

EXTENDED PERFECT AND IMPERFECT REPEATS FINDER USING MOLECULAR DNA SEQUENCES

¹ DR. RAJU BHUKYA,

¹ Department of Computer Science and Engineering, National Institute of Technology, Warangal, INDIA

E-mail: ¹drrajunitw@gmail.com

ABSTRACT

Microsatellites are ubiquitous short tandem repeats found in all known genomes and are known to play a very important role in various studies and fields including DNA fingerprinting, paternity studies, evolutionary studies, virulence and adaptation of certain bacteria and viruses etc. Therefore, it is of importance to study distribution, enrichment and polymorphism of microsatellites in the genomes of interest. For this, the prerequisite is the availability of a computational tool for extraction of microsatellites (perfect as well as imperfect) and their related information from whole genome sequences. Examination of available tools revealed certain lacunae in them and prompted us to develop a new tool.

Keywords: *DNA Sequence, Microsatellites, Perfect, Imperfect, Extension of repeats.*

1. INTRODUCTION

Bioinformatics is a multi-disciplinary science that uses methods and principle from mathematics, computer science and statistics for analyzing biological data. DNA sequencing plays a key role in various applications in computational biology for data analysis like feature extraction, searching, disease and structural analysis. Microsatellites or simple sequence repeats (SSRs) are the nucleotide sequences arising out of tandem repeating of short sequence motifs of the size 1–6bp [18]. Microsatellites have been found in all the known genomes so far and are widely distributed both in coding and non-coding regions [1], [2]. This paper deals with repeated perfect and imperfect repeats repeating more than 10bp in number. Mutations occurring at microsatellite loci within or near certain genes have been implicated to be responsible for some human neurodegenerative diseases [18]. Furthermore, microsatellite instability has also been implicated in the induction of cancer [19]. Owing to their high mutability, it is thought that the microsatellites are one of the sources of genetic diversity [10]. Imperfect microsatellites are more stable than perfect microsatellites as they are less prone to slippage mutations [17] and are known to play a role in gene regulation [14]. A large body of microsatellite data

from several genome sequences still remains unexplored. Studies pertaining to distribution, enrichment, mutational dynamics of microsatellites along with their role in gene function and expression are very essential to understand the processes that underpin the evolution and diversity of genomes. However, a large body of microsatellite data from several genome sequences still remains unexplored. Studies pertaining to distribution, enrichment, mutational dynamics of microsatellites along with their role in gene function and expression are very essential to understand the processes that underpin the evolution and diversity of genomes.

We made a survey of existing software tools for identification and extraction of microsatellites from nucleotide sequences. these tools can be divided into the two groups: those which can identify only perfect microsatellites (e.g. SSRF [20], Poly [5], SSRIT [21]) and the others which can identify perfect as well as imperfect microsatellite (which can identify perfect as well as imperfect microsatellites (e.g. TRF [4], ATR Hunter [22], 2004), Sputnik [23], and IMEx [16]. Our survey also revealed certain ‘lacunae’ in the tools. Programs such as ‘mreps’ [11] and TandemSWAN [7] consider only substitutions but not indels. The algorithms of TRF [4], ATR Hunter [22] and

STRING [15] have been designed to find tandem repeats of large-size motifs as large as 2000 bases and hence large numbers of microsatellites go unidentified by these methods. Programs like TRbase [6] a database for tandem repeats in disease genes. Tools like TROLL [8] generates perfect and imperfect repeats, special programs like [9] for simple sequences with complex evolution, to know the differential distribution of repeats [12], microsatellite with in gene structure [13] for the analysis of gene structure. Many of these programs, [9] do not generate alignments between perfect and imperfect microsatellites. We develop new algorithm, which is fast, highly sensitive and also flexible where user can set the limit of imperfection (for perfect microsatellite and imperfect microsatellite both). The output comprises of a list of microsatellite, each of which with information content, sequences alignments (starting and ending point).

