# EXTRACTING SENSE FROM STRUCTURE: AN APPLICATION TO FUNCTIONAL NON-CODING RNA POLYMERS

F. Costa

Bioinformatics Group
Department of Computer Science
Albert-Ludwigs-University Freiburg, Germany

13-17 February 2012

## What is this talk about?

Visualization of folding hypothesis landscape for a ncRNA family.
Semi-automatic construction of a vocabulary for structures.

- Allows finer grain view than clustering
- Useful to get an idea on the plasticity of a ncRNA family

## Map of the talk

1. **Introduction:** Clustering induces an implicit prototype (cluster center) with which to measure the typicality of its members. Consensus structures allow to visualize the average agreement on different parts.

2. **Question:** Can we decompose the character of a RNA family into meaningful traits?

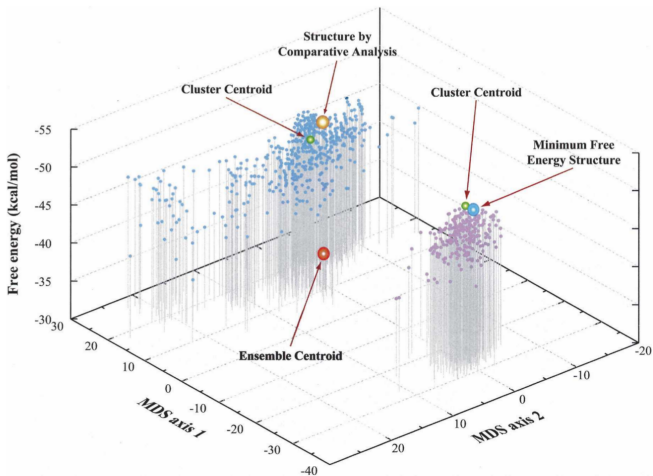3. **Answer:** Represent the principal directions of change identify the parts characteristic for different directions.

- The identification of a single folding structure to characterize a functional ncRNA family is a difficult and ill-posed problem
- **Idea:** characterize the entire set of probable structures

## ISSUES

How to best represent:

- multiple sub-groups
- continuous structural variation

MDS for all folding configurations for a single sequence[1].

[1]Source: Ding, Y. & Lawrence, C. E., *A statistical sampling algorithm for RNA secondary structure prediction*, Nucl. Acids Res. 2003

## ISSUES WITH RNA STRUCTURAL REPRESENTATIONS

- The MFE yields a single folding hypothesis that can be (at times) non representative
- Partition function based dot plots represents only statistics on all folding structures
- Accessibility information marginalizes base-pairedness in an aggregate with loss of structural information
- Suboptimal sampling is expensive and requires an additional (heuristic) clustering step

## PROPOSAL

- Use shape approach to derive a set of representative folding structures
- Represent each structure fully (i.e. as a labeled graph)
- Process set with graph kernels or explicit subgraph fingerprint techniques

FREIBURG

## SAMPLING REPRESENTATIVE STRUCTURES

- Sample all folding hypothesis
- ...which exhibit <u>significantly different</u> structure
- ...and are in a <u>small energy range</u> above the minimum free energy $\mapsto$ representative structures: *shapes*[a]

---

[a]Giegerich, B, Voß and M. Rehmsmeier, *Abstract shapes of RNA*, NAR 2004



(a) (b) (c)

(d)

| Shape | | |
|---|---|---|
| [] | GGGCCCAUAGCUCAGUGGUAGAGUGCCUCCUUUGCAAGGAGGAUGCCCUGGGUUCGAAUCCCAGUGGGUCCA | |
| [] | ((((((((((((((.(((.....((((((...)))))).)))))))))))........)))))))). | -35.9 kcal/mol |
| [[] []] | (((((((.....(.((((...(((((((...))))))).)))))))(((.......))).)))))))). | -32.2 kcal/mol |
| [[] [] []] | (((((...(((.......)))).(((((((...)))))))....(((((.......)))))))))))). | -31.7 kcal/mol |

## Representing RNA structure as graphs

Neighborhood Subgraph Pairwise Distance Kernel (NSPDK)[a]
Features as all pairs of near small neighborhood subgraphs
$\approx$ a generalization of *k-mers* with gaps

[a]F. Costa, K. De Grave, *Fast Neighborhood Subgraph Pairwise Distance Kernel*. ICML 2010

A C U U G G C U G U U C A A G U
( ( ( ( ( . . . . . . ) ) ) ) )

A

B

r=2

r=1

d=3

r=0

u

v

## From Linear Model to Importance Signal

Given a binary classification task, induce linear models:

