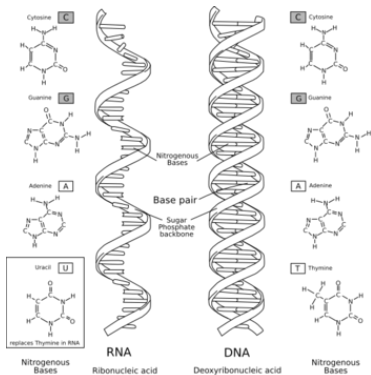


LocARNA-P: Accurate Boundary Prediction and Improved Detection of Structured RNAs

Sebastian Will

CSAIL, MIT

RNA



chain of building blocks A,C,G,U

once upon a time . . .

RNA was considered only an intermediate in protein synthesis

Changed Picture: RNA plays central role



Multitude of Non-coding RNAs
with various functions

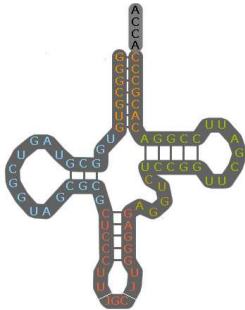
- gene regulation
- catalysis of reactions
(“ribozymes”)

Versatility due to RNA Structure

primary structure = sequence (of bases A,C,G,U)

GGGCGUGUGGCGUAGUCGGUA ... GUUCGAUUCGGACACGCCACCA

secondary structure



base pairs: C-G,A-U,G-U

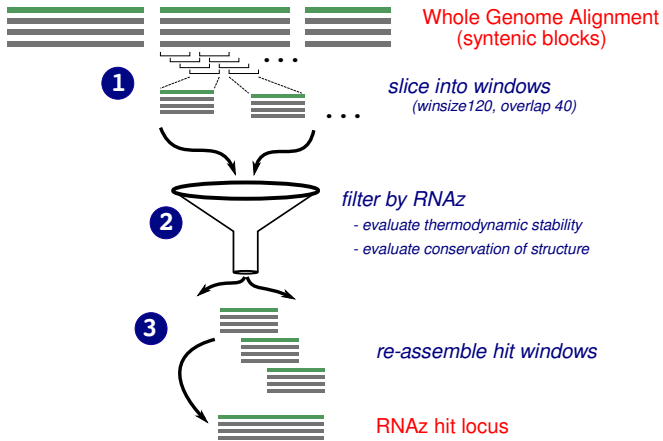
tertiary structure



all atoms 3D

De-novo Prediction of Structural RNAs

RNAz Washietl *et al*, Fast and reliable prediction of noncoding RNAs, *PNAS*, 2005



RNAz detects $\geq 40,000$ ncRNA candidate loci in Fly

Current Limits of De-novo ncRNA Prediction

Main Problems:

- coarse RNA window boundaries
- high false discovery rate (FDR) $\approx 50\%$
experimental analysis expensive

Main Limitations:

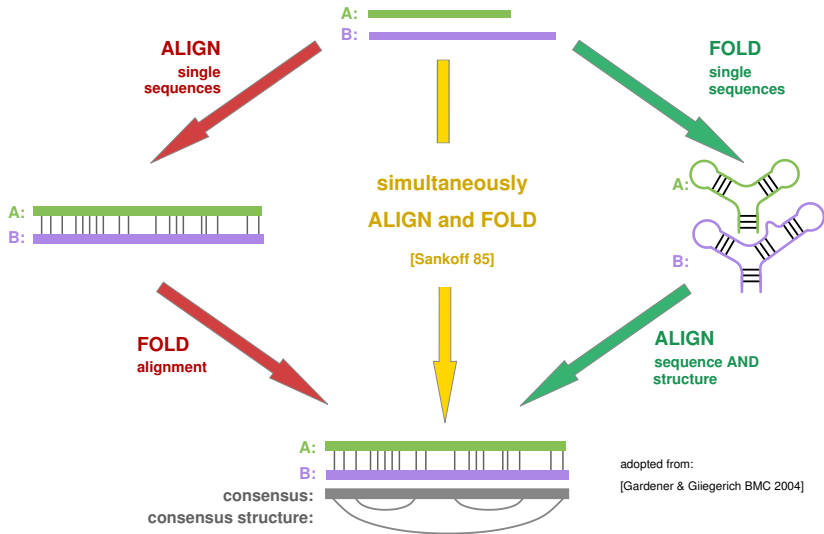
- no structural re-alignment (WGA sequence-based)
- no information about alignment reliability
- no information about alignment space