Perfect Repeat:

In a given sequence, a tandem repeat of a size *n* a given sequence, a perfect repeat of a size *n* is a subsequence which repeats continuously twice or more in the sequence (specified by the user). DNA molecules are subject to a variety of mutational events. One of the less well understood is perfect repeat duplication in which a stretch of DNA, which we call the pattern, is converted into two or more copies, each following the preceding one in a contiguous fashion. For example we could have ... TCGGA ... → ... TCGGCGGCGGA ... in which the single occurrence of triplet CGG has been transformed into three identical, adjacent copies. So here, CGG is the perfect.

TACGAGTACGGCGGCGGATGCCGTAT

Figure.1. This is three consecutive occurrence of the pattern 'CGG'.

In a given sequence, after certain intervening nucleotides, the repeat motif does not contain any imperfection (i.e. *k*=0).

TACGAGTACGGCACCGGATGCCGT

Figure 2. Here, nucleation sites characterized by two identical motifs inverted by 3 nucleotides. The intervening CAG is an iteration of CGG with *c* → *G* operation (*k*=1).

Imperfect nucleotides:

In a given sequence, imperfect repeat is the extension of the nucleation sites of the motif (with imperfection less than '*k*' value) as long as some termination criteria is satisfied. The number of imperfections between the individual repeat copy and the perfect repeat motif is more than the limit (denoted by '*k*' parameter set by the user) and (ii) the percentage of imperfection is more than the limit set by the user (denoted by '*p*' parameter). The percentage imperfection is calculated as follows:

$$p = \frac{\text{number of point mutations in the observed tract}}{\text{total number of bases in the equivalent perfect tract}} \times 100$$

The user can set the value of '*k*' between 0 to '*m*' where *m* is the repeat motif size. Once the termination criteria is satisfied, only those candidate microsatellites that are more than the minimum repeat number of that repeat size set by the user (denoted by '*n*' parameter) are reported.

2. METHODS

2.1 Discovering perfect repeats (Exact Repeats):

The process to discover perfect Repeats in the given sequence file consists of two phases: Initially start by looking for a subsequence of length one (mono nucleotide) and check for its continuous repetition. If there is a repetition of repeats, then increase the value of count for every repeats till the same repeat is found and if the count value is equal or greater than the user value (denoted by the '*n*'), then write the sequences into the file for every repeat. Follow the procedure for the next subsequence from the character just after the ending index. If a subsequence of length one (mono nucleotides) does not repeat continuously (given value by the user for mono nucleotides) up to end of the file, then increase the length of the subsequence (means Di nucleotides) and search for the repetition from starting point of the file to end of the file. Repeat the same procedure up to deca nucleotide. Always compare the count value to the user specified value if the count value is greater than or equal to then write the sequence into the file otherwise the algorithm goes for the next sequence and so on.

2.2 Discovering imperfect Repeats (Approximate perfect repeats):

The process to discover approximate perfect repeats (imperfect repeats) in the given sequence file.

(i) The no of imperfections between the individual repeat copy and the perfect repeat motif is more then, the limit (denoted by the 'k' parameter) then we can write into the file for that repeat.

(ii) The percentage of imperfection is calculated for every repeats. If the percentage of imperfection is less than the limit set by the user (denoted by 'p' parameter) then we can write into the file for that repeat.

Initially start by looking for a subsequences of length one (mono nucleotide) and check for its continuous repetition with k -mismatched ('k' value set by the user 1,2,3,...). If there is a repetition of repeats, then look for the index till where the subsequences repeats itself continuously with k -mismatched. Here we compare the count value with user value (specified by the user 'n'). For finding imperfect repeats, use normal string matching algorithm. The algorithm stores starting & ending indices for mono-imperfect repeats in a file and follow the procedure for the next subsequences from the character just after the ending index. If a subsequence of length mono does not repeat continuously, then increase the length of the subsequence to di-nucleotides and search for the repetition of all subsequences of DNA file starting point to end point of the DNA file and so on.

The process of finding exact compound tandem repeats is depicted in the Figure 3.

