

1

# Supplementary to Red: an intelligent, rapid, accurate tool for detecting repeats de-novo on the genomic scale

2 Hani Z Girgis\*<sup>1</sup>

<sup>1</sup> Computational Biology Branch, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, 8600 Rockville Pike, Bethesda, MD 20894

<sup>2</sup> Tandy School of Computer Science, University of Tulsa, 800 South Tucker Drive, Tulsa, OK 74104

Email: Hani Z Girgis\* - girgishz@mail.nih.gov;

\*Corresponding author

## 3 **Supplementary Methods**

4 In this section, I discuss the following: (i) how the labeling module delineates candidate repetitive regions and  
5 potential non-repetitive regions, (ii) the run-time analysis of Red, (iii) the default values of Red's parameters,  
6 (iv) Red's parameters on the unassembled *Drosophila melanogaster* genome, (v) how the related tools were  
7 executed, and (vi) the source of the data used for evaluating the performance of Red and the related tools.

### 8 **Delineation of candidates and potential non-repetitive regions**

9 Once the local maxima have been found by the labeling module, the boundaries of the candidate regions  
10 are determined. Locating candidate repetitive regions involves the interleaving non-repetitive regions. A  
11 region is considered non-repetitive (repetitive) if the percentage of low ( $\leq$  a user-specified threshold) scores  
12 in this region is greater (less) than the expected percentage assuming a uniform background distribution.  
13 The expected percentage is the percentage of the scores that are  $\leq$  the score of the threshold in the genome.  
14 For example, suppose that 55% of the scores of a genome are  $\leq 2$ . If these low scores were distributed  
15 uniformly, then one would expect the percentage of low scores in a random segment of this genome to be  
16 about 55%. However, due to the fact that non-repetitive regions include a high percentage of low scores,  
17 the observed percentage should be greater than 55%. In the current implementation of Red, if the expected  
18 percentage is low, it is increased to 52.5%.

19 Properties of repetitive regions are utilized for reducing the number of false maxima. The original scores

20 in the small region flanking a local maximum are examined to ensure that this maximum is in a repetitive  
21 region; the percentage of the low scores in the targeted region must be lower than the expected percentage.  
22 The targeted region has the same size as the Gaussian mask.

23 The presence of local maxima and of high scores is characteristic of repeats, whereas non-repetitive  
24 regions consist mainly of low scores. Two definitions are required. A separator and a core illustrate how the  
25 boundaries of repetitive regions are determined. A separator is defined as a non-repetitive region located  
26 between two consecutive local maxima; consequently, a separator does not include any local maximum.  
27 A core is defined as a repetitive region including at least one local maximum, and it is bounded by two  
28 separators. To begin the delineation of repetitive regions, the separators are identified and, in turn, the  
29 cores. Then, the boundaries of each core are adjusted.

30 The boundaries of a core are expanded or eroded in two stages. The first stage is a step-by-step expansion.  
31 During this stage, the small region adjacent to the start or the end of the core is added to the core if this  
32 region is a repetitive region. The size of this region, i.e. the step, is half the width of the Gaussian mask.  
33 This stage is repeated until a non-repetitive region is encountered. In the second stage, the start and the end  
34 of the core are further expanded or eroded nucleotide by nucleotide. The one-by-one expansion is executed  
35 if the score adjacent to the core is greater than that of the threshold. In this case, the boundaries of the core  
36 are adjusted to include this score. Expanding the core is repeated until a score that is less than or equal  
37 to the score of the threshold is encountered. The one-by-one erosion is executed if the original score at the  
38 start or the end of the core is less than or equal to that of the threshold. Then the boundaries of the core  
39 are adjusted to exclude this score. Eroding the core is repeated until a score that is greater than that of the  
40 threshold is encountered. During the second stage, the start or the end of the core is either expanded or  
41 eroded; however, it is not subject to the two operations combined.