Recall talk title

“LocARNA-P: Accurate Boundary Prediction and Improved Detection of Structured RNAs for Genome-wide Screens”

Comparative RNA Analysis



LocARNA: more efficient Alignment and Folding

- **LocARNA** does a simplified variant of Simultaneous and Folding due to Hofacker

Sankoff

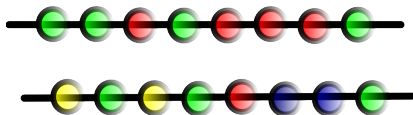
Zuker \times Needleman-Wunsch

full energy model

Hofacker

Nussinov \times Smith-Waterman +
McCaskill

use full energy model
via base pair probabilities



LocARNA: more efficient Alignment and Folding

- **LocARNA** does a simplified variant of Simultaneous and Folding due to Hofacker

Sankoff

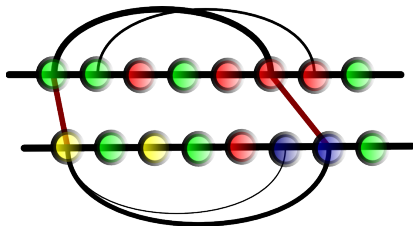
Zuker \times Needleman-Wunsch

full energy model

Hofacker

Nussinov \times Smith-Waterman +
McCaskill

use full energy model
via base pair probabilities



LocARNA: more efficient Alignment and Folding

- **LocARNA** does a simplified variant of Simultaneous and Folding due to Hofacker

Sankoff

Zuker × Needleman-Wunsch

full energy model

Hofacker

Nussinov × Smith-Waterman +
McCaskill

use full energy model
via **base pair probabilities**

- **LocARNA** is fast and robust
- **LocARNA** improves scoring and adds features (e.g. locality)
- **LocARNA** reduces space complexity to $O(n^2)$



LocARNA-P adds Probabilities

- match probabilities for sequence **and** structure

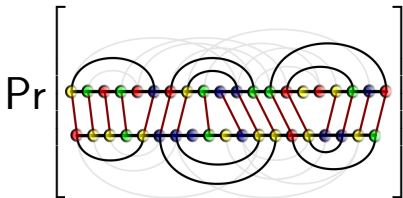


- reduced space complexity: $O(n^2)$ “as in LocARNA”
- make use of very accurate LocARNA scoring

Alignment Probabilities

LocARNA scoring function: $\text{score}(\mathcal{A}, \mathcal{S})$

Given sequences A, B : define **probability**



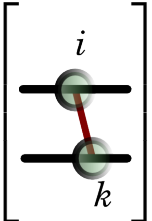
i.e.: $\text{Pr}[\text{ pair of alignment } \mathcal{A} \text{ and consensus structure } \mathcal{S}]$

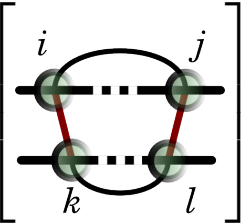
as probability in **Boltzmann ensemble**

$$\text{Pr}[(\mathcal{A}, \mathcal{S})] \sim \exp(\beta \text{score}(\mathcal{A}, \mathcal{S}))$$

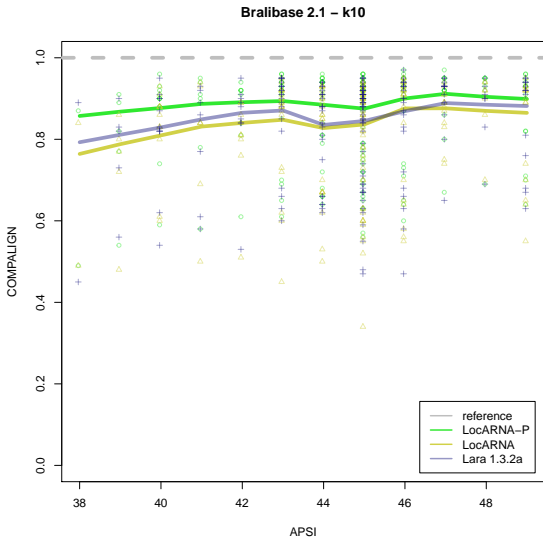
needs **partition function** $Z = \sum_{(\mathcal{A}, \mathcal{S})} \exp(\beta \text{score}(\mathcal{A}, \mathcal{S}))$

