



1

Comprehensive reconstruction and visualization of non-coding regulatory networks

Vincenzo Bonnici¹, Francesco Russo^{2,3}, Nicola Bombieri¹, Alfredo Pulvirenti^{4,*} and Rosalba Giugno^{4,*}

¹Department of Computer Science, University of Verona, Verona, Italy

²Laboratory of Integrative Systems Medicine (LISM), Institute of Informatics and Telematics (IIT) and Institute of Clinical Physiology (IFC), National Research Council (CNR), Pisa, Italy

³Department of Computer Science, University of Pisa, Pisa, Italy

⁴Department of Clinical and Molecular Biomedicine, University of Catania, Catania, Italy

Correspondence*:

Rosalba Giugno

Department of Clinical and Molecular Biomedicine, University of Catania, Viale Andrea Doria 6, Catania, 95125, Italy, giugno@dmi.unict.it

Alfredo Pulvirenti

Department of Clinical and Molecular Biomedicine, University of Catania, Viale Andrea Doria 6, Catania, 95125, Italy, apulvirenti@dmi.unict.it

Bioinformatics of Non-Coding RNAs with Applications to Biomedicine: Recent Advances and Open Challenges

2 ABSTRACT

3

4 **Keywords:** microRNAs, lncRNAs, non-coding RNAs, networks, cytoscape, gene expression

1 INTRODUCTION

5 After the sequencing of the human genome it became evident that only 20,000 genes are protein-coding,
6 while over 98% of all genes are untranslated non-protein-coding RNAs (ncRNAs) (**Consortium et al.**
7 (2012)). During the last years thousands of ncRNAs were identified in the eukaryotic transcriptome (**Bu**
8 **et al.** (2011); **Khalil et al.** (2009)). Usually, ncRNAs are divided into two groups according to their length:
9 short ncRNAs (e.g. microRNAs (miRNAs)), consisting of less than 200 nucleotides, and long non-coding
10 RNAs (lncRNAs), sized from 200 nucleotides up to 100 kb (**Mattick** (2001)).

11 miRNA is the best known class of small RNAs. They are endogenous non-coding regulatory RNAs
12 (17–25 nucleotides), which play important roles in post-transcriptional gene regulation. miRNAs bind
13 the 3'-untranslated regions (3'-UTR) of various target mRNAs leading to direct mRNA degradation or
14 translational repression.

15 They regulate gene expression and contribute to development, differentiation and carcinogenesis. Many
16 studies have shown that over 30% of human genes are regulated by miRNAs and a single miRNA can
17 control over hundreds of mRNAs. The aberrant expression or alteration of miRNAs also contributes to

18 a range of human pathologies, including cancer (**Lu et al. (2005)**). Moreover, a significant amount of
19 miRNAs has been found in extracellular human body fluids (**Mitchell et al. (2008)**; **Hanke et al. (2010)**)
20 and some circulating miRNAs in the blood have been successfully revealed as biomarkers for several
21 diseases including cardiovascular diseases (**Gupta et al. (2010)**) and cancer (**Mitchell et al. (2008)**).
22 An emerging class of non coding RNAs consists of long non-coding RNAs (**Fatica and Bozzoni (2013)**).
23 They are both nuclear and cytoplasmic. Nuclear lncRNAs function by guiding chromatin modifiers to
24 specific genomic loci (**Rinn and Chang (2012)**; **Batista and Chang (2013)**; **Guttman and Rinn (2012)**;
25 **Khalil et al. (2009)**) while many others have been identified in the cytoplasm (**Batista and Chang**
26 **(2013)**). These lncRNAs are involved in gene regulation and often show sequence complementarity with
27 transcripts that originate from either the same chromosomal locus or independent loci.
28 Interesting, the expression of lncRNAs has been quantitatively analysed in several tissues and cell types
29 by RNA-seq experiments, and it was generally found to be more cell type specific than the expression
30 of protein-coding genes (**Rinn and Chang (2012)**; **Derrien et al. (2012)**; **Guttman and Rinn (2012)**;
31 **Mercer et al. (2008)**; **Cabili et al. (2011)**; **Pauli et al. (2012)**).
32 One of the last classes recently discovered is the circular RNA (circRNAs) (**Memczak et al. (2013)**). They
33 are an enigmatic class of RNA with unknown function. Numerous circRNAs form by head-to-tail splicing
34 of exons, suggesting previously unrecognized regulatory potential of coding sequences. Recent results
35 (**Memczak et al. (2013)**) have been shown that thousands of well-expressed, stable circRNAs, often have
36 tissue/developmental-stage-specific expression. Moreover, human circRNAs are bound by miRNAs such
37 as the miR-7 showing a potential role of circRNAs as post-transcriptional regulators.
38

