MiGaL

RNA secondary structure modelling and comparison.

J. Allali & M.F. Sagot

IGM, Marne la Vallée, INRIA Rhône-Alpes



Introduction about RNA.

- Introduction about RNA.
- Previous Work.

- Introduction about RNA.
- Previous Work.
- MiGaL: a new modelling scheme.

- Introduction about RNA.
- Previous Work.
- MiGaL: a new modelling scheme.
- An algorithm for the "fusion" of two ordered trees.

- Introduction about RNA.
- Previous Work.
- MiGaL: a new modelling scheme.
- An algorithm for the "fusion" of two ordered trees.
- Conclusion.

Introduction

MiGaL: introduction



MiGaL: introduction

RNA folding:



Base pairing

- Watson Crick:(A-U) (C-G)
- Wobble: (G-U)
- non-canonical: Holbrook, 91 (C-U) Noller, 84 (G-A)



Helix.



Helix.Hairpin Loop.



- Helix.
- Hairpin Loop.
- Multiloop.



- Helix.
- Hairpin Loop.
- Multiloop.
- Bulge.



- Helix.
- Hairpin Loop.
- Multiloop.
- Bulge.
- Internal Loop.



- Helix.
- Hairpin Loop.
- Multiloop.
- Bulge.
- Internal Loop.
- Stem.



- Helix.
- Hairpin Loop.
- Multiloop.
- Bulge.
- Internal Loop.
- Stem.
- Pseudo knot.

RNA secondary structure comparison taking into account pseudo knots.

Motif inference for RNA secondary structures.

RNA secondary structure comparison taking into account pseudo knots.

Motif inference for RNA secondary structures.

we need a data structure for RNA secondary structures

Previous Work

MiGaL: parenthesis sequence



Also called "arc annotated sequencies"

MiGaL: parenthesis sequence



Zhang (99) NP-Hard for crossing pairs.



Several different representations using trees.



internal node=base pair leaf=unpaired base



Zhang (90)

internal node=helix,bulge,internal loop or multiloop leaf=unpaired bases



Shapiro (89)

internal node=helix,bulge,internal loop or multiloop leaf=hairpin loop By adding pseudo knots to this tree structure, we change it into a graph.

Comparison, edition . . . are difficult problems with graphs.

MiGaL: grammar

Stochastic Context Free Grammars : Sakakibara (93)

 $S \to (S)|()|s|SS$

Used for alignment or folding algorithms.

S-attribute Grammars Lefebvre (96)

Ideas:

• Start from a tree representation.

Ideas:

- Start from a tree representation.
- Add edges for pseudo knots.

Ideas:

- Start from a tree representation.
- Add edges for pseudo knots.
- Use simultaneously various levels of representation.

Ideas:

- Start from a tree representation.
- Add edges for pseudo knots.
- Use simultaneously various levels of representation.

- Multiloop network is the backbone of the RNA structure.
- Nucleotide conservation is not necessary.

MiGaL: MultIGrAphLayer



- Detailed structure at the deepest level.
- Abstraction : the top level.
- Relation between adjacent layers.

MiGaL: On RNA



























Our model takes into account pseudo knots.

The requested amount of memory is linear in the size of the RNA.

MiGaL can represent either a model for various RNAs or a single RNA.

Our model takes into account pseudo knots.

The requested amount of memory is linear in the size of the RNA.

MiGaL can represent either a model for various RNAs or a single RNA.

We now need an algorithm to compare (two) MiGaLs

Our model takes into account pseudo knots.

The requested amount of memory is linear in the size of the RNA.

MiGaL can represent either a model for various RNAs or a single RNA.

We now need an algorithm to compare (two) MiGaLs

We now need an algorithm to compare (two) Layer0s

An *Edit-like* algorithm based on our observations concerning RNA multiloop network and adapted to structural tree matching (level 0,1 and 2).

MiGaL: Classical edit operations

An edit algorithm is based on three operations:



MiGaL: Classical edit operations

An edit algorithm is based on three operations:



MiGaL: Classical edit operations

An edit algorithm is based on three operations:





Formally:



Formally:





Formally:



Formally:



In a "classical" editing algorithm we have:

- Relabeling
- Deletion
- Insertion
- In a "fusion" algorithm we have:
 - Relabeling
 - Node fusion
 - Edge fusion
 - Deletion/Insertion













This new edition algorithm is better to do a zoning than a exact matching. Currently we are searching for:

• The exact complexity of this algorithm (probably $O(n^k * n^2)$).

Score function for our edit operations.

MiGaL: Algorithm

We apply this algorithm to make a top-down matching across the differents MiGaL layers:



MiGaL: Algorithm

We apply this algorithm to make a top-down matching across the differents MiGaL layers:



The first results seems to show a great improvement in time computation:

- The direct edit algorithm on layer 3 on sample take around 4 minutes.
- The Top-Down algorithm make a few seconds to be compute.

Our approach seems to be original and pertinent in relation to the biological aspects.

The final algorithm should be polynomial in the length of the sequence (we hope so :)

Extension to multiple comparison/inference.