## Orthogonal nonnegative factorization-based analysis of nineteen protein-RNA interaction CLIP data sets

Martin Stražar<sup>1</sup>, Marinka Žitnik<sup>1</sup>, Blaž Zupan<sup>1</sup>, Jernej Ule<sup>1,2</sup> and Tomaž Curk<sup>1</sup>

<sup>1</sup>Bioinformatics Laboratory, University of Ljubljana, Faculty of computer and information science, Večna pot 113, 1000 Ljubljana <sup>2</sup>Department of Molecular Neu oscience, UCL Institute of Neurology, Queen Square, London, WC1N 3BG, UK

## Abstract

RNA binding proteins (RBPs) play a major role in post-transcriptional processes: splicing, transport and polyadenylation<sup>1</sup>. Computational modelling of protein-RNA interactions depends on precise structural knowledge, which is often difficult to obtain<sup>2</sup>. We propose an approach to model protein-RNA interactions using publicly available data. We combine data sources to achieve accurate prediction and a comprehensive model based on integrative, orthogonal nonnegative matrix factorization (iONMF).



## **Combinations of data sources** improve prediction



Data sources are matrices describing features in the interval [-50, 50] nt positions around crosslink sites.

Ranking of individual data sources:

- **R:** RNA structure, probability of double-stranded RNA, which affects nucleotide accessibility.
- C: co-binding; combinatorial protein-RNA interactions, which compete or cooperate for the same RNA.
- **T:** region type, **K:** RNA k-mers, **G:** Gene Ontology.

$$J = \sum_{i=1}^{N} (\|\mathbf{X}_i - \mathbf{W}\mathbf{H}_i^T\|_2^2 + \alpha \|\mathbf{H}_i^T\mathbf{H}_i - \mathbf{I}\|_2^2)$$
 loss function

iONMF: integrative orthogonal nonnegative matrix factorization

📕 ELAVL1 📕 hnRNPC 📕 TIAL1

0 -50-40-30-20-10 0 10 20 30 40 50

**U2AF2** 

Co-binding with U2AF2.

hnRNPL

ELAVL1-MNase

Aqo2

The original data matrices  $X_i$  are approximated **simultaneously** as a product of: • Common matrix **W**: similarity of positions based on all data sources. • Matrices **H**<sub>i</sub>: commonly occuring data source-specific patterns.

54.3%

hnRNP

hnRNPC

## **Orthogonality constraints induce** non-overlapping patterns



hnRNPC

86.5%

QKI

FUS

/ TIAL1

5'UTR 3'UTR CDS exor 46.2% 49.2%





Gene

G

TDP-43

Ontology

Ago2-MNase



Splicing regulators bind U-rich (intronic) motifs.

hnRNPC motifs hnRNPC binding site CAPS regulated exon

CCA**GGCTGG**TATGCGGTGGTGTGATCGTAGCTCACTGCAGTCTCGAACTCCTGGGTTCA**A** GCGATCCTTCCACTTCAGCCTCCCAAGTAGCTGGTACTACAGgtgtgtgccacgacacccc gctaagtttttgaaatttatttttgtagagacaggattttcctatgttgcccaggctggt tttcaaact**cctggc**c

Common motifs at antisense Alu elements (GGCUGG, GCCCAG, CCUGCC, GCCGGG) are associated with hnRNPC.

-50-40-30-20-10 0 10 20 30 40 50 -50-40-30-20-10 0 10 20 30 40 50

> Argonaute binding sites are known to contain and unstructured region **upstream** of the crosslink sites<sup>6</sup>.



Discovered motifs are in accordance with *in vitro* assay RNAcompete<sup>8</sup>.

Data was obtained from servers iCount (http://icount.biolab.si) and DoRiNA <sup>3</sup> G. Hutvagner and M. J. Simard, Nat. Rev. Mol. Cell Biol., Jan. 2008. References <sup>4</sup>I.L. Hofacker et.al Monatshefte f. Chemie (1994). (dorina.mdc-berlin.de). <sup>1</sup> König et. al, Nature review Genetics (2013) <sup>5</sup> H. Kazan and Q. Morris, PLoS Comput. Biol, Jan. 2010.

<sup>2</sup> Puton et al., Journal of structural biology (2013)

<sup>6</sup> Li et al., Genome Res. (2014.)

<sup>7</sup> D. Ray, T. R. Hughes, et al. Nature, Jul. 2013.

This work was supported by grants from the Slovenian Research Agency (P2-0209, J7-5460) and the European Research Council (206726-CLIP).