39 Understanding the complex system derived from the interactions of regulators and possible targets gives
40 a clue on the dynamics and causes of disorders (**Couzin (2007)**). In this direction, platforms to visualize
41 networks such as Cytoscape (**Shannon et al. (2003)**) together with tools to visualize the topology of
42 regulatory networks are increasing their influence in science.
43 miRScope (**Ferro et al. (2009)**) is one of the first Cytoscape plugin visualizing protein-protein interaction
44 networks annotated with miRNAs. It uses a web knowledge base (**Laganà et al. (2009)**) to infer
45 associations between genes and phenotypes through miRNAs. CyTargetLinker (**Kutmon et al. (2013)**)
46 is a recent Cytoscape app which builds biological networks annotated with miRNAs, transcription factors
47 and drugs.
48 Several tools are designed to analyze the regulatory effect of miRNAs in protein coding genes and to
49 export the results in a Cytoscape network format. For example, Magia (**Sales et al. (2010)**) is a web
50 tool to statistically analyze miRNA and gene expressions. TSmir is a web based tool to browse regulatory
51 network of tissue-specific miRNAs with transcription factors. mirConnX **Huang et al. (2011)** extend a given
52 network of genes, transcription factors and miRNAs with further TF and miRNA-gene intersections inferred
53 by user expression data. miRTrail **Laczny et al. (2012)** statically analyzes the role of miRNAs and genes
54 deregulated in a disease by using an extensive miRNA-gene networks and expression data.
55

2 CONSTRUCTION AND CONTENT

2.1 DATA SOURCE

56 In the following sections we describe data sources considered in ncRNAdb.

57 *2.1.1 HGNC* The HUGO Gene Nomenclature Committee (HGNC) is responsible for approving unique
58 symbols and names for human loci, including protein coding genes, ncRNA genes and pseudogenes, to
59 allow unambiguous scientific communication (<http://www.genenames.org/>) (**Gray et al. (2012)**).
60 It contains 19060 protein coding genes, 5714 non-coding RNAs (of those 2546 long non-coding RNAs
61 and 1879 microRNAs), 598 phenotypes and 12621 pseudogenes.
62 We used HGNC as official database of approved names and alias for each entry in ncRNAdb.

63 **2.1.2 *lncRNAdb*** lncRNAdb provides a central repository of known long non-coding RNAs (lncRNAs)
64 in eukaryotic cells (including those derived from viruses), their aliases and published characteristics
65 (**Amaral et al.** (2011)). lncRNAdb is available online at <http://www.lncrnadb.org>.
66 The published version of the database contains over 150 lncRNAs identified from the literature in around
67 60 different species.

68 **2.1.3 *circ2Traits*** Circ2Traits (<http://gyanxet-beta.com/circdb/>) is a comprehensive database for
69 circular RNA potentially associated with disease and traits (**Ghosal et al.** (2013)).
70 Circular RNAs, formed by covalent linkage of the ends of a single RNA molecule, are newly discovered
71 RNAs that sponge miRNAs to block their function (**Memczak et al.** (2013)).
72 Circ2Traits uses the circular RNA dataset from Memczak et al. (**Memczak et al.** (2013)). This dataset
73 consists of 1953 predicted human circular RNAs along with their genomic coordinates, annotation, and
74 predicted miRNA seed matches. The disease related miRNA data is taken from miR2disease (**Jiang et al.**
75 (2009)). From Starbase (**Yang et al.** (2011)) the authors collected the miRNA-mRNA interaction data
76 predicted by miRanda (**Betel et al.** (2008)), TargetScan (**Lewis et al.** (2005)), PiTA (**Kertesz et al.**
77 (2007)), PicTar (**Krek et al.** (2005)), and RNA22 (**Loher and Rigoutsos** (2012)).
78 Moreover, dataset for predicted miRNA and long non-coding RNA interaction pairs is collected from the
79 miRCode database (**Jeggari et al.** (2012)).

80 **2.1.4 *HMDD*** The Human microRNA Disease Database (HMDD) (**Li et al.** (2013)) is a database
81 of curated experiment-supported evidence for human microRNAs (miRNAs) and disease associations
82 (<http://www.cuilab.cn/hmdd>).
83 Currently, HMDD collected 10368 entries that include 572 miRNA genes, 378 diseases from 3511 papers.
84 The database contains detailed and comprehensive annotations of human miRNA-disease association
85 data, including miRNA-disease association data from the evidence of genetics, epigenetics, circulating
86 miRNAs, and miRNA-target interactions.

