

InSyBio ncRNASeq: a Web tool for analysing non-coding RNAs

Aigli Korfiati^{1✉}, Konstantinos Theofilatos¹, Christos Alexakos¹, Seferina Mavroudi^{1,2}

¹InSyBio Ltd, United Kingdom

²UK & Department of Social Work, School of Sciences of Health and Care, Technological Educational Institute of Western Greece, Koukouli, Patra, Greece

Competing interests: AK, KT, CA: employer of InSyBio Ltd. InSyBio Ltd participates in the NBG Business Seeds program by the National Bank of Greece and financed this study. However, InSyBio and NBG Business Seeds program were not involved in the study design, collection, analysis & interpretation of data and writing of the paper. They only participated in the process of selecting the suitable journal for submitting the paper. SM: external collaborator of InSyBio Ltd and an employer in the Department of Social Work, School of Sciences of Health and Care, Technological Educational Institute of Western Greece. InSyBio Ltd participates in the NBG Business Seeds program by the National Bank of Greece and financed this study. However, InSyBio and NBG Business Seeds program were not involved in the study design, collection, analysis & interpretation of data and writing of the paper. They only participated in the process of selecting the suitable journal for submitting the paper.

Abstract

Non-coding RNA (ncRNA) genes encode non-protein-coding RNAs, which are classified into infrastructural and regulatory ncRNAs, depending on their role. Regulatory ncRNAs are involved in a variety of key cellular processes, and are hence associated with several diseases. For this reason, their accurate and efficient identification, and the identification of their targets, is a promising research area and an open topic for the bioinformatics community. We present InSyBio ncRNASeq, a cloud-based Web platform that assists users to analyse ncRNA sequences, and predict whether they are miRNAs, pseudo-hairpins or whether they belong to another ncRNA category. Additionally, InSyBio ncRNASeq offers a unique miRNA target-prediction pipeline, which results in scored miRNA target sites in 3'-untranslated regions (3'-UTR) of messenger RNAs (mRNAs). We show this tool to have the highest accuracy metrics compared with other state-of-the-art methods and tools.

Introduction

Traditionally, most RNA molecules were regarded as information-carrying intermediates from the gene to the translation machinery, with exceptions being transfer RNA (tRNA) and ribosomal RNA (rRNA), both of which are infrastructural RNAs that are central to the protein-synthesis process. However, since the late 1990s, evidence (Liu *et al.*, 2005) suggests that other types of non-protein-coding RNA molecules are present in organisms ranging from bacteria to mammals, and are involved in a variety of cellular processes, including plasmid replication, phage development, bacterial virulence, chromosome structure, DNA transcription, RNA processing and modification, development control, regulation of gene expression and others. In humans, it is estimated that about 98% of the genome can be transcribed, out of which only ~2% code for proteins, suggesting the possibility that a large percentage of the genome may encode ncRNAs. ncRNAs have been reported as biomarkers for disease and response to

treatment in cancer, liver and cardiovascular diseases, and central nervous system disorders, among many others (Lopez *et al.*, 2015).

Although the importance of ncRNAs in cellular activities is well known, and new members and classes of ncRNA are continuously being reported, our current knowledge about the collection of all ncRNAs present in a particular genome is very limited because of the lack of effective computational or experimental methods. Indeed, their computational identification represents one of the most important and challenging issues in computational biology. Existing methods for ncRNA-gene prediction rely mostly on homology information, thus limiting their applications to ncRNA genes that have homologues.

Evidence linking ncRNAs with diseases seems to be attracting the attention of the pharmaceutical and biotechnology industries. Companies and institutions are developing ncRNA-based strategies against cancer, cardiovascular, neurological and muscular diseases. Thus, continued investigation of ncRNA biogenesis and function will provide a new framework for a more comprehensive understanding of human diseases.

InSyBio ncRNASeq enables users to analyse ncRNAs. Users can search and analyse both a specific

Article history

Received: 11 December 2016

Accepted: 7 April 2017

Published: XX XX 2017

© 2017 Korfiati *et al.*; the authors have retained copyright and granted the Journal right of first publication; the work has been simultaneously released under a Creative Commons Attribution Licence, which allows others to share the work, while acknowledging the original authorship and initial publication in this Journal. The full licence notice is available at <http://journal.embnet.org>.

RNA sequence of interest and a full-sequence data-set derived from available online databases or produced by experimental or computational procedures. With InSyBio ncRNASeq, users can predict and analyse ncRNAs and miRNA target genes. InSyBio ncRNASeq also allows result storage in its knowledge-base, equipped with information retrieval tools, to allow users to produce and extract their own data-sets.

Methodologies

RNA characterisation and miRNA prediction

InSyBio ncRNASeq enables users to characterise ncRNA sequences. A set of 58 informative features is calculated by supplying the sequence in FASTA format. The features include sequence, thermodynamic and structural properties of the RNA sequences; they were first introduced in (Kleftogiannis *et al.*, 2015), and are presented in [Supplementary Table 1](#). Batch calculations of many sequences are allowed. This enables users to analyse data-sets derived from online databases, experimental sequencing or computational techniques. The results are presented in a browseable table in the Web interface, and can also be downloaded as a tab-delimited text file. For each ncRNA, its sequence and values of its 58 features are presented.