Match Probabilities

$$Pr \left[\begin{array}{c} i \\ \text{---} \\ \text{---} \\ k \end{array} \right] = \sum_{(\mathcal{A}, \mathcal{S}) \text{ with } i \sim k \in \mathcal{A}_S} Pr[(\mathcal{A}, \mathcal{S})]$$
A diagram showing two horizontal black lines representing paths. The top line has a green circular node labeled 'i' on it. The bottom line has a green circular node labeled 'k' on it. A red line connects the two nodes, representing a match between them.

$$Pr \left[\begin{array}{c} i \quad j \\ \text{---} \quad \text{---} \\ \text{---} \quad \text{---} \\ k \quad l \end{array} \right] = \sum_{(\mathcal{A}, \mathcal{S}) \text{ with } (i,j) \sim (k,l) \in \mathcal{S}} Pr[(\mathcal{A}, \mathcal{S})]$$
A diagram showing two horizontal black lines representing paths. The top line has two green circular nodes labeled 'i' and 'j' on it, with a dashed line between them. The bottom line has two green circular nodes labeled 'k' and 'l' on it, with a dashed line between them. Two red lines connect the nodes: one from 'i' to 'k' and one from 'j' to 'l'. A black oval encloses the entire diagram, representing a match between the two paths.

LocARNA-P Improves Multiple Alignment



10-fold alignment, low sequence identity $\leq 50\%$

Local Alignment Quality: Reliability Profiles

Goal measure reliability of alignment columns

Method sum pairwise match probabilities

Result Reliability profile (structure # + sequence *)

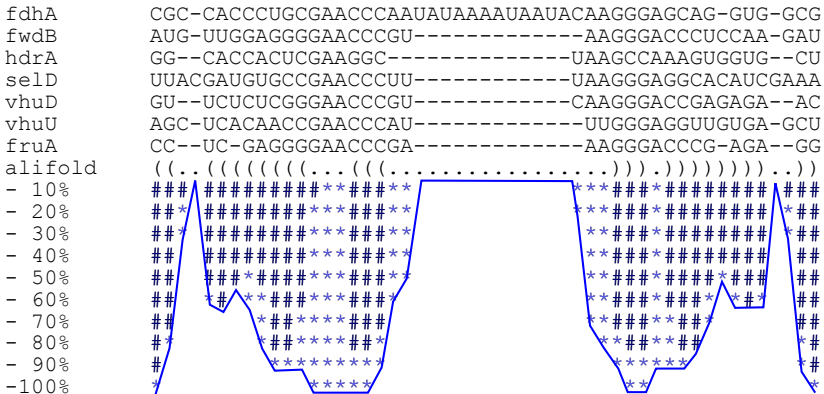
```
fdhA      CGC-CACCCUGCGAACCCAAUAUAAAAUAAUACAAGGGAGCAG-GUG-GCG
fwdB      AUG-UUGGAGGGGAACCCGU-----AAGGGACCCUCCAA-GAU
hdrA      GG--CACCACUCGAAGGC-----UAAGCCAAAGUGGUG--CU
selD      UUACGAUGUGCCGAACCCUU-----UAAGGGAGGCACAUCGAAA
vhuD      GU--UCUCUCGGGAACCCGU-----CAAGGGACCGAGAGA--AC
vhuU      AGC-UCACAACCGAACCCAU-----UUGGGAGGUUGUGA-GCU
fruA      CC--UC-GAGGGGAACCCGA-----AAGGGACCCG-AGA--GG
alifold   ((..((( ((( ( ((..((( ( (.....)))))))))..))
- 10%     ### #####**##**          ***###*#####  ##
- 20%     ##* #####**##**          ***###*##### *##
- 30%     ##* #####**##**          *###*##### *##
- 40%     ## #####**##**          *###*#####  ##
- 50%     ## ###*#####*#####*    *###*#####*##
- 60%     ## *# **#####*#####*   **#####* *##
- 70%     ##      *#####         **#####*   ##
- 80%     #*        *#####         **#####*   *##
- 90%     #          ****          ***#####         *##
-100%    *            *****          **#####         *
```

Local Alignment Quality: Reliability Profiles

Goal measure reliability of alignment columns

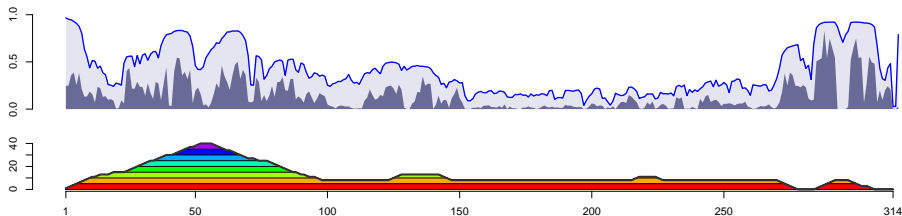
Method sum pairwise match probabilities

Result Reliability profile (structure # + sequence *)



