

## Sequence analysis

**SITEBLAST—rapid and sensitive local alignment of genomic sequences employing motif anchors**

Morris Michael\*, Christoph Dieterich† and Martin Vingron

Computational Molecular Biology, Max Planck Institute for Molecular Genetics, Ihnestraße 63-73, D-14195 Berlin, Germany

Received on November 10, 2004; revised on December 6, 2004; accepted on December 13, 2004

Advance Access publication December 14, 2004

**ABSTRACT**

**Motivation:** Comparative sequence analysis is the essence of many approaches to genome annotation. Heuristic alignment algorithms utilize similar seed pairs to anchor an alignment. Some applications of local alignment algorithms (e.g. phylogenetic footprinting) would benefit from including prior knowledge (e.g. binding site motifs) in the alignment building process.

**Results:** We introduce predefined sequence patterns as anchor points into a heuristic local alignment strategy. We extended the BLASTZ program for this purpose. A set of seed patterns is either given as consensus sequences in IUPAC code or position-weight-matrices. Phylogenetic footprinting of promoter regions is one of many potential applications for the SITEBLAST software.

**Availability:** The source code is freely available to the academic community from <http://corg.molgen.mpg.de/software>

**Contact:** christoph.dieterich@molgen.mpg.de

**1 INTRODUCTION**

We expand on the idea of heuristic alignment algorithms like BLAST (Altschul *et al.*, 1997). These tools follow a strategy where a collection of ‘matching’  $n$ -mers defines anchor points for building local pairwise alignments. Especially in the context of phylogenetic footprinting, it is often desirable to include prior information to guide the alignment building process. Comparative approaches in the domain of promoter analysis will benefit from an approach where alignments are extended from seed pairs whose identity is known a priori. In the context of promoter analysis, prior information is given as consensus patterns or weight matrices, which model a regulatory unit (e.g. a transcription factor binding site). The main idea is to generalize the rigid concept of seeds from identical or similar  $n$ -mers to a controlled set of seeds. We implement our concept on top of the existing rapid alignment solutions, namely the BLASTZ software by Schwartz *et al.* (2003).

**2 APPROACH**

Initially, our SITEBLAST software identifies all potential seeds. The user has two possibilities to specify the patterns that are considered as seeds. First, a list of consensus sequences in IUPAC code can be

given together with a distance  $D$ . Each occurrence of an  $n$ -mer, which matches a consensus sequence with at most  $D$  errors is considered as seed. Second, a list of position weight matrices (PWMs) can be used. Here, the user is free to set constraints on the proportion of false positives ( $P$ -value) or true positives (power). To avoid unspecific matrix matches, there is an option to set a power-Limit given a  $P$ -value cut-off or a  $P$ -value limit given a power cut-off. Matrices that do not meet these criteria are excluded from the seed finding process. For example, a matrix with power  $t$  and a preset cut-off level  $T$  (power-Limit) for some fixed  $P$ -value cut-off is selected only if  $t \geq T$ . Seeds may occur in any orientation (also reversed, complemented or reverse complemented).

After all possible seeds in both sequences are found, alignments are computed by BLASTZ (Schwartz *et al.*, 2003). Therefore, each seed in the first sequence for each consensus sequence or PWM is paired with each equivalent occurrence in the second sequence, respectively. These seed pairs are used as anchor points for BLASTZ.

Computed alignments are annotated with all matching seed pairs.

**2.1 Finding the seeds**

**2.1.1 Specified by consensus sequences** If seed patterns are given by a list of IUPAC consensus sequences, the user can choose between two different search strategies. The first one needs no preprocessing and scans the two sequences for matches of any consensus sequences in any orientation in a trivial way. This option is faster for short sequences and high distance  $D$ . The second strategy generates all matching patterns for all consensus sequences and inserts them in all desired orientations in an Aho–Corasick pattern matching machine (Aho and Corasick, 1975). This matching machine is used to rapidly retrieve all seeds in the two input sequences. This is considerably faster for long sequences and low distance  $D$ .

**2.1.2 Specified by PWMs** Two position specific score matrices (PSSMs) are computed for each PWM (for details see Rahmann *et al.*, 2003). One PSSM is tailored to the GC-content of the first sequence and another to the GC-content of the second sequence. We compute the score threshold  $S$  for a fixed  $P$ -value or power respectively. Then, the attainable power or  $P$ -value is computed to test whether they match a given power or  $P$ -value limit. For these computations, we employ a modified version of the PATSER (version 3b) algorithm by Hertz and Stormo (1999). If the two PSSMs do not meet the power or  $P$ -value limits, the pair is discarded. With the remaining PSSMs the two input sequences are scanned. All corresponding

\*To whom correspondence should be addressed.

†Present address: Department of Computer Science, PO Box 68 (Gustav Hällströminkatu 2b), FIN-00014 University of Helsinki, Finland.