An additional functionality provided by InSyBio ncRNASeq is the prediction of pre-miRNAs, discriminating them from pseudo-hairpins and other molecules. The prediction of pre-miRNAs and pseudo-hairpins is accomplished through the application of a novel methodology, which combines a Genetic Algorithm (GA) (Holland, 1995) with epsilon-Support Vector Regression-SVR (Smola and Schölkopf, 2004) techniques. GAs, a state-of-the-art meta-heuristic optimisation technique, are used to optimise the feature subset to be used as input for the classification model and the regularisation parameters (C), the width of the Gaussian distribution (sigma) of the Radial Basis Function, and the epsilon threshold of e-SVR models. The accuracy of this technique in predicting pre-miRNAs is 95%. A sequence is predicted as 'other' if the minimum free energy is more than -15 kcal/mol, or the number of base pairs is less than 18. The supported input file comprises RNA sequences in FASTA format; batch calculations of many sequences are also allowed. For each ncRNA, the system provides the sequence, its calculated confidence score, the prediction whether it is an miRNA, a pseudo-hairpin or 'other', and its 58 features.

miRNAs and transcripts knowledge-base

The InSyBio ncRNASeq knowledge-base includes stem-loop and mature miRNAs. For each stem-loop, InSyBio ncRNASeq provides the following information: accession ID, name, species, length, description and

comments (if any). For the sequences, the FASTA format is downloadable, and the sequence, its description and secondary structure (in dot-bracket notation) are presented. Visualisation of the secondary structure is performed with FornaContainer (<https://github.com/pkerpedjiev/fornac>), which shows the Minimum Free Energy (MFE) structure. Stem-loops are linked to their corresponding mature miRNAs. For each mature miRNA, its accession ID, name and sequence are presented. The miRNA sequence is downloadable in FASTA format. Additional information about each mature miRNA is evidence, which may be experimental, similarity to another stem-loop structure, or found in the literature. References for the miRNA of interest are also available, as well as external links to other databases, such as MIRBASE (Griffiths-Jones *et al.*, 2006), ENTREZGENE (Maglott *et al.*, 2007), HGNC (Bruford *et al.*, 2008), RFAM (Griffiths-Jones *et al.*, 2005), and to publications.

The InSyBio ncRNASeq knowledge-base also includes genes and transcripts. For each gene, the information reported is the name, source, description, location, biotype, annotation method and version, as well as the list of its transcripts and links to each transcript. For each transcript, the information provided is the name, source, description, corresponding gene, corresponding protein, location, transcription start site, length, biotype, annotation method, version and 3'-UTR sequence. The 3'-UTR sequence is visualised with FornaContainer, and available for download in FASTA format.

miRNA target analysis and prediction

InSyBio ncRNASeq enables users to calculate 124 features for every miRNA and potential target-site pair within an mRNA. These features include sequence, thermodynamic and structural properties of the miRNA:mRNA pair; a detailed description of them can be found at (Korfiati *et al.*, 2015). Users can supply a file of miRNA sequences and a file of the respective mRNA binding-site sequences, both in FASTA format. Feature calculation of many miRNA:mRNA pairs is allowed in batch. The results are presented in a browseable table in the Web interface, and can also be downloaded as a tab-delimited text file. For each miRNA:mRNA pair, the miRNA sequence, the mRNA binding-site sequence and the 124 miRNA:mRNA pair features are provided. A description of the supported features for the characterisation of miRNA:mRNA pairs is available in [Supplementary Table 2](#).

With InSyBio ncRNASeq, users can computationally validate miRNA target sites with an intelligent computational technique (hybrid combination of GAs and e-SVRs) and 124 informative features, based on the method presented in (Korfiati *et al.*, 2015), by supplying a file of miRNA sequences and a file of the respective mRNA binding-site sequences in FASTA format. The ncRNASeq tool then creates the miRNA:mRNA pairs by combining all the mRNA target sites from the first

