

REPuter: fast computation of maximal repeats in complete genomes

Stefan Kurtz and Chris Schleiermacher

Technische Fakultät, Universität Bielefeld, Postfach 10 01 31, D-33501 Bielefeld, Germany

Received on January 11, 1999; revised on February 22, 1999; accepted on February 23, 1999

Abstract

Summary: A software tool was implemented that computes exact repeats and palindromes in entire genomes very efficiently.

Availability: Via the Bielefeld Bioinformatics Server (<http://bibiserv.techfak.uni-bielefeld.de/reputer/>).

Contact: {kurtz,icschlei}@techfak.uni-bielefeld.de

Introduction

Computer scientists have developed methods to locate repeated substrings of different kinds, dating back to the pioneering work of Martinez (1983). Several software tools have been developed to locate repeated substrings, for example Devereux *et al.* (1984), Agarwal and States (1994) and Rivals *et al.* (1997). However, to our knowledge all available software tools for repeat analysis have strict limits on the maximal length of the input sequence they allow to process. For example, the repeat-finder of the GCG-package (version 7.0) only allows input sequences of length up to 350 000 bases.

We have developed a software tool *REPuter* which allows to determine all exact repetitive substrings contained in complete genomes. Exact repeats are only a small fraction of all repeats of biological interest. However, since exact repeats usually form the core blocks of approximate repeats, *REPuter* can also serve as a fast subroutine for programs that detect these, see for example Leung *et al.* (1991). The running time and space requirement of *REPuter* is linear in the length of the genome and in the size of the output. The method we use is thus asymptotically optimal. The main advantage over previous programs is the reduced time and space complexity of our tool, and the guaranteed linear worst case behavior. This allows one to handle much longer input sequences. The current version is able to process input sequences consisting of up to 67 million bases. For example, to compute all maximal repeats of length at least 20 contained in the *S. cerevisiae* genome (12 147 818 bases) takes about 46 s on a Pentium 350 MHz-computer, using 160 Mbyte of space.

REPuter consists of two programs, a search engine and a visualizing component. The search engine processes a DNA sequence given by the user in Fasta-format, and returns a representation of all maximal repeats in a simple ASCII-format. The visualizing component processes the output of the search engine

and produces an overview of the number, the length, and the location of the repeated substrings. *REPuter* is available via the WWW on the Bielefeld-Bioinformatics Server (see above). For user convenience we have precomputed the maximal repeats for some genomes. The search engine can be downloaded as an executable binary for several platforms.

Methods

Consider some sequence s over the DNA alphabet. A *repeat* is a substring in s which occurs at least twice. Suppose w is a repeat of length l in s which occurs at the starting positions i and j , i.e. $s_i \dots s_{i+l-1} = s_j \dots s_{j+l-1}$. Let $i < j$. Then we say that (l, i, j) is a *forward repeat* of length l , F-repeat, for short. Note that w is contained in a longer repeat if the bases to the left or to the right of both occurrences of w are identical. To reduce redundancy, we restrict attention to maximal F-repeats: (l, i, j) is *maximal* if $(s_{i-1} \neq s_{j-1})$ and $(s_{i+1} \neq s_{j+1})$ (whenever these positions exist in s). Note that each F-repeat in s can easily be deduced from a maximal F-repeat. Maximal palindromes are defined in a similar way: (l, i, j) , $i \leq j$, is a *palindromic repeat*, P-repeat, for short, if $s_i \dots s_{i+l-1} = \overline{s_j \dots s_{j+l-1}}$, where \overline{w} denotes the reverse complement of a DNA-sequence w . (l, i, j) is *maximal* if the complement of base s_{i-1} and s_{j-1} is different from s_{j+1} and s_{i+1} , respectively, whenever these positions exist in s .

Example 1

gtcaca contains one maximal F-repeat (2,2,4) (*ca*) and one maximal P-repeat (2,0,3) (*gt*), of length ≥ 2 .

For a given threshold $l > 0$, our search engine computes all maximal repeats of length at least l , using a variation of an algorithm described in Gusfield (1997). Two phases are necessary to deliver all maximal F- and P-repeats for a sequence s :

In the first phase, the suffix tree for the sequence $t = xsy\overline{sz}$ is computed. x , y , and z are unique symbols not occurring in s . They allow to conveniently handle boundary cases. The suffix tree is a well known data structure in string processing (McCreight, 1976). It can be computed in $O(n)$ time and space, where n is the length of t . The tool of Rivals *et al.* (1997) is also based on suffix trees, but it computes non-overlapping F-repeats, and the worst case running time is quadratic.