87 **2.1.5 *lncRNADisease database*** The lncRNADisease database (<http://210.73.221.6/lncrnadisease>)
88 (**Chen et al.** (2013)) is a resource for the experimentally supported lncRNA-disease association data and
89 a platform that integrates tools for predicting novel lncRNA-disease associations.
90 Moreover, lncRNADisease contains lncRNA interactions in various levels, including proteins, RNAs,
91 miRNAs, and DNA.
92 The current version of the database has more than 1000 lncRNA-disease entries and 475
93 lncRNA interaction entries, including 321 lncRNAs and 221 diseases from about 500 publications.
94 lncRNADisease also provides the predicted associated diseases of 1564 human lncRNAs.

95 **2.1.6 *miRandola*** miRandola (<http://atlas.dmi.unict.it/mirandola/>) (**Russo et al.** (2012); **Russo et al.**
96 (2014)) is the first manually curated database of extracellular circulating miRNAs. It is a comprehensive
97 classification of different extracellular miRNA types and a collection of non-invasive biomarkers for
98 several diseases (e.g. cancer and cardiovascular disease).
99 The last updated version of the database contains 139 papers, 2366 entries, 599 unique mature miRNAs
100 and several tools for cellular and extracellular miRNA analysis.

101 **2.1.7 *NONCODE*** NONCODE (<http://www.noncode.org/>) is a database of all kinds of non-coding
102 RNAs (except tRNAs and rRNAs), in particular it contains 210831 lncRNAs for several species (**Bu**
103 **et al.** (2011)).
104 NONCODE also provides an ID conversion tool from RefSeq or Ensembl ID to NONCODE ID and a
105 service of lncRNA identification.

106 2.1.8 *NPInter* NPInter (<http://www.bioinfo.org/NPInter/>) (Wu et al. (2006)) reports functional
107 interactions between non-coding RNAs (except tRNAs and rRNAs) and biomolecules (proteins, RNAs
108 and DNA) which are experimentally verified.
109 The authors collected primarily physical interactions, although several interactions of other forms also are
110 included.
111 Interactions are manually collected from publication, followed by an annotation process against known
112 databases including NONCODE (Bu et al. (2011)), miRBase (the miRNA registry) (Kozomara and
113 Griffiths-Jones (2013)) and UniProt (the database of proteins) (Consortium et al. (2013)).
114 The second version of NPInter contains 201107 interactions of ncRNA with other biomolecules from 18
115 organisms.

2.2 DATA SCHEMA

116 Inside public databases, biological entities (e.g. genes) are catalogued via common used nomenclatures,
117 they can be human readable names or alphanumeric identifiers (that include accession numbers), with
118 the intent to assign a keyword that can identify the specific entity. For example, genes are classified by
119 their names (usually, a summary description of their function), by their symbols (short abbreviations) or
120 database-specific identifiers. The *breast cancer 1* gene can be identified by its assigned symbols BRCA1,
121 BRCC1 and PPP1R53, or by its specific identifier among gene databases like HGNG:1100, Entrez Gene
122 672 and UCSC uc002ict.3. Unfortunately, some of these nomenclatures are not disambiguous, thus
123 distinct biological entities can be identified by the same nomenclature. The gene example is a lucky
124 one, like the disease nomenclature, since, in the last years, researches have made an effort for their
125 disambiguation, like the HUGO Gene Nomenclature Committee. Non-coding RNAs, as we present here,
126 are a relatively recent discovery and there is not yet a comprehensive database including all of them and
127 where the researcher can find an unambiguous nomenclature systems. Moreover, the non-coding RNA
128 knowledge is spread among several databases where new discovered entities are named with internal
129 identifiers and poor mapping schemas between such databases exist. The largest collection of ncRNA is
130 NONCODE v4 but the mapping, from internal to external ID systems, that it provides is still not extensive
131 also because most the reported ncRNAs are uniquely to NONCODE and they are not reported in any other
132 database.
133 For our purpose, we avoided the use of an our internal nomenclature system but we introduced a generic
134 resource identifier system. Each identifier is composed by three parts, also called levels: the entity type,
135 the data source and the external name, respectively.

