds r and r' , we will say r' is a slack substring of r (denoted $r' \prec r$) if there is a substring r'' of r such that (1) $|r'| = |r''|$, and (2) for all $0 \leq i < |r'|$, either $r'_i = r''_i$, $r'_i = *$, or $r''_i = *$.*

(In otherwords, its a substring, with a potentially different $*$ pattern.) **Note:** I'm not sure if we should be allowing $r'_i = *$.

Definition 5 *We say an refd r matches a genome G at position i if, for all j such that $r_j \neq *$, $r_j = G_{i+j}$.*

Example: If $r = AA*TT$, and $G = AACTTGGAAGTT$, then r matches G at positions $i = 0$ and $i = 7$. If $G = AAATTT$, then r matches G at positions $i = 0$ and $i = 1$.

Definition 6 *A spaced seed s matches an refd r at position i if, for every $0 \leq j < |s|$, $s_j = 0$ whenever $r_{i+j} = *$.*

Example: If $s = 11011$ and $r = AAAAA*GGG$, then s matches r at positions $i = 0$ and $i = 3$, but not at any other i .

Observation 1 *s matches r at i if $\sigma_s(r[i : i + |s|])$ does not contain any $*$ symbols.*

$\sigma(r[3 : 8]) = AAGG$ (hence a match at $i = 3$), but $\sigma(r[2 : 7]) = AA*G$ (hence no match at $i = 2$).

Definition 7 *Given a seed s and refd r , let $M_s(r)$ be the set of values i such that s matches r at i .*

Let $s = 11011$ and $r = AA * CC * GGGGGG$. Then $M_s(r) = [0, 3, 6, 7]$.

Definition 8 A spaced seed s is consistent with an refd r if for every i , $0 \leq i < |r|$, there is some $i - |s| \leq j \leq i$ such that s matches r at position j .

In other words: for any position i of the refd, we must be able to match the seed to a position of r such that it then covers position i .

Example: The seed 11011 is compatible with $r = AA * AA * AA$. (The seed matches at positions 0 and 3, and all positions are covered by these two.) But it is not consistent with $AAAAA * AAAAA$, and there is no seed that can match this string at any position that can cover $i = 5$ or $i = 6$.

Observation 2 Let L be the sortest sequence of the values in $M_s(r)$. Then s is consistent with r if and only if $\max_{0 \leq j < |s|-1} L[j+1] - L[j] \leq |s|$.

That is, in the sorted list, every pair of adjacent elements must be within $|s|$ of each other.

Definition 9 Given a fixed genome G , and fixed spaced seed s , and a fixed value f , we define a elementary repeat family as a set S of genome coordinates, $|S| \geq f$, such that there exists an refd r where:

- r matches the sequence of length $|r|$ starting at each element of S . (S is the set of all instances.)
- s is consistent with r . (r corresponds to the seed.)
- There does not exist an refd r' , $r' \prec r$, such that r' is consistent with s and $M_s(r') - M_s(r) \neq \emptyset$. (You cannot have a proper substring of r that describes sequences outside of the instances described by r – minimality.)
- There does not exist an refd r' , $r \prec r'$, such that s is consistent with r' , such that the instances defined by $M_s(r')$ contain all the instances of $M_s(r)$. (Maximality.)

Comment: I'm not sure if the $r \prec r'$ is the right relationship. Maybe just straight substring? Or perhaps the definition of \prec isn't quite right?

Definition 10 Given the refd r of an elementary repeat family with set S , we say that r is tight if, for each i such that $r_i = *$, there exists two sequences defined by S that have different bases in position i .

In other words: r is tight if it only uses $*$ symbols where it must to match everything sequence defined by S .

Comments: I had it in mind that the algorithm would always return a tight refd. But in Carly's defense she gave the example: $s=11011$, $AAAAATCCCC$, $AAGAATCCGCC$, where ends up with $AA*AA*AA*AA$, which is not tight. Interestingly, if we have $AAAAATTCCCC$ and $AAGAATCCGCC$, then we get $AA*AATTCC*CC$ – to that extra T makes it tight. Not sure if this is significant.

2 From Thesis

In the following I'm going through definitions / parallels from Nate's thesis and seeing if I can create an analogy for spaced seeds.

The following are things I wanted to try to prove that may or may not be useful. In all of these I'm assuming a fixed genome G , a fixed frequency requirement f , and a fixed seed s with length l and weight w .

Lemma 1 *For any l -mer x , the seeded l -mer $\sigma_s(x)$ can be a member of at most l different families.*

My gut is that this is true, but not necessarily tight. Maybe its w , or $l - w$?

Definition 11 *Let x and y be two strings such that $x = a\dot{b}$, $y = b\dot{c}$, and $|b| = i$. Then $x \circ_i y = a\dot{b}\dot{c}$.*

$AAACC \circ_2 CCGGG = AAACCGGG$. $AAACC \circ_3 CCGGG$ is undefined.

This is a modification of Nate's $x \circ y$ operator. May need to be adapted for seeds.

3 Random Lemmas