## 42 **Run time analysis**

43 In order to determine the total run time of Red, I calculate the run time depending on the implementation,  
44 i.e. the code, of each of the four modules. To start, consider a genome of length  $n$ , a background Markov  
45 chain of the  $o^{th}$  order, a mask of size  $m$ , and an HMM with  $s$  states. The following is an analysis of the time  
46 required by each module.

- 47 • **The scoring module:** This module takes  $n$  to count the k-mers,  $o \times n$  to train the background Markov  
48 chain, and a constant time  $c$  to adjust the counts of k-mers. Thus, the total time taken by the scoring  
49 module is  $(o + 1) \times n + c$ .

50 • **The labeling module:** It takes  $n$  to score a sequence,  $m \times n$  to smooth the scores,  $n$  to calculate the  
51 first derivative,  $n$  to calculate the second derivative,  $n$  to find the local maxima, at most  $n$  to find the  
52 separators, and at most  $n$  to expand the candidate regions. Therefore, the total time required by the  
53 labeling module is  $(6 + m) \times n$ .

54 • **The training module:** The module takes  $n$  to score a sequence,  $n$  to calculate the states using a  
55 logarithmic function of the scores, at most  $n$  to calculate the prior probabilities, and  $n$  to calculate the  
56 transition probabilities. The total time required by the training module is  $4 \times n$ .

57 • **The scanning module:** It takes  $n$  to score a sequence,  $n$  to calculate the states, and  $s \times n$  to calculate  
58 the optimal series of states according to Viterbi's algorithm. The total time for the scanning module is  
59  $(2 + s) \times n$ .

60 In sum, the total run time of Red is  $(13 + o + s + m) \times n + c$ . Because  $c$ ,  $o$ ,  $s$ , and  $m$  are constants, the  
61 upper bound of the run time is  $O(n)$ , i.e. linear with respect to the length of the genome.

## 62 **System defaults**

63 The system allows the user to specify the value of five parameters. However, an average user is likely to use  
64 Red with the default settings. Therefore, the default values of the parameters are important. Red has the  
65 following parameters.

### 66 *Word length*

67 The default value is equal to the floor of the logarithmic (base 4) value of the effective length of the genome.  
68 The effective length of a genome is the number of all nucleotides except the ones labeled as "N." The  
69 minimum default value is 12, and the maximum default value is 15. The user can adjust the length of the  
70 word using the "-len" parameter.

### 71 *Order of the background Markov chain*

72 The default value equals  $\lfloor \text{word length} \div 2 \rfloor - 1$ . On one hand, if the order is too small, the background Markov  
73 chain may not be realistic. On the other hand, if the order is too large, the background Markov chain may  
74 memorize the genome. Therefore, a Markov chain that has an order of roughly half of the word length  
75 provides practical estimation of the expected count of a word in a genome with a similar composition of  
76 short words. The user can adjust the order using the "-ord" parameter.

77 ***Half width of the Gaussian mask***

78 The value of this parameter controls the width of the mask. The size of the mask is twice the value of this  
79 parameter. The standard deviation of the Gaussian distribution equals the value of this parameter divided  
80 by 3.5 because almost all samples fall within 3.5 standard deviations on each side of the mean. There are  
81 two default values: 20 and 40. If the G-C content is greater than 67% or less than 33%, the default value is  
82 40. Otherwise, the default value is 20. The user can adjust the half width using the “-gau” parameter.

83 ***Score of the threshold of low scores in non-repetitive regions***

84 The default score of the threshold of low scores is 2. A score of 2 indicates that the observed count of a k-mer  
85 in the genome is greater than the expected count by 2. Non-repetitive regions including duplicated segments  
86 consist mainly of less abundant k-mers. Therefore, a threshold of 2 suits the definition of non-repetitive  
87 regions. The user can adjust the score of the threshold using the “-thr” parameter.

88 ***Minimum observed copy number of a word***

89 For a word having observed copy number less than or equal to the value of this parameter, the adjusted  
90 count is zero. The default value of this parameter is 2 to avoid duplicated segments. This parameter is very  
91 important for processing unassembled genomes (see the next section). The user can adjust the value of this  
92 parameter using the “-min” parameter.

