

RNAcontext: A new method for Learning the Sequence and Structure Binding Preferences of RNA-Binding Proteins

Presented by: Ria X. Peschutter

OPEN O ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

RNAcontext: A New Method for Learning the Sequence and Structure Binding Preferences of RNA-Binding Proteins

Hilal Kazan¹, Debashish Ray², Esther T. Chan³, Timothy R. Hughes^{2,3,4}, Quaid Morris^{1,2,3,4}*

1 Department of Computer Science, University of Toronto, Toronto, Ontario, Canada, 2 Banting and Best Department of Medical Research, University of Toronto, Toronto, Toronto, Ontario, Canada, 3 Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada, 4 Donnelley Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, Ontario, Canada, 4 Donnelley Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, Ontario, Canada, 4 Donnelley Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, Ontario, Canada

Received September 23, 2009; Accepted May 25, 2010; Published July 1, 2010

Overview

- 1) Biological Background
- 2) RNACompete
- 3) RNAContext
 - a) RNA annotation
 - b) Motif model
 - c) Fitting of the model
- 4) Results

RNA binding proteins - RBPs

- Critical roles in numerous cellular processes
- Essential during germline and early embryonic development

Challenges:

- Binding preferences not well characterized
- Multiple RNA targets
- Identification of mechanism behind function



http://en.wikipedia.org/wiki/File:Zinc_finger_rendered.png

nature biotechnology

Rapid and systematic analysis of the RNA recognition specificities of RNA-binding proteins

Debashish Ray^{1,4}, Hilal Kazan^{2,4}, Esther T Chan³, Lourdes Peña Castillo¹, Sidharth Chaudhry³, Shaheynoor Talukder¹, Benjamin J Blencowe^{1,3}, Quaid Morris^{1–3} & Timothy R Hughes^{1,3}



© 2009 Nature America, Inc. All rights reserved.

RNACompete

- Estimates binding affinity data
- Pool of unique short RNA sequences (29- to 38-nt)



RNACompete

- <u>Eukaryotic RBPs recognize:</u>
 - Short unstructured sequences
 - → Loop sequences in RNA stem-loop



<u>Mircoarray design:</u>

- independent duplicate sets of unstructured and stem-loop RNAs
- Set A stem-loop and unstructured sequences
- Set B stem-loop and unstructured sequences



RNACompete

<u>Caution:</u>

- Unstructured sometimes of stem-loops
- Stem-loops unexpected structures



<u>Constraints to minimize:</u>

- → Folding of unstructured RNAs
- misfolding of structured RNAs
- → Extensive base-pairing among any two RNAs
- microarray cross-hybridization

RNAcontext

- Motif model
- Predicts sequence and structure preferences of RBPs
- Considers multiple structural preferences simultaneously
- Webserver (restriction on input size)
- Terminal application (no restriction on input size)



Methods

Three steps:

- 1. How to annotate RNA in terms of structural context
- 2. Details and mathematical formulation of the motif model
- 3. How to fit the RNAcontext motif model



Methods - RNA annotation

- Annotation of each base in each structure by context alphabet
- Predicting secondary structure \rightarrow SFOLD
- RNA multiple distinct stable secondary structures
- Distribution over all its possible context

Alphabet A:

- P paired
- L hairpin loop
- **U** unstructured or external region
- M miscellaneous

<u>Input :</u>

- Set of sequences together with their estimated binding affinities (S)
- RNA secondary structure annotations of the sequences (P)







estimate of probability that RBP binds $s_{t+1:t:k}$ in ideal structural context

estimate of probability that RBP prefers structural context $p_{t+1:t*K}$

- S sequence set
- P annotation profiles
- K width of binding site
- <u>Θ model parameter:</u>
- β s sequence affinity bias
- βp structural context bias
- Γ structure annotation

 $p_{\alpha,k}$ – probability that base at posk of shas structural context α

$$\underbrace{N^{seq}(s, \Theta) = \sigma(\beta_s + \sum_{k=1}^{K} \Theta_{s_k,k})}_{estimate of probability that RBP binds s_{i+1:::k} in ideal structural context} \qquad \underbrace{C(p, \Theta) = \sigma(\beta_p + \sum_{\alpha \in A} \Gamma_\alpha * \sum_{k=1}^{K} p_{\alpha,k})}_{estimate of probability that RBP prefers structural context p_{i+1:::K}}$$
Logistic function: $\sigma(x) = (1 + \exp(-x))^{-1}$

$$\int_{-6}^{0} -4 -2 \int_{-2}^{0} \int_{0}^{0} \frac{1}{2} \int_{-4}^{0} \frac{1}{6}$$
http://en.wikipedia.org/wiki/Logistic_function

$$\underbrace{\widetilde{f(s,p,\Theta)}}_{t=0} = 1 - \prod_{t=0}^{|s|-K} 1 - N(s,p,\Theta) = 1 - \prod_{t=0}^{|s|-K} 1 - (N^{seq}(s,\Theta) * C(p,\Theta))$$