Reliability Profile: Case Study 7SK

- Here: realign hand curated alignment
- Reliability profile fits with experience from manual alignment



dark blue structure reliability

light blue sequence reliability

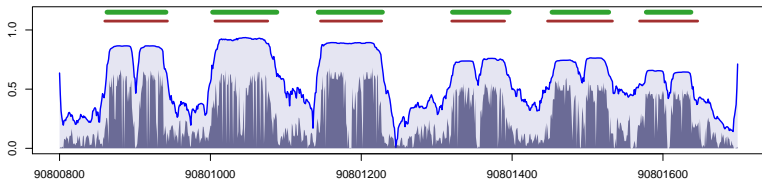
blue line total reliability = structure + sequence

rainbow colored plot mountain plot, shows RNA structure

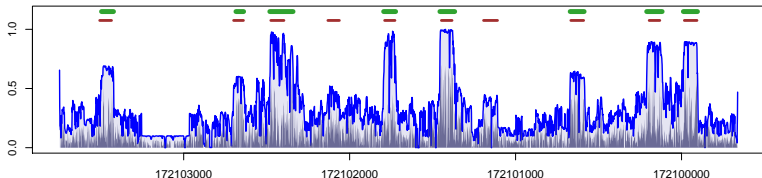
Reliability Profiles

- genomic cluster with known ncRNAs
- align corresponding regions in 10/5 vertebrates
- show reliability profile for human DNA

cluster of 6 micro RNAs, length ≈ 900

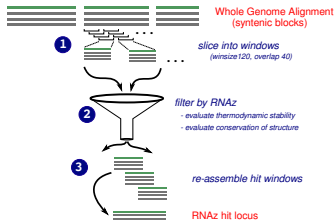


cluster of 10 CD-Box snoRNAs 'GAS5', length ≈ 4000



green = LocARNA-P prediction; red = ncRNA annotation

Reliabilities for Refining the Drosophilids ncRNA Screen



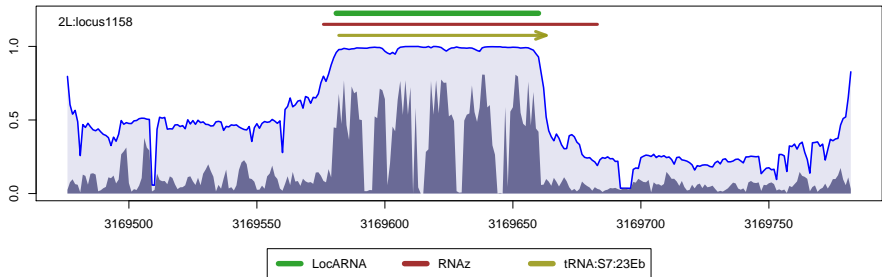
Rose *et al.* Computational RNomics of Drosophilids, BMC Genomics, 2007.

12 Drosophilid genomes alignment, RNAz: 120nt windows in 40nt, combine windows into loci

NEW: Analyze hits using reliabilities

- Realign with context (100 up, 100 down) and reliability profile
- Predict boundaries of ncRNA ⇒ **exact location**
- Compute reliability score ⇒ **improve predictive power**
- Benchmark: 301 RNAz loci annotated as ncRNA