93 **Red's parameters on the unassembled *Drosophila melanogaster* genome**

94 I obtained the short reads of a *Drosophila melanogaster* genome from the *Drosophila* 1000 genomes project.  
95 The average coverage of each genome in the project is 10 times. I used the forward reads only in this test.  
96 This data set consists of 21,806,206 76-bp-long sequences totaling 1,657,271,656 bp. The parameters of Red  
97 were adjusted to process the unassembled genome. Specifically, the word length was 16 and the minimum  
98 number of the observed words was 30 (3× the average coverage). I set the default minimum number of the  
99 observed words in an assembled genome to 3 to avoid duplicated segments. In the case of an unassembled  
100 genome, there are 30 (3× the average coverage) short reads, on average, covering each nucleotide of a segment  
101 that occurs in the genome 3 times.

102 **Executing other tools**

103 To evaluate the performance of Red, I compared it to RepeatScout, ReCon, and WindowMasker. In this  
104 study, I tried to evaluate other de-novo tools. However, RepeatScout, ReCon, and WindowMasker have been  
105 the only tools capable of processing the human genome or at least one human chromosome. The rest of the  
106 tools I tested have been unsuccessful in accomplishing the same task. The sensitivities of the tools were  
107 assessed with respect to the repeats located by RepeatMasker.

108 *RepeatMasker*

109 The ground truth consisted of repeats detected by RepeatMasker. RepeatMasker searches a sequence for  
110 instances of repeats in RepBase [1], which is a library of manually annotated repeats. BLAST was the search  
111 engine. Default values were used for all of the parameters. The following command was used for locating  
112 instances of known repeats:

113 RepeatMasker sequenceFile -s -species 'Species name' -dir outputDirectory

114 *RepeatScout*

115 RepeatScout was used with the default parameters. The program comes with some scripts to filter out  
116 simple repeats and low complexity regions. None of these filters was applied because the task at hand  
117 is the identification of all repeats, rather than the identification of transposable elements only. First the  
118 length of the word,  $k$ , is determined using the following equation that is suggested by the the inventors of  
119 RepeatScout:

120  $k = \text{round}(\log(\text{chromosomeLength}) / \log(4)) + 1$

121

122 Next, RepeatScout processes a chromosome to generate a library of repeats using the following com-  
123 mands:

124 `> build_lmer_table -l k -sequence sequenceFile -freq tableFile`

125 `> RepeatScout -sequence sequenceFile -output libraryFile -freq tableFile -l k`

126

127 Then RepeatMasker utilizes this library for locating repeats in the same sequence file using the fol-  
128 lowing command:

129 `> RepeatMasker sequenceFile -nolow -no.is -lib libraryFile -dir outDir`

130 **ReCon**

131 I used ReCon with the default settings. Hundreds of 20-kbp-long segments were sampled from all  
132 chromosomes, or contigs, of a genome. The total number of the samples collected from the genomes of the  
133 *Glycine max*, the *Zea mays*, the *Drosophila melanogaster*, the *Homo sapiens*, the *Dictyostelium discoideum*,  
134 the *Plasmodium falciparum*, and the *Mycobacterium tuberculosis* are 12970, 4443, 3633, 1703, 1476, 1156,  
135 and 220. BLAST [2] processes the samples to produce all-vs-all alignments using the following commands:

```
136 > makeblastdb -in samplesFile -dbtype nucl -out blastDbFile  
137 > blastn -db blastDbFile -query samplesFile -outfmt 6 -out blastOutFile
```

138

139 Next, the output of BLAST is further processed as recommended by the inventors of ReCon; a com-  
140 plete match between a query sequence and itself is removed. Each sequence is used as a query sequence  
141 and as a subject sequence. Thus, there are two identical pairs of matches; one of them is removed. The  
142 processed file is called ‘mspFile.’ A file including the names of the sequences is called ‘namesFile.’ ReCon  
143 processes the ‘mspFile’ using the following command:

```
144 > recon.pl namesFile mspFile 1
```