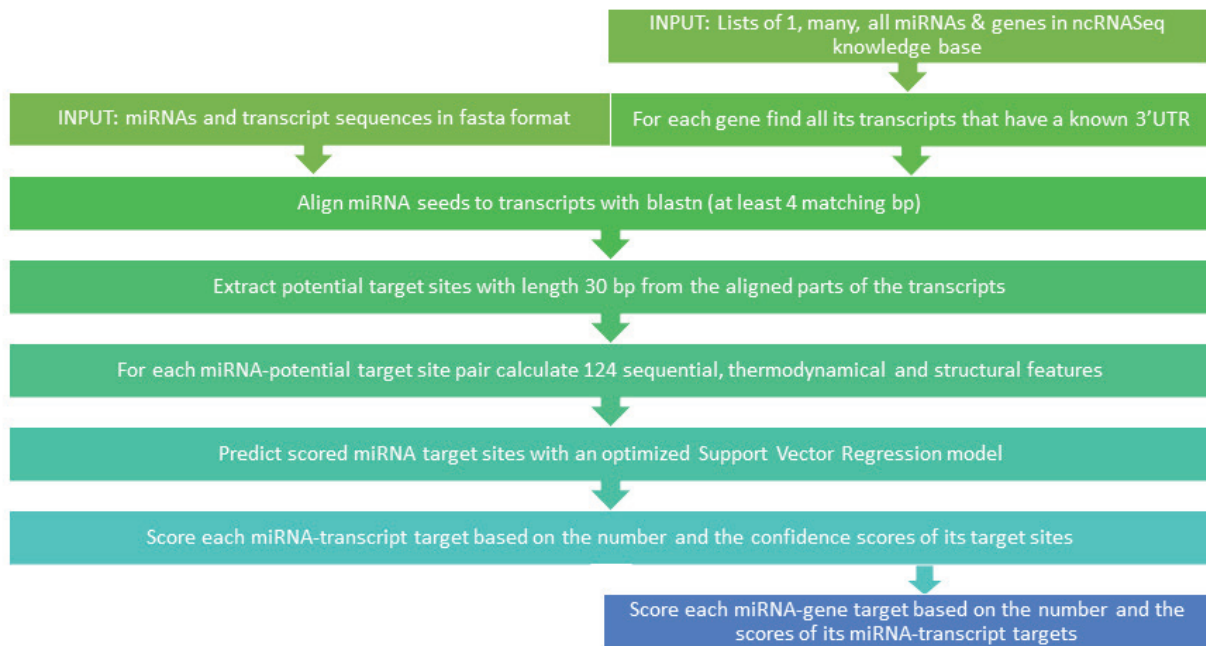


Figure 1: InSyBio ncRNASeq miRNA target-prediction workflow.

file with all the miRNAs from the second file. For each miRNA:mRNA pair, InSyBio ncRNASeq provides information about the miRNA sequence, the mRNA binding-site sequence, whether the miRNA:mRNA pair shares a targeting relation or not, the confidence score of the prediction, and the 124 miRNA:mRNA pair features grouped in categories (Thermodynamic, Positional, Motif and Structural). A link from the given miRNA to a protein in the InSyBio Interact database is also presented, indicating that the specific miRNA regulates the expression of the respective mRNA. InSyBio Interact contains information on protein functions, clusters, protein-protein interactions, etc.

Apart from miRNA target-site prediction, InSyBio ncRNASeq enables users to perform miRNA target prediction, either selecting genes or transcripts and miRNAs from the InSyBio ncRNASeq knowledge-base or providing their own miRNAs and transcript sequences in FASTA format. miRNA targets can be predicted in the 3'-UTR sequence of all transcripts of a gene for which a confidence score for the miRNA-gene interaction and additional scores for the miRNA-transcript interactions are provided. Figure 1 shows the miRNA target-prediction workflow.

Results and Conclusions

As a case study, InSyBio ncRNASeq was used to predict the miRNAs that target the 24 biomarkers reported in (Theofilatos *et al.*, 2016) for Parkinson's Disease (PD). For these biomarkers, InSyBio ncRNASeq yielded 18,359 interactions from 1,594 miRNAs (Table 1). miRTarBase (Chou *et al.*, 2016), the experimentally validated miRNA-target interaction database, contains 585 interactions for the same genes with 489 miRNAs. The results of InSyBio ncRNASeq overlap those of MiRTarBase in 190 interactions (32.5%). The miRNA target prediction tool miRTar (Hsu *et al.*, 2011), which utilises four of the most well-known target-prediction algorithms ((TargetScanS (Lewis *et al.*, 2005), miRanda (Bino *et al.*, 2004), RNAhybrid (Krüger and Rehmsmeier, 2006) and PITA (Kertesz *et al.*, 2007), match only 16 interactions (2.7%) of the experimentally validated interactions in miRTarBase. Figure 2 shows the interaction network among the highest-scoring interactions (10%) predicted by InSyBio ncRNASeq.

These results demonstrate that InSyBio ncRNASeq outperforms other tools in predicting miRNA targets. The predicted interactions are scored, which helps biologists to better design their experiments. Additional

Table 1: miRNA target-prediction results for 24 biomarkers of Parkinson's Disease obtained using InSyBio ncRNASeq and miRTar, compared with experimentally validated interactions in miRTarBase

	Interactions	miRNAs	Genes	Matching results in miRTarBase	Matching results in miRTarBase / miRTarBase total interactions	Matching results in miRTarBase / total predicted interactions
miRTarBase	585	489	24	–	–	–
miRTar	1,430	683	24	16	0.027	~0.011
ncRNASeq	18,359	1,594	24	190	0.325	~0.010