$$\textit{entitytype} : \textit{datasource} : \textit{alias} \quad (1)$$

136 The entity type describes the biological classification of the element, while the data source describes the
137 external data resource from which the information come from. An entity type can be NCRNA, DISEASE,
138 GENE or RNA (which does not include ncRNA). A data source, or resource, names is composed by
139 the name of the external resource plus its version (e.g. HMDD_2). We also consider two additional data
140 source identifiers, GENERIC and UNKNOWN, for dealing with generic names (e.g. the symbol BRCA1
141 is a generic used one that can not be related to any specific resource) and unknown (to us) resources,
142 respectively. The resource identifiers is also used to store the origin of a relation, namely the databases
143 from where it has been extracted. Table ?? shows the complete set of resource identifiers used in our
144 system. Moreover, the database engine (OrientDB) used for out system labels each item with a unique
145 system-scope identifier, called ORID, and every record (entity or relation) can be retrieved by its ORID.
146 Thus, our record identify format allows to retrieve and entity by its ORID in place of a three level identifier.
147

148 A system representing a set of biological entities and their relations (physical interactions, functional
149 relations, etc...) can be naturally modelled within a network, or formally described as a graph, a
150 mathematical object composed by vertices (entities) and edges (relations). Classical relation database
151 management systems (RDMS) are widely used to store biological data. Such systems are based on the
152 relational model which is designed to represent an entity-relation system, and often the SQL (Structure

153 Query Language) language is used to provide a querying interface to data. However, newly models have
154 been proposed to solve performance issues or to provide additional feature in specific applications. Some
155 of them are grouped under the name of NoSQL (Not only SQL) databases and they are increasingly used
156 in big data and web applications, since they can provide schema-less representation for non-structured
157 data. Nowadays, there a re plenty of such systems, each one with specific features and data model for data
158 representation. They can be summary classified by four main data models, even if some of them include
159 more then one model: column model where data are represented by tuples, document-oriented databases
160 for storing, retrieving, and managing document-oriented information, also known as semi-structured
161 data, key-value model where data are stored as a collection of key-value pairs stored using associative
162 arrays, maps, symbol tables of dictionaries, and finally graph databases where data are modelled using
163 a graph structure. Moreover, they often implements the object-oriented model since they are used as
164 persistence layer in developing software using object-oriented languages, thus they model concepts like
165 class, inheritance and polymorphism. Functionality, complexity, flexibility, scalability and performance
166 depend on the suitability of the implemented data models to specific applications. Usually, RDMS systems
167 store data into tables and provide a single table index (based on the key columns of the table) for querying
168 data in an efficient way. Instead, a pecuniary feature of NoSQL systems is that several indexes can be built
169 on the same data container, proving different access point to the same data.
170

171 Among the several NoSQL databases, OrientDB has been chosen as persistence layer of our system. It
172 implements a graph model and an object-oriented model, on top of a document model. Moreover, it offers
173 an SQL interface in addition to several language specific interfaces. It is developed in Java and provides
174 native Java API (Application Programming Interface) for accessing the database, which is suitable for
175 developing Cytoscape applications. Figure 1 shows the schema of our proposed systems. We preferred
176 to keep separated the concepts of biological entity and its related aliases. The abstract class *BioEntity*
177 represents biological entities and it specializes in the four biological entity types ncRNA, Gene, Disease
178 and RNA. Aliases are represented by the abstract class *Alias* which is specialized in four different sub-
179 classes related to the four biological entity types. Each alias stores the coming resource and the internal ID
180 (in such resource), thus the complete three levels identifier can be reconstruct combining those two fields
181 with the specific alias sub-class where they are stored. In a graph model, elements on both classes (and
182 sub-classes) represent vertices of the graph having distinct properties. The naming of a biological entity
183 by an alias is represented by adding an edge between the corresponding graph vertices. Aliases serve as
184 access points to the biological entity collection. Two different indexes are built on top of each alias sub-
185 class. The first one is a composite key dictionary working on the two fields, resource and ID, to provide an
186 unique identification of an external reference, and it implements the fully three levels alias identification.
187 The second one is a single field not-unique map to efficiently query the set of IDs of an alias sub-class
188 without combining them to the external resource identification. This index results useful when the origin
189 (the resource where it comes from) of an ID is not known, thus the system outputs all the aliases having
190 such ID regardless the external resource system that produced them. In this way, we can implement a
191 two levels (ambiguous) identification system based on the biological type and a generic name. A further
192 not-unique index works on the entire *Alias* class hierarchy to implement a single level identification when
193 neither the biological entity type ans the external resource are known. This index was introduced because
194 sometimes some of the imported databases do not specify the biological entity type and we were not able
195 to catalogue some elements because no matches in other databases where found or because they were
196 ambiguous. Interactions are modelled as graph edges between two elements of the class *BioEntity*. For
197 each interaction ad additional set of information is stored, when the imported databases included them. It
198 include a PubMed ID of the referencing article plus the phrase supporting the relation, the database from
199 where the interaction was extracted and the interaction level. The interaction level is composed by a pair
200 of elements within $\{RNA, DNA, Protein, TF\}$, on *NA* (when it is not specified). For example, the pair
201 *RNA – TF* specifies that a ncRNA is interacting with the transcription factor of a gene, while the pair
202 *RNA – DNA* specifies that the ncRNA is interacting with the coding genomic region of the gene. If the
203 same relation is stored in two (or more) distinct data sources, than two (or more) interaction edges are
204 present inside our system.