3. RESULTS:

Discovering perfect and imperfect repeats of proposed algorithm technique are implemented in Python programming language. For experiment, we used genome sequences for discovering perfect as well as imperfect repeats. The proposed algorithm finds the perfect repeat which is able to discover up to 20 or more in size. Here the technique which is used is simple string matching algorithm for finding perfect and imperfect repeats. To discuss the capabilities of our code, we analyzed the human atrophin1 gene (BC051795) and compared the result obtained with those obtained using Tandem Repeat Finder (TRF) [4] and Sputnik [23] and IMEx [16]. TRF was initially tested with the parameters used in the earlier studies [2], [6],

[24] which yielded very few microsatellites. Hence, we used the most relaxed set of parameters (Match: +2, Substitution: -7, Indel: -7, min Score: 2) which yielded substantial number of microsatellites. This is because using TRF, the length of the microsatellite

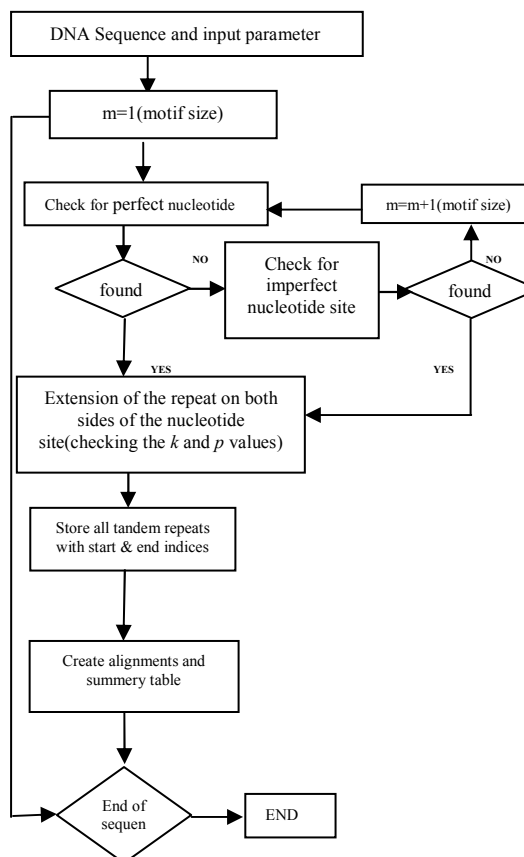


Figure 3: Process for Exact Compound Tandem Repeats

we used the least stringent parameters (Match: +1, Mismatch: -3, Min Score: -5), and for IMEx [16], we set the 'p' value of all tracts to 10%; 'k' value for each pattern size: Mono: 1, Di: 1, Tri: 1, Tetra: 2, Penta: 2, Hexa: 3 and further restricted to report only those microsatellites with minimum repeat copy number (Mono:5, Di: 3, Tri: 2, Tetra: 2, Penta: 2, Hexa: 2) to match those reported by TRF [4] and Sputnik [23]. TRF [4] and Sputnik [23] identified 50 and 19 repeats respectively, whereas IMEx [16] identified 146 microsatellite tracts. In our program we take the input as same as IMEx [16] but the program restricted to report only those microsatellite with minimum repeat copy number (Mono:5, Di: 3, Tri: 2, Tetra: 2, Penta: 2, Hexa: 2, Octa: 2, Enea: 2, Deca:2) to match those reported by TRF [4], Sputnik [23] and IMEx [16]. IMEx



| | | | | | |
|-----|-----------|------|------|---|-------------------|
| 203 | GAGCGC | 2727 | 2738 | 2 | GAGCGGAGCGC |
| 204 | GCACCT | 3428 | 3439 | 2 | GCACTGCACCT |
| 205 | AAGGCC | 3987 | 3998 | 2 | AAGCCAAAGGCC |
| 206 | AACCAA | 4321 | 4332 | 2 | AACCAAAACCAA |
| 207 | AAAAAA | 4361 | 4378 | 3 | AAAAA 3 times |
| 208 | CCCTTC | 3212 | 3225 | 2 | CCCTTCCCTTTC |
| 209 | TGCCCT | 4247 | 4260 | 2 | TGCCCTTGCCCT |
| 210 | AAAAAA | 4361 | 4381 | 3 | AAAAA 3 times |
| 211 | CTTCCAG | 1442 | 1457 | 2 | CTTCCAGCTTCCAG |
| 212 | AAAAAAA | 4361 | 4376 | 2 | AAAAAAA 2 times |
| 213 | CAGCAGCAG | 1710 | 1754 | 5 | CAGCAGCAG 5 times |
| 214 | AAAAAAA | 4361 | 4378 | 2 | 9 As 2 times |
| 215 | AAAAAAA | 4361 | 4380 | 2 | 10 As 2 times |