- $\Gamma \alpha \rightarrow$ favored structure
 - → N(s,p,Θ) ≈ N(s,Θ)
 - → determined by s
- $\Gamma \alpha \rightarrow$ unfavored structure
 - → C(p, Θ) ≈ 0
 - → N(s,p, Θ) ≈ 0
 - → determined by p



http://en.wikipedia.org/wiki/Logistic_function

Methods - Motif model fitting

- Model affinity of sequence as linear model
- Minimizing least squares cost function by L-BFGS-B
 - Constrain slope of linear model to only take positive values
 - Making use of position weight matrix (PWM)
 - Small constant to ensure unique global minimum

- Cost function multi-modal
 - -> different initializations \rightarrow different outputs
 - → Take the best

RNAcontex-Output



Results

- Compared MEMERIS(MEME), RNAcontext and MatrixReduce
- Each method trained individually
- making use of each possible parameter setting
- best model

Proteins	RNAcontext	MEMERIS	MatrixREDUCE	
RBM4	0.91	0.43	0.63	
FUSIP1	0.53	0.31	0.32	
Vts1p	0.65	0.58	0.56	
YB1	0.17	0.07	0.11	
SLM2	0.81	0.49	0.77	
SF2	0.70	0.50	0.66	
U1A	0.30	0.27	0.21	
HuR	0.96	0.74	0.94	
PTB	0.69	0.26	0.67	

Results



Results - RNAcompete

Protein	<u>Domain(s)</u>	Known motif/ _binding site	Our highest- correlating motif	
Vts1	One SAM domain	الله الله الله الله الله الله الله ال	$\begin{bmatrix} 2.0 \\ 1.0 \\ 0.0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \end{bmatrix}$	
SLM2	One KH domain	NA	ZIO AUAAA	5
YB1	Full-length; one cold-shock domain		^{2.0} ^{21.0} 0.0 1 2 3 4 5 6 7 8	¢
RBM4	Full-length; two RRM domains	NA	2.0 普1.0 0.0 1 2 3 4	
SF2/ASF	Two RRM domains	E GAAGAAC	2.0 2 1.0 0.0 1 2 3 4	1
FUSIP1	One RRM domain	ACAAAGACAAA	2.0 2.0 2.0 AGAG	
HuR	Full-length; three RRM domains		$\begin{array}{c} 2.0\\ \begin{array}{c} \\ \end{array}\\ 1.0\\ 0.0 \end{array} \\ 1 2 3 4 5 6 7 8 9 10 \end{array}$	
U1A	One of two RRM domains	器1.0 0.0 1 2 3 4 5 6 7		
РТВ	Full-length; four RRM domains	^{2.0} 登1.0 0.0 1 2 3 4	20 # <u>f</u> 1.0 0.0 1 2 3 4 5 6 7 8	

Results - RNAContext

_	Domain(s)	Previously reported binding site	RNAcontext predicted motifs
Vts1p	one SAM domain	in stem-lo source: [17]	context: P L L L L L
SLM2	one KH domain	Source: [5]	context: U U U U U U U
RBM4	full-length two RRM domains	source: [5]	context: U U U U U U U
SF2/ASF	two RRM domains	Source: [43]	Context: M M M M P
FUSIP1	one RRM domain	source:[5]	context: U U U U U U U U U
HuR	full-length; three RRM domains	AUUUAUUA GAUUAUUAG source: [12]	context: U U U U U U U P
РТВ	full-length; four RRM domains	source: [44]	context: U U U U

Summary

- New motif model of RBP binding preferences
 - → Sequence and structure
 - Corresponding algorithm for fitting model
- Recovers previously reported sequences and binding preferences
- Predicts new structure binding preferences
- Interprets RNAcompete data better than correlation analysis
- Initially designed for short RNA sequences can be used for long:
 - RNAplfold

Thank you for your attention!

Questions?

Methods – L-BFGS-B

BFGS :

- iterative method
- solves unconstrained non-linear optimization problems
- Approximates Newton's method by Hessian matrix
- Hessian
 - → square matrix of second-order partial derivatives of a function
- stores and computes whole matrix

Methods – L-BFGS-B

L-BFGS:

- Vectors implicitly model Hessian
- Optimization algorithm for parameter estimation in machine learning

L-BFGS-B:

- extends L-BFGS
- handle simple box constraints (aka bound constraints) on variables
- identifies fixed and free variables at every step
- L-BFGS method only free variables
 - → higher accuracy