

Hairpins in a Haystack

Jana Hertel

Bioinformatics Group,
Department of Computer Science and Interdisciplinary Center for Bioinformatics,
University Leipzig

Leipzig, February 22, 2006

Outline

Introduction

- Background

- Detection

- Basis

- Purpose

Methods

- Hairpin filter

- Descriptors and SVM

- SVM training

Application

- Homo sapiens

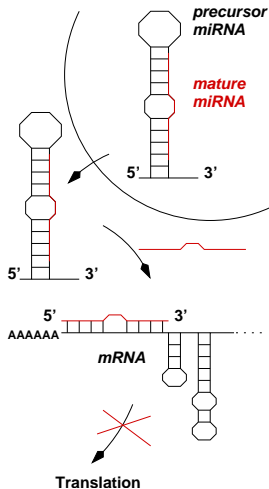
- Nematodes

- Seasquirts

Discussion

MicroRNAs - Background

- ▶ class of noncoding RNAs
- ▶ important regulatory functions
- ▶ longer transcripts (pre-miRNAs)
→ ~ 100nt
- ▶ functional mature miRNA in one stem side → ~ 22nt
- ▶ mature miRNA highly conserved
- ▶ bind to 3'UTRs of mRNA targets
 - ▶ suppress expression
 - ▶ mark for degradation



MicroRNAs - Detection

- ▶ candidates homologous to miRNAs
- ▶ candidates adjacent to known miRNAs

Problem:

- ▶ candidates that do not feature these facts can not be found

MicroRNAs - Detection

- ▶ several approaches for detecting novel miRNA genes
- ▶ secondary structure, 3' and 5' patterns in stem loop

Examples:

- ▶ miRscan¹ (nematodes), miRseeker² (insects), miralign³ (vertebrates)
- ▶ candidate search, classifying by features partly machine learning

¹Lim *et al.* 2003

²Lai *et al.* 2003

³Wang *et al.* 2005

Basis of this approach

- ▶ genome-wide screens for ncRNAs
- ▶ RNAz^4 → evolutionary conserved secondary structure in multiple sequence alignments
- ▶ automatic tools necessary to assign candidates to ncRNA classes

⁴Washietl *et al.* 2005

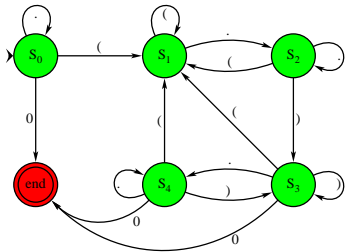
Purpose of RNAmicro

- ▶ RNAmicro works as "sub-screen" for ncRNA surveys
- ▶ SVM based tool for classifying miRNAs among ncRNAs

Particular steps

1. detect *almost* hairpins
2. computation of descriptors
3. SVM classification

Detecting *almost* hairpins



- ▶ specify window of alignment
- ▶ consensus sequence and structure
- ▶ dot-parantheses string read by automaton
- ▶ start and length of each stem loop stored
- ▶ accept structure if
 - ▶ exactly 1 stem loop $> 10\text{nt}$
 - ▶ other smaller stem loop $\leq 4\text{ nt}$

Computation of Descriptors

stem length

loop length

G+C content

z-score

mean single mfe

adjusted mfe

mfe index

consensus mfe

3'/5' stem, loop entropy

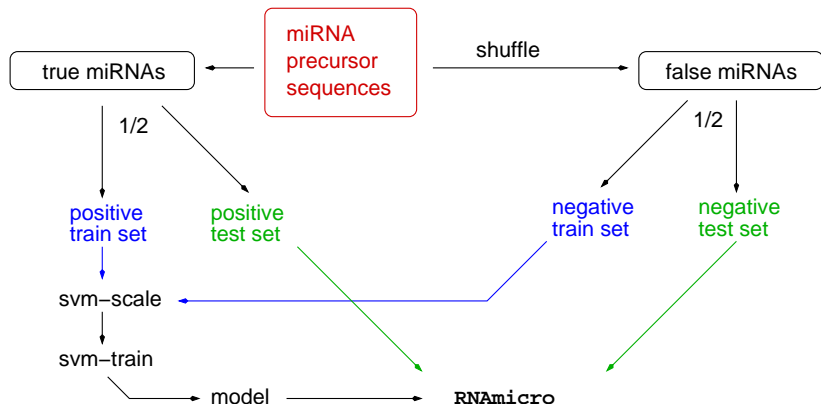
23nt block entropy

SVM implementation

- ▶ SVM from `libsvm`⁵
- ▶ scale descriptors to $[-1, +1]$
- ▶ model: rbf kernel, probability estimates

⁵Chang and Lin, 2001

Initial training



Results of initial training

- ▶ 134/147 (90%) sensitivity, 381/383 (99%) specificity
- ▶ train SVM again with entire datasets
- ▶ test on *RNAz* screens of nematodes and seasquirts
 - significant number of known ncRNAs false classified
 - initial negative set not sufficiently good