145

146 The output of ReCon is processed by Dialign [3] to construct a library of consensus sequences. Di-  
147 align aligns the longest 10 members of each family to produce a consensus sequence using a “majority vote.”  
148 The longest 10 members are written to a file called ‘familyFile.’ The consensus sequence representing a  
149 family is generated using the following command:

```
150 > dialign2-2 -n -fa familyFile
```

151

152 The consensus sequences are written to a file called ‘libraryFile.’ Finally, RepeatMasker searches for  
153 repeats in a chromosome using the following command:

```
154 > RepeatMasker sequenceFile -nolow -no_is -lib libraryFile -dir outDir
```

155 **WindowMasker**

156 WindowMasker was used with the default settings. Unlike RepeatScout and ReCon, the output of  
157 WindowMasker is the repeats themselves not a library of consensus sequences. WindowMasker processes  
158 a genome in two stages. At the first stage, it counts k-mers occurring in a genome using the following  
159 command:

160 > windowmasker -mk\_counts -fa\_list true -in chromListFile -mem 5000 -out countFile

161

162 In the second stage, WindowMasker searches for repeats in a particular chromosome. The following  
163 command executes this task:

164 > windowmasker -ustat countFile -in file -out wmFile

## 165 ***BLAST***

166 BLAST is used for validating the novel repetitive sequences discovered in the human genome. First, a  
167 database of the human genome is made using the following command:

168 > makeblastdb -in hg38.fa -dbtype nucl -out hg.db -title "Hg38"

169

170 Next, for segments of lengths of 20-49 bp, BLAST is executed with the following command:

171 > blastn -task blastn-short -db hg.db -query short.fa -evalue 1e-4 -perc\_identity 90.0 -outfmt 6 -out short.txt

172

173 Finally, for segments that have at least 50 bp, BLAST is executed with the following command:

174 > blastn -db hg.db -query long.fa -evalue 1e-4 -perc\_identity 80.0 -outfmt 6 -out long.blast

## 175 **Data**

176 Red and the related tools described above were evaluated on the sequences of seven genomes. The following  
177 list provides the version and the source of each genome:

178 • *Zea Mays* (AGPv3.21):

179 [http://ensembl.gramene.org/Zea\\_mays/Info/Index](http://ensembl.gramene.org/Zea_mays/Info/Index)

180 • *Homo sapiens* (HG38):

181 <https://www.ncbi.nlm.nih.gov/Ftp/>

182 <http://hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/>

183 • *Drosophila melanogaster* (DM6):

184 <http://hgdownload.soe.ucsc.edu/downloads.html#fruitfly>

185 • *Drosophila melanogaster* (SRR350908 - unassembled short reads):

186 <http://www.ncbi.nlm.nih.gov/sra/?term=SRR350908>

187 • *Plasmodium falciparum* (3D7):  
188 <http://www.sanger.ac.uk/resources/downloads/protozoa>

189 • *Dictyostelium discoideum*:  
190 <http://dictybase.org/Downloads/>

191 • *Mycobacterium tuberculosis* (CDC1551)  
192 <http://genome.tbdb.org>

193 • *Glycine Max* (Gmax\_109)  
194 <http://www.plantgdb.org/XGDB/phplib/download.php?GDB=Gm>

195 RepeatMasker's repeats of the human genome were obtained from the following site:  
196 <http://www.repeatmasker.org/species/hg.html>

## 197 **References**

- 198 1. Jurka J: **Repbase Update: a database and an electronic journal of repetitive elements.** *Trends Genet*  
199 2000, **16**(9):418–420.
- 200 2. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.** *J Mol Biol* 1990,  
201 **215**(3):403 – 410.
- 202 3. Al Ait L, Yamak Z, Morgenstern B: **DIALIGN at GOBICS—multiple sequence alignment using various**  
203 **sources of external information.** *Nucleic Acids Research* 2013, **41**(W1):W3–W7.