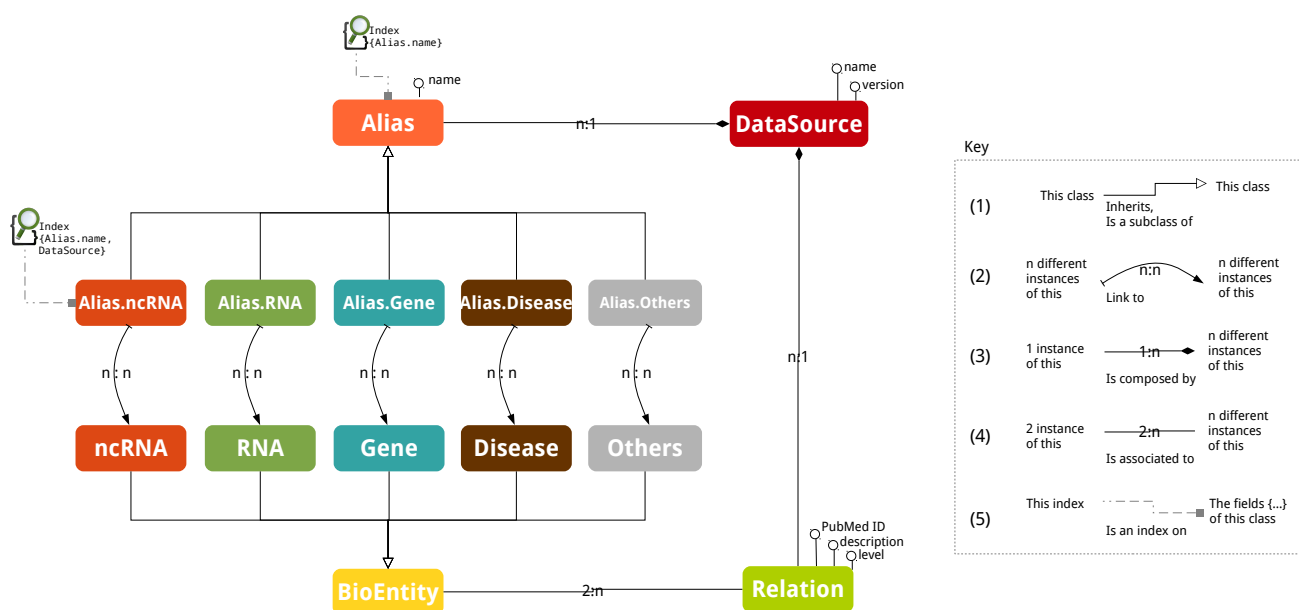


Figure 1. The database schema.

2.3 DATA IMPORT

205 In this section, we describe the import phase from the external databases listed in section ?? to our system.
 206 The main issue of such operation was the alias ambiguity and the entity type miss-knowledge. The first
 207 one is described in section 2.2 and it has been resolved by a merging procedure. Formally, given three
 208 aliases A, B and C such that two different mappings $A \mapsto B$ and $B \mapsto C$, possibly coming from different
 209 data sources, have been found. Then the three aliases are considered as alternative nomenclatures to the
 210 same biological entity. This procedure does not prevents that the first two aliases are referred to a different
 211 biological entities of the other two, for example when a gene name is used to identify its transcript, but
 212 reduces the number of mismatching where querying the system. The miss-knowledge issue was due to
 213 the fact that the imported interaction data were in the form ncRNA-(something else), and the biological
 214 type of the second interaction actors was not specified, neither the indication of the external nomenclature
 215 system. Thus, we were not able to fully classify (with a complete three level identification) the entity,
 216 and it has been catalogued as *UNKNOWN* type. Some of such miss-classification were resolved by
 217 applying a procedure after the import step, when all the data from the external sources were imported and
 218 more information was available, by trying to reclassify the *UNKNOWN* type entities.
 219