- **performance:** good generalization guarantees
- **fast and scalable:** linear in practice;
  can manage $> 10^5$ instances
- **interpretable:** model $\mapsto$ set of feature-weight pairs

Interpret the weight as importance score for each feature

## From Importance Signal to Important Parts

1. Compute the importance for each vertex $v_i =$
   cumulative importance of all subgraphs that involve $v_i$
2. Visualize regions with high vertex importance...

decomposition
of graph
in feaures

FIGURE: Cumulative vertex importance

decomposition
of graph
in feaures

decomposition
of graph
in feaures

FIGURE: Cumulative vertex importance

decomposition
of graph
in feaures

FIGURE: Cumulative vertex importance

decomposition
of graph
in feaures

FIGURE: Cumulative vertex importance

weight
vector

decomposition
of graph
in feaures

single
feature

FIGURE: Cumulative vertex importance

weight
vector

decomposition
of graph
in feaures

single
feature

FIGURE: Cumulative vertex importance

FIGURE: Consensus RF00029 Intron gpll (Ribozyme)

1. Important parts (=connected components with importance > threshold) are <u>structural motifs</u> that can be clustered for characterization and insights

2. Importance score can complement energetic score in <u>folding</u> algorithms

3. Important parts can be constrained to match in <u>alignment</u> procedures even when dissimilar at sequence level

Using NSPDK we can represent graphs in a <span style="color:red">very high</span> dimensional vector space.

But how to map graphs onto a <span style="color:red">plane</span> for visual inspection?

## DIMENSIONALITY REDUCTION TECHNIQUES

1. Multi Dimensional Scaling (MDS)
   *Determin 2D coordinates so to maximally preserve the pairwise distances that instances originally had.*
   - - Non/trivial identification of directions of change
   - - Non-convex optimization problem $\rightarrow$ locally optima
   - + Non-linear embedding

2. Singular Value Decomposition (SVD)
   or equivalently Principal Component Analysis (PCA)
   *Rank ortogonal directions that induce the best reconstruction of the original vectors.*
   - + Trivial identification of directions of change
   - + Convex optimization problem $\rightarrow$ global optimal solution
   - - Linear embedding

FREIBURG

MDS for face image set[2].

Along the red line the expression moves from sad to happy.

SVD for digit image set[3].

$X \mapsto$ length of lower trait; $Y \mapsto$ thickness.

[3]Source: T.Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning. 09

Given a set of RNA sequences belonging to a functional class:

1. materialize *m* folding hypothesis for each sequence



FIGURE: Color-code: G-C=blu-cyan A-U=red-orange

## Proposal in a nutshell

Given a set of RNA sequences belonging to a functional class:

1. materialize *m* folding hypothesis for each sequence
2. structural graph $\mapsto$ vector representation

Given a set of RNA sequences belonging to a functional class:

1. materialize *m* folding hypothesis for each sequence
2. structural graph $\mapsto$ vector representation
3. SVD $\mapsto$ compute 2 main components and embed

Given a set of RNA sequences belonging to a functional class:

1. materialize *m* folding hypothesis for each sequence

2. structural graph $\mapsto$ vector representation

3. SVD $\mapsto$ compute 2 main components and embed

4. induce discriminative model on binary classification task:
   *instances in half space vs. instances in the other half*

Given a set of RNA sequences belonging to a functional class:

1. materialize $m$ folding hypothesis for each sequence
2. structural graph $\mapsto$ vector representation
3. SVD $\mapsto$ compute 2 main components and embed
4. induce discriminative model on binary classification task: *instances in half space vs. instances in the other half*
5. partition instances into $k \times k$ tiles in 2D plane

1. **Plot 1:** plot only one representative shape per tile (choose highest frequency shape)
2. **Plot 2:** plot importance signal on each vertex
3. **Plot 3:** plot consensus structures

UNI
FREIBURG

FIGURE: Artificial example: sequences with pattern $[U]^m GGGCCC[A]^n$

FIGURE: Artificial example: the part in common to all sequences that cannot be used to discriminate is white.

FIGURE: RF00005: tRNA

FIGURE: The poles represent GC vs AU content. Loop parts are white
$\mapsto$ more interesting.

0 _____ 1
Sequence conservation

FIGURE: RF00013: 6S

FIGURE: The presence of a (until recently unknown) functional hairpin is white $\mapsto$ important, and present in many consensus alignments.

0 ⟶ 1
Sequence conservation

FIGURE: RF00012: U3

FIGURE: The opposite poles represent the variant vertebrate vs. non-vertebrate, characterized by one stem vs. two stems.