LocARNA-P Refines Drosophilids ncRNA Screen

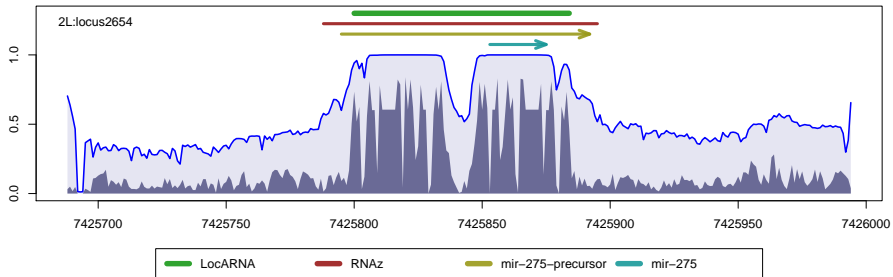


dark blue structure reliability

light blue sequence reliability

blue line total reliability = structure + sequence

LocARNA-P Refines Drosophilids ncRNA Screen

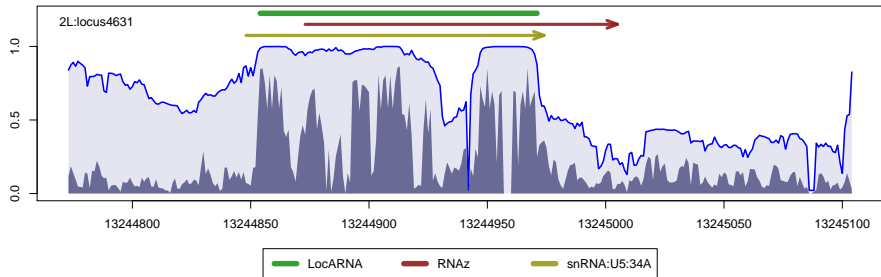


dark blue structure reliability

light blue sequence reliability

blue line total reliability = structure + sequence

LocARNA-P Refines Drosophilids ncRNA Screen

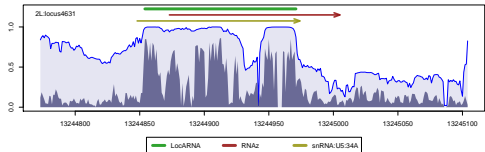
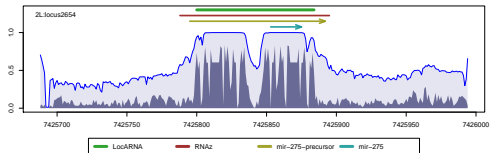
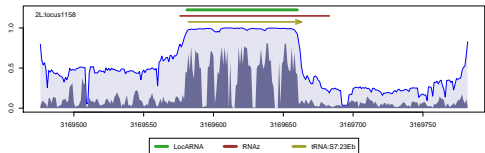
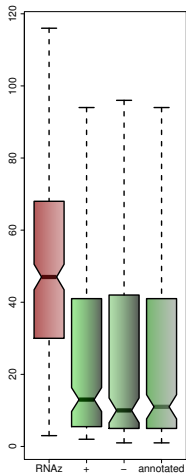


dark blue structure reliability

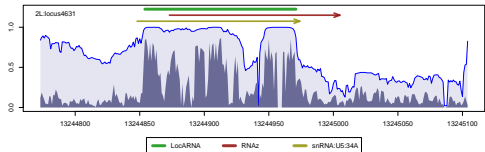
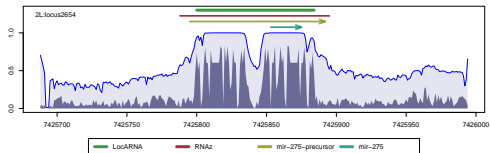
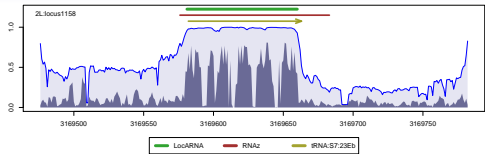
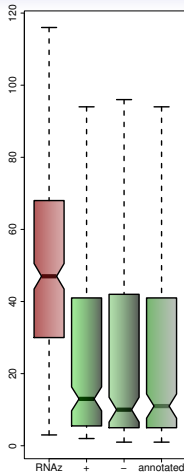
light blue sequence reliability