220 Below, a list of the importing data source and the type of information, extracted from them, is given.
 221 The order is the same in which with extract and add up information from a source to our collection. For
 222 all the importing source, we considered only those information related to the species *Homo sapiens*.
 223 From HGNC (the database of human gene names) we import a list of non-coding RNA, protein-coding
 224 genes, pseudogenes and phenotypes (considered as diseases) and their approved aliases in a series of well-
 225 established databases.
 226 From lncRNADB we import a list of non-coding RNA and their common names (aliases, symbols). This
 227 source includes also some briefly description about expression, function, conservation and misc that we
 228 do not take into account.
 229 From circ2traits, we import a set of interacting lncRNA, circRNA and messenger RNA together with the
 230 disease, relation and article where the interaction is reported.
 231 From HMDD, we import a list of disease names together with a set of genes that interact with a ncRNA.

232 In this database, interactions are listed in the forms ncRNA-disease or ncRNA-gene-disease, and the
233 referencing article together the support sentence are reported. The ncRNA-gene-disease multi-relations
234 are split into distinct relation of type ncRNA-gene and ncRNA-disease.
235 From lncRNAdisease, we import a list of lncRNA (and their symbol aliases) and diseases and their
236 (co-)relation supported by a pubmed reference, a dysfunction type description and a support sentence.
237 Moreover, from this database, we import also a set of entities interactions supported by the article,
238 dysfunction type, support sentence and interaction level. In this case, since the entity type the actors
239 of the interaction is not provided, but only the entity level (i.e. RNA-protein) at first we query our system
240 (filled with the information extracted until this step) searching for a corresponding entity to the given alias,
241 and if no entry is found then the new alias is stored and connected to a new entity of unknown type.
242 From the Mirandola database, we import a set of miRNA (and their aliases) and their related disease
243 together with the supporting article and sentence.
244 From NONCODE, we import a list of non-coding RNA aliases and a mapping from the NONCODE v4
245 DI system to other ID systems. Moreover, we import also a set of associated GO terms related to each
246 listed ncRNA.
247 From NPInter, we import a set of ncRNA and their interactions with other entities supported by the
248 referencing article, support sentence and interaction level. Similar to the lncRNAdisease case, the entity
249 type of the second actor of the interaction is not specified thus at first we search for an entry in our database
250 and if the entity alias is still not reported we insert a new unknown-type entity and its alias.
251

252 At the end of the import procedure, 853,543 three-level identifiers (alias) were created, a total of
253 222,970 biological entities were imported, and 889,675 edges between *Alias* and *BioEntity* classes were
254 created. Moreover, 238,524 entity relation were extracted. Table ?? shows the total number of imported
255 biological entities, sorted by they biological type, and how many of them are actually involved in a
256 relation.

3 UTILITY AND DISCUSSION

257 As said before, OrientDB is supported by plenty of language connectors, beside the native Java API. One
258 a server instance is up, the final user can connect to it by a series of programming language binding, or
259 by using the OrientDB SQL console. It also implements technology standard like HTTP REST/JSON,
260 TinkerPop Blueprints (for graph computing) and JDO (Java Data Object for object persistence).
261 However, we propose three alternative interfaces available to final users for querying our public database
262 instance. A CytoScape (version 3) plug-in for importing data in a network environment, plus a web
263 interface and a command-line interface for raw resource queries. The CytoScape plug-in and the
264 command-line applications can be downloaded from the ncRNA-DB website, while the web interface
265 is part of it. A complete schema of the proposed system, from the import phase to the user-end interfaces,
266 is shown in figure 2.

3.1 CYTOSCAPE INTERFACE

267 The CytoScape interface is composed by three main functions, integrated in a single CytoScape plug-
268 in, for expanding an input network with ncRNA-DB data, searching for specific biological entities, or
269 importing the whole ncRNA-DB in a CytoScape network environment. Figure 3 shows the graphical user
270 interfaces of such applications, that can be accessed through the *ncRNA-DB* menu item, once the plug-in
271 have been installed.
272

273 The *Expand Network* functionality takes in input a user built network and expand (add to) it ncRNA-DB
274 entities and relations. It is required that the input network table must have two columns for specifying
275 the biological entity type of an element and for listing its set of known aliases in their ncRNA-DB