- We propose to:
    - visualize a set of folding structures
    - embed them in a plane...
    - ...whose coordinate system is aligned to the directions of major sequence-structural changes
- automatically learn how to discriminate between extreme cases
- ... to identify common regions

## HOW CAN ALL THIS BE USED?

- Give biologists a new way to look into a ncRNA family
- Help them identify in a semi-automatic way interesting parts or sub-families
- Help them characterize and give a name to the structural *(and possibly functional)* traits
- Give a way to use biological knowledge to select a subset of meaningful sequences-structures to make better models

FREIBURG

## Thanks
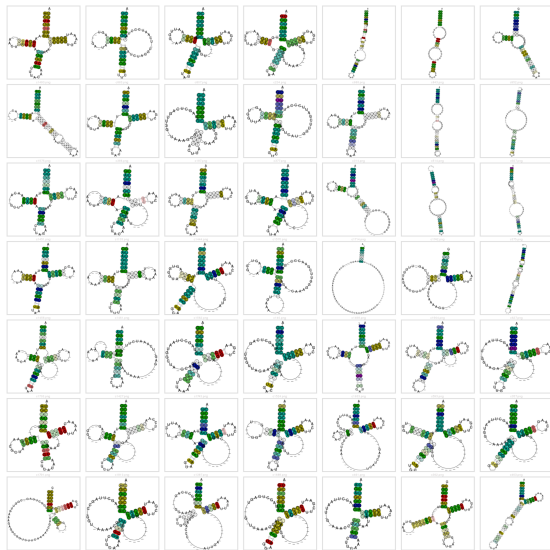
*Acknowledgments:*
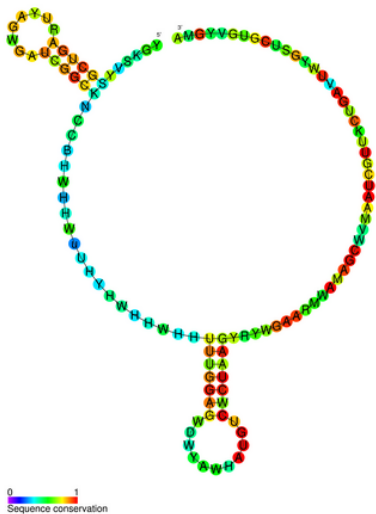Steffen Heyne
Sita Lange
Robert Kleinkauf
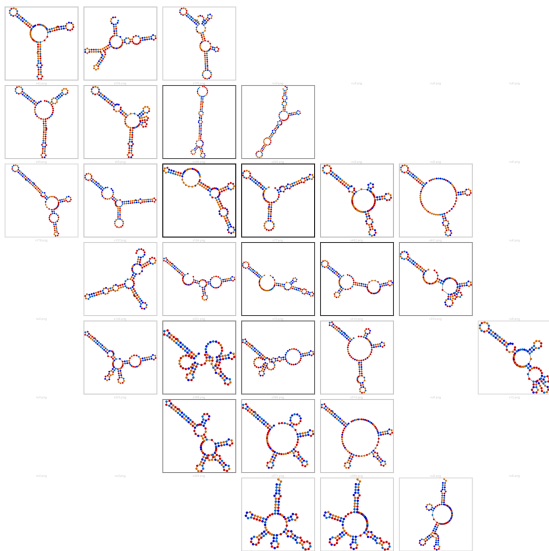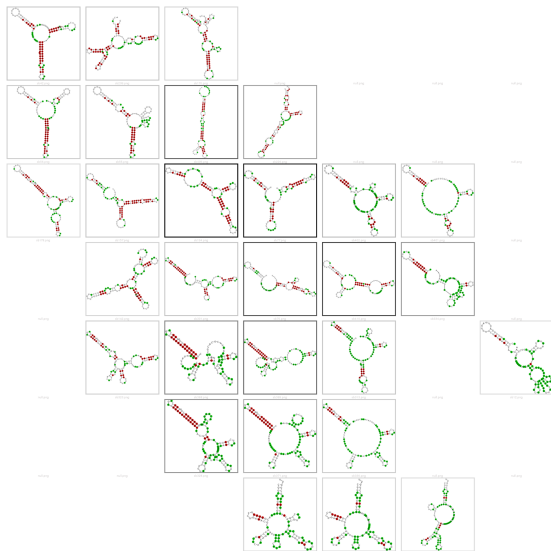Manja Marz
Rolf Backofen
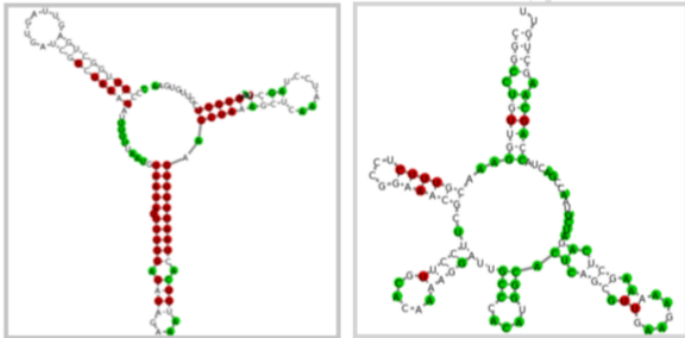
FIGURE: RF00114: Ribosomal S15 Leader