blue line total reliability = structure + sequence

Significant Improvement of Boundaries

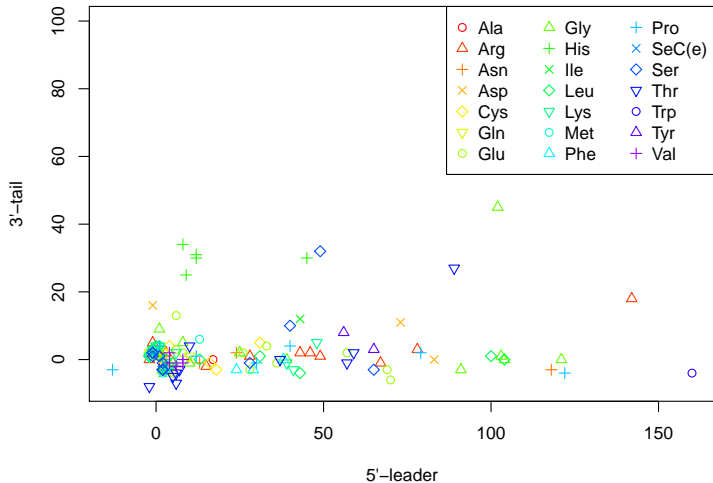


Significant Improvement of Boundaries



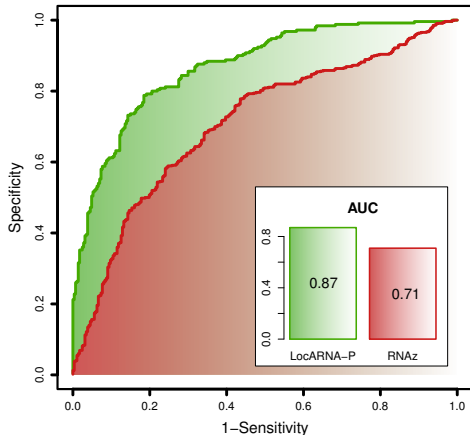
Moreover: many deviations have a biological reason
for example: tRNA

Deviation Reveals 5' Signal for tRNAs



Non-random difference of deviation distributions in header and tail indicates signal!

LocARNA-P Improves Predictive Power of Predictions

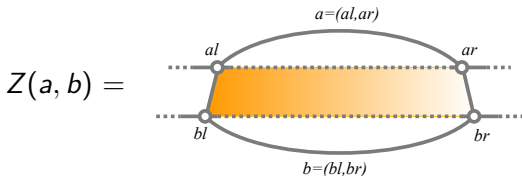


- positive set = Rfam annotated RNAz hits
- negative set by shuffling of RNAz hits
- LocARNA-P: reliability score vs. RNAz max. P score

Partition Functions by Inside Algo

Partition Function $Z = \sum \exp(\beta \text{score}(\mathcal{A}, \mathcal{S}))$
: \mathcal{A} alignment, \mathcal{S} structure

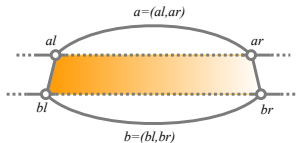
Inside PFs $Z(a, b) =$ Partition function, where
 $a = (al, ar) \sim b = (bl, br)$
inside $[al..ar]$ and $[bl..br]$



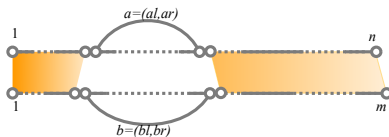
Computation of Match Probabilities

Requires inside and outside partition functions

Inside $Z(a, b)$:



Outside $Z'(a, b)$:

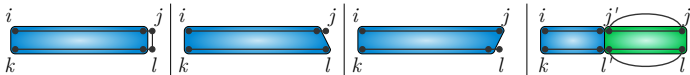


Structural Match probability $Pr[a \sim b] = Z(a, b) \cdot Z'(a, b) / Z$

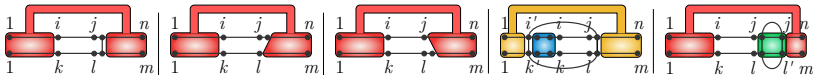
Sequence Match probability slightly more complex
(sum over lots of cases)

Inside and Outside by DP

Inside Decomposition



Outside Decomposition



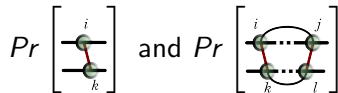
Complexity

naïve: $O(n^6)$ time, $O(n^4)$ space

LocARNA-P: $O(n^4)$ time, $O(n^2)$ space

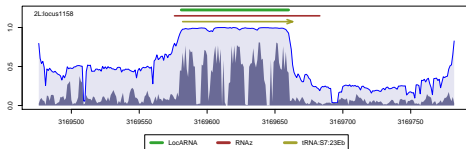
Take home: LocARNA-P

- Reliability Profiles due to sequence-structure alignment

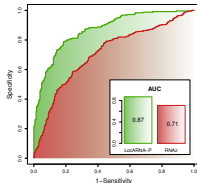


- Refines de-novo prediction of ncRNA

- Boundaries



- Predictive Power



- Potential for broad, general application to RNA analysis

Thank you for your attention

and thanks to these people for working with me on LocARNA(-P)

Freiburg

Rolf Backofen

Steffen Heyne

Tejal Joshi

Leipzig

Peter Stadler

Kristin Reiche

Wolfgang Otto

Michael Siebauer

Wien

Ivo Hofacker

More Info+Download:

<http://www.bioinf.uni-freiburg.de/Software/LocARNA/>

Web Server:

<http://rna.informatik.uni-freiburg.de>