LTR retrotransposon detection via deterministic finite automata

Author: Lucie Klímová Supervisor: doc. Mgr. Lukáš Holík Ph.D.

MOTIVATION

LTR retrotransposons make up a significant part of the human genome (8.3%)

They can influence gene expression (the amount of protein that is synthesized)

They are highly nested and therefore hard to detect



TE-GREEDY NESTER



- Recursively removes the best matching LTR elements
- Due to the recursion, the algorithm appears relatively slow (80% of its runtime is taken up by BLASTX)

Figure 3 - Ilustration of transposon nesting

Figure 4 - Representation of the TE-greedy nester algorithm

FASTA

LTR FINDER RETURNED 03

RUN LTR

FINDER

SCORE CANDIDATES

COUNT INSERTIONS ADJUST POSITIONS AND FRAGMENTED

FEATURES

EXPORT GFF FILE REMOVE BEST NON-

OVERLAPPING LTR TEs

ABOVE TRESHOLD

SCORE

Precisely models the desired sequence using match, insert and delete states

Non-deterministic, time complexity O(LM²)

PROFILE HMM



Figure 5 - Profile hidden Markov model structure

DOMAIN SEARCH USING DFAs

- The main idea is to transform a profile HMM into a deterministic model while minimizing the loss of accuracy
- Direct determinization is unfeasible, so the model must be simplified
- A profile HMM can be converted into a bounded counting automaton (BCA), which uses a counter to determine how many consecutive transitions below a threshold *t* were taken





Figure 6 - Process of generating the deterministic model

Figure / - Execution time comparison



Figure 8 - Sensitivity comparison