FIGURE: While one hairpin is common to the whole family (top-left white), the second hairpin (bottom green) seems to represent only one of the extreme cases.
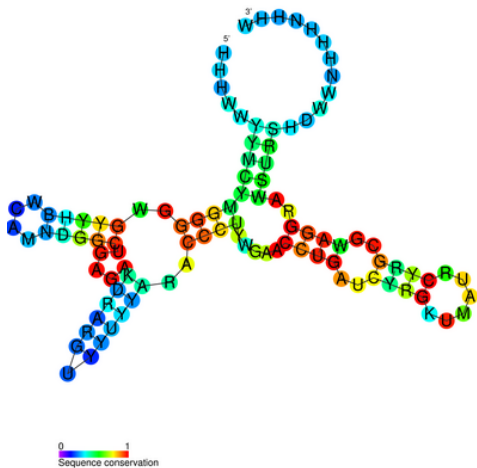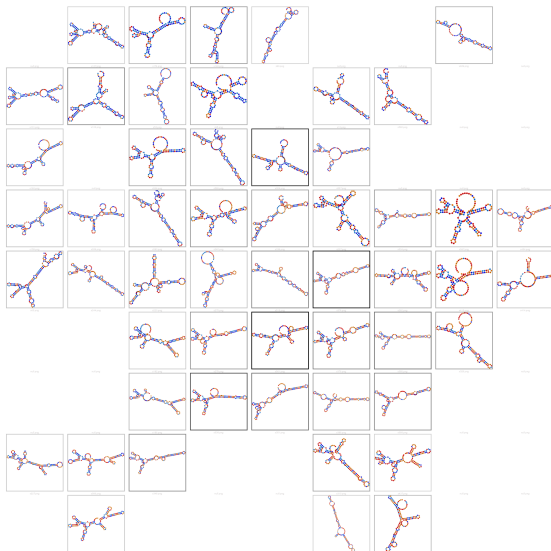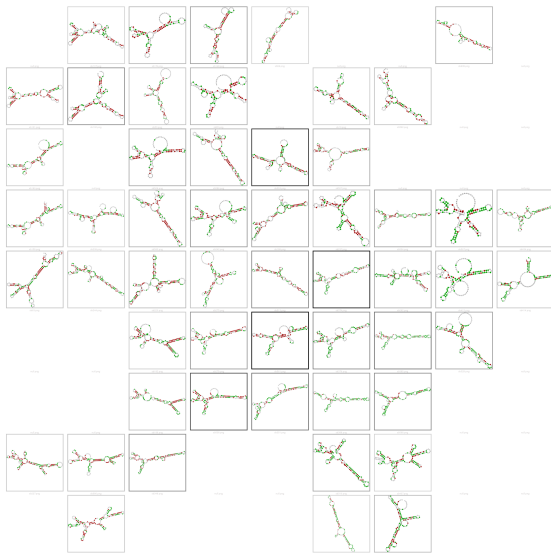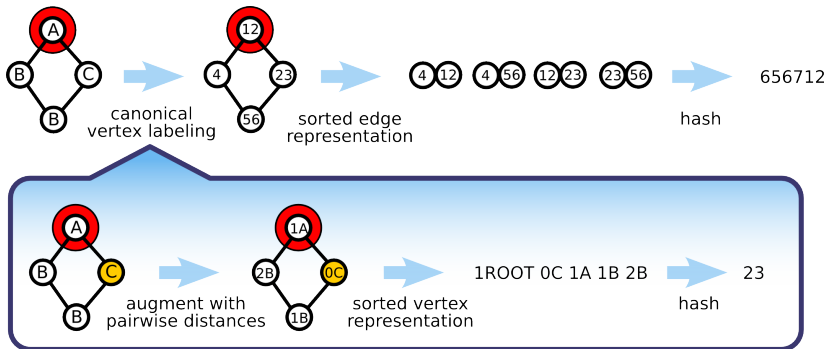
FIGURE: RF00059: TPP Ryboswitch

Given graph as a (multi)set of pairs of near small subgraphs
compute the explicit sparse representation via hashing techniques



Complexity dominated by edge sorting or all-pairwise-distance
computation in small subgraphs $\mapsto$ efficient (linear) in practice