As we can see from the result, our program reports many more tracts which are missed by the other three program (TRF [4], Sputnik [23], IMEx [16]). It is important to mention that Sputnik does not report the mononucleotides and after that the penta nucleotides(i.e. hexa nucleotide and more) and IMEx also not report the some of the repeats size mononucleotide to hexa nucleotides and it does not report after the hexa nucleotides to deca nucleotide and more. The IMEx tool dosent reports some of the repeats which is starting from mono to hexa nucleotide repeats which is reported by the proposed algorithm. Our program reports all the possible repeats from mono nucleotides to deca nucleotide. Using our algorithm we can show repeats up to 25, up to 50, up to 75 and up to 100. For 25 it generates 221 repeats, for 50 it generates 221 and so on. We also run the program on four whole genome sequences Plasmodium falciparum chromosome IV (NC_004318.1), yeast chromosome IV (NC_001136.8), Mycobacterium tuberculosis H37Rv genome (NC_000962.2) and E.coli K12 genome (NC_000913.2). TRF uses the probabilistic algorithm which includes a 'detection step' to identify the candidate repeats and an 'analysis' step that uses different statistical criteria to filter the candidate repeats. Sputnik uses a recursive algorithm and the performance depends on the recursion depth of the program. Hence, Sputnik's execution time seems to be dependent on the sequence composition. IMEx uses the simple string-matching algorithm that scans the entire sequence using sliding window approach and reports the results in a single run. Hence, the processing time of IMEx is dependent on the length of the DNA sequence and not on the number of microsatellites. Our program uses the simple string matching algorithm and report the results in a single run. Hence processing time is not dependent on the DNA sequence file, it depends upon the size of the motif. The program has been designed keeping in view of the limitations we encountered

with the other available tools. Using the proposed program the user can search perfect microsatellite and also imperfect microsatellite, generate alignments, set the imperfection percentage threshold of the entire tract of each repeat size, search the repeats of a particular size. From the result, we generate more number of repeats as shown in above Table.1.

Table 2: Comparison of execution times(in seconds) of TRF, Sputnik, IMEx, and proposed program

| Sequences | TRF | | Suptnik | | IMEx | | EPI | |
|-------------------------|---------|------|---------|-------|---------|------|---------|-------|
| | Repeats | Time | Repeats | Time | Repeats | Time | Repeats | Time |
| Plasmodium Chr4(1204Kb) | 25601 | 69.8 | 10810 | 89.1 | 54232 | 2.9 | 111695 | 5.69 |
| Yeast Chr4 (1531 Kb) | 7308 | 4.4 | 2831 | 287.2 | 39759 | 4.0 | 54768 | 7.35 |
| MTB H37Rv (4411Kb) | 16439 | 25.5 | 9412 | 17.7 | 111113 | 11.6 | 131290 | 21.03 |
| E.coli K12 (4639 Kb) | 12043 | 8.8 | 5387 | 8.5 | 105392 | 12.3 | 129229 | 22.09 |

TRF: Match: ≥ 2 Subs: 8 Indel: 8 Min. Score: 20 pM: 0.80 pI: 0.10
 Max.Period:6.Sputnik: Match: ≥ 2 Mismatch: 6 Min. Score: 8.
 IMEx: 'k' value: Mono: 1, Di: 1, Tri: 1, Tetra: 2, Penta: 2, Hexa: 3; 'p' value:10% for all repeat sizes; repeat length: 10 bases or more XXXX: 'k' value: Mono: 1, Di: 1, Tri: 1, Tetra: 2, Penta: 2, Hexa: 3; 'p' value:10% for all repeat sizes; 'n' value : Mono: 5, Di: 3, Tri: 2, Tetra: 2, Penta: 2, Hexa: 3 repeat length: 10 bases or more.

The above comparison shows that our program reports more number of repeats than the mentioned three tools (TRF, Sputnik, and IMEx). So we can say that our program is more efficient, more accurate and flexible than the other ones.

