

MMsat—a database of potential micro- and minisatellites

Andrew Shelenkov ^{*}, Alexander Korotkov, Eugene Korotkov

Bioinformatics Department of Bioengineering Centre of Russian Academy of Sciences, Prospect 60-tya Oktyabrya, 7/1, Room 303, 117312, Moscow, Russia

Received 20 July 2007; received in revised form 8 October 2007; accepted 16 November 2007

Available online 28 November 2007

Received by M. Di Giulio

Abstract

We present MMsat—a database of DNA sequences from GenBank possessing the latent periodicity at high level of statistical significance and having the period length in a range from 2 to 100 bases. The periodicity was found by analytical method of information decomposition. These sequences can be considered as potential micro- and minisatellites and thus can be useful for PCR analysis and evolutionary studies. Distribution, properties, and potential functions of periodicity are discussed.

Availability: <http://victoria.biengi.ac.ru/mmsat>

© 2007 Elsevier B.V. All rights reserved.

Keywords: Bioinformatics; Databases; Nucleic acid; Information decomposition; Latent periodicity

1. Introduction

The presence of repeated sequences is a common feature for both eukaryotic and prokaryotic genomes. It has been suggested that the repeats themselves produce unusual physical structures in the DNA, causing polymerase slippage and the resulting amplification (Weitzmann et al., 1997; Wells, 1996). The other potential role for tandem repeats is gene regulation, in which the repeats may interact with transcription factors, alter the structure of the chromatin or act as protein binding sites (Richards et al., 1993; Lu et al., 1993). Also, repeat regions are often found in coding regions (Tomba, 2003), so they are directly involved in genome functioning. When they fall in regulatory regions, they may have direct influence on phenotype (Fondon et al., 2004). In the last few years, tandem repeats have been increasingly recognized as markers of choice for genotyping a number of pathogens (Keim et al., 2000; Frothingham and Meeker-O'Connell, 1998; Supply et al., 2000). The rapid evolution of

these structures appears to contribute to the phenotypic flexibility of pathogens. Further, the studying of repeat regions is important for population genetic and forensic applications as well (Estoup et al., 2001; Blouin et al., 1996).

Tandem repeats are usually classified among satellites (spanning megabases of DNA, associated with heterochromatin), minisatellites (repeat units in the range 6–100 bp, spanning hundreds of base-pairs) and microsatellites (repeat units in the range 2–5 bp, spanning a few tens of nucleotides) (Le Fleche et al., 2001). Although the perfect tandem repeats can be found in genomes of different organisms, it is possible, especially in bacterial genomes, to find the very old or “ancient” microsatellites possessing very fuzzy periodicity, so they can be passed through by the mathematical methods of tandem repeat finding (Korotkov et al., 2003). Yet these ancient sequences are of great biological interest since they are usually highly polymorphous and thus can be used as genetic markers (Le Fleche et al., 2001; van Belkum et al., 1997; Adair et al., 2000).

We used the method of information decomposition (Korotkov et al., 2003) to make a search for periodic sequences through the whole GenBank database. That is, we have scanned GenBank using the software developed by us and put the results into our database. Information decomposition (ID) is a spectrum representing the statistical significance of mutual information

Abbreviations: ID, information decomposition; PCR, polymerase chain reaction.

^{*} Corresponding author. Tel./fax: +7 499 135 2161.

E-mail address: fallandar@gmail.com (A. Shelenkov).

Download English Version:

<https://daneshyari.com/en/article/2819464>

Download Persian Version:

<https://daneshyari.com/article/2819464>

[Daneshyari.com](https://daneshyari.com)