Retraining the SVM

- ▶ ncRNA alignments extracted from Rfam database
- ▶ add known false positives to negative train set
- ▶ iterate process of adding false positives and retraining until no significant improvement on the Rfam dataset

Data

- ▶ vertebrate genomes⁶, nematode and urochordate⁷ ncRNA alignments
- ▶ screen with 70,100 and 130 nt window
- ▶ retaining best (p!) non-overlapping hits of each alignment

⁶Washieta *et al.* 2005

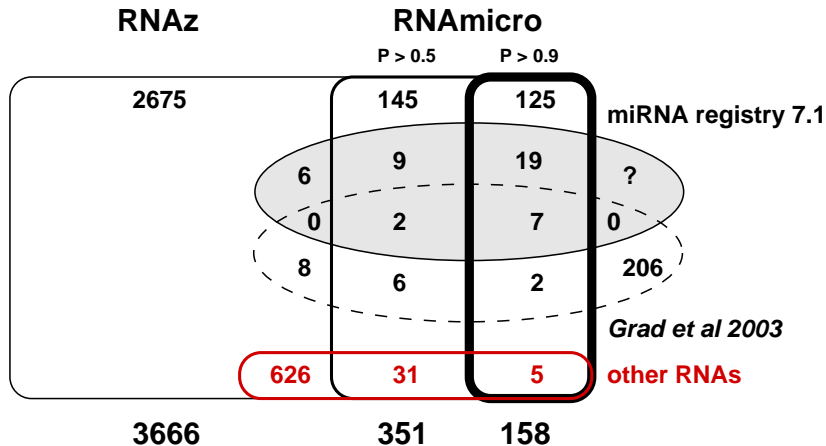
⁷Missal *et al.* 2005

Verification

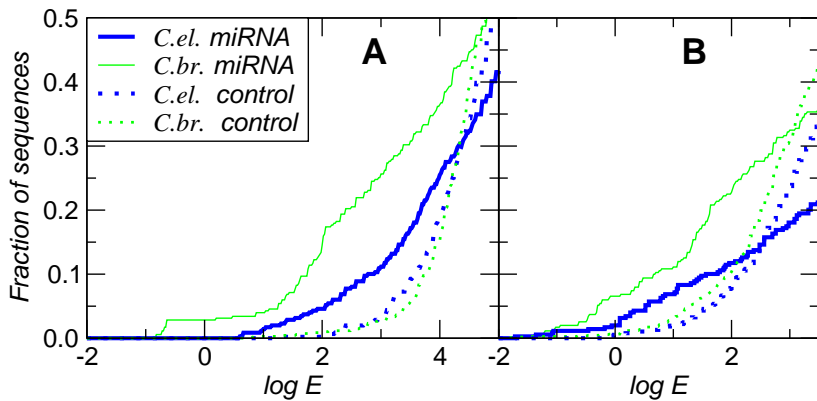
- ▶ re-evaluate `RNAmicro` candidates with other SVM approach⁸
- ▶ very restrictive hairpin filter
 - ▶ 3077 / 5440 with $P > 0.5$ passed filter, 1590 recognized
 - ▶ 953 / 1481 with $P > 0.9$ passed filter, 657 recognized
- ▶ 4245 / 5440 candidates not associated with protein coding genes
- ▶ 1107 candidates located within introns (36 known)

⁸Xue *et al.* 2005

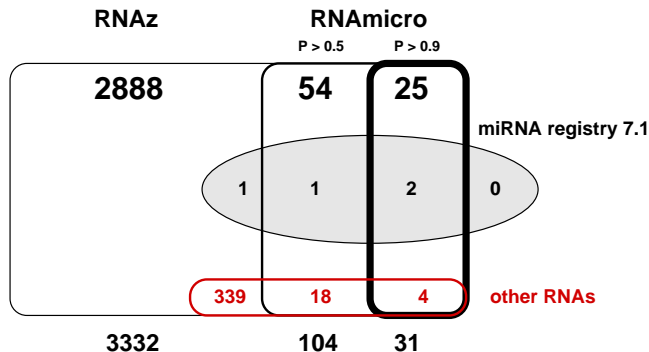
General results



Upstream motif



General results



- ▶ 5 clusters with 10 members

Summary

- ▶ RNAmicro designed for classification of ncRNA alignments
- ▶ applied to 3 recent RNAz based studies
- ▶ large number of novel miRNA candidates
- ▶ verification through
 - ▶ comparison with other approaches and annotations
 - ▶ analysing genomic location to other candidates
 - ▶ location in introns
- ▶ large number of RNAmicro predictions correspond to real miRNAs
- ▶ only small fraction of true miRNA repertoire known