4. CONCLUSION:

In this paper we have presented a new algorithm for finding perfect and imperfect microsatellite repeats in DNA sequences. In this first we have found all perfect repeats in DNA sequences and then also we have to find imperfect. And then we have stored them in to a text file. for finding tandem repeats is wide ranging and non-standardized. One has to be careful in understanding the tools' inherent constraints to select the right tool for the right purpose. It is hence important to be able to compare the repeat search tools and understand their behavior and inherent limitations.

REFERENCES:

- [1] Anwar,T. and Khan ,A.U. (2006) SSRscanner: a program for reporting distribution and exact



- location of simple sequence repeats
Bioinformatics, 1, 89–91.
- [2] Archak, S. et al. (2007) InSatDb: a microsatellite database of fully sequenced insect genomes. *Nucleic Acids Res.*, 35, D36–D39.
- [3] Ross, Charles L., et al. "Rapid divergence of microsatellite abundance among species of *Drosophila*." *Molecular Biology and Evolution* 20.7 (2003): 1143-1157.
- [4] Benson, G. (1999) Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.*, 27, 573–580.
- [5] Bizzaro, J.W. and Marx, K.A. (2003) Poly: a quantitative analysis tool for simple sequence repeat (SSR) tracts in DNA. *BMC Bioinformatics*, 4, 22.
- [6] Boby, T. et al. (2005) TRbase: a database relating tandem repeats to disease genes in the human genome. *Bioinformatics*, 21, 811–816
- [7] Boeva, V. et al. (2006) Short fuzzy tandem repeats in genomic sequences, identification, and possible role in regulation of gene expression. *Bioinformatics*, 22, 676–684.
- [8] Castelo, A. et al. (2002) TROLL – Tandem repeat occurrence locator. *Bioinformatics*, 18, 634–636.
- [9] Ellegren, H. (2004) Microsatellites: simple sequences with complex evolution. *Nat. Rev. Genet.*, 5, 435–445.
- [10] Kashi, Y. and King, D.G. (2006) Simple sequence repeats as advantageous mutators in evolution. *Trends Genet.*, 22, 253–259.
- [11] Kolpakov, R. et al. (2003) mreps: efficient and flexible detection of tandem repeats in DNA sequences. *Nucleic Acid Res.*, 31, 3672–3678.
- [12] Katti, M.V. et al. (2001) Differential distribution of simple sequence repeats in eukaryotic genome sequences. *Mol. Biol. Evol.*, 18, 1161–1167.
- [13] Li, Y.C. et al. (2004) Microsatellites within genes: structure, function, and evolution. *Mol. Biol. Evol.*, 21, 991–1007.
- [14] Meloni, R. et al. (1998) A tetranucleotide polymorphic microsatellite, located in the first intron of the tyrosine hydroxylase gene, acts as a transcription regulatory element in vitro. *Hum. Mol. Genet.*, 7, 423–428.
- [15] Parisi, V. et al. (2003) STRING: finding tandem repeats in DNA sequences. *Bioinformatics*, 19, 1733–1738.
- [16] Suresh B. et al. (2007) IMEx: Imperfect Microsatellite Extractor *Bioinformatics*, 1181–1187.
- [17] Sturzeneker, R. et al. (1998) Polarity of mutation in tumor-associated microsatellite instability. *Hum. Genet.*, 102, 231–235.
- [18] Tautz, D. and Schlotterer, C. (1994) Simple sequences. *Curr. Opin. Genet. Dev.*, 4, 832–837.
- [19] Thibodeau, Stephen N., G. Bren, and D. Schaid. "Microsatellite instability in cancer of the proximal colon." *Science* 260.5109 (1993): 816-819.
- [20] Sreenu VB, et al. MICAS: a fully automated web server for microsatellite extraction and analysis from prokaryote and viral genomic sequences. *Appl. Bioinformatics* 2003;2:165-168.
- [21] Temnykh S, et al. Computational and experimental analysis of microsatellites in rice (*Oryza sativa* L.): frequency, length variation, transposon associations, and genetic marker potential. *Genome Res.* 2001;11:1441-1452.
- [22] Wexler Y, et al. RECOMB 2004. 2004. Finding approximate tandem repeats in genomic sequences.
- [23] Abajian C. Sputnik. <http://espressosoftware.com/pages/sputnik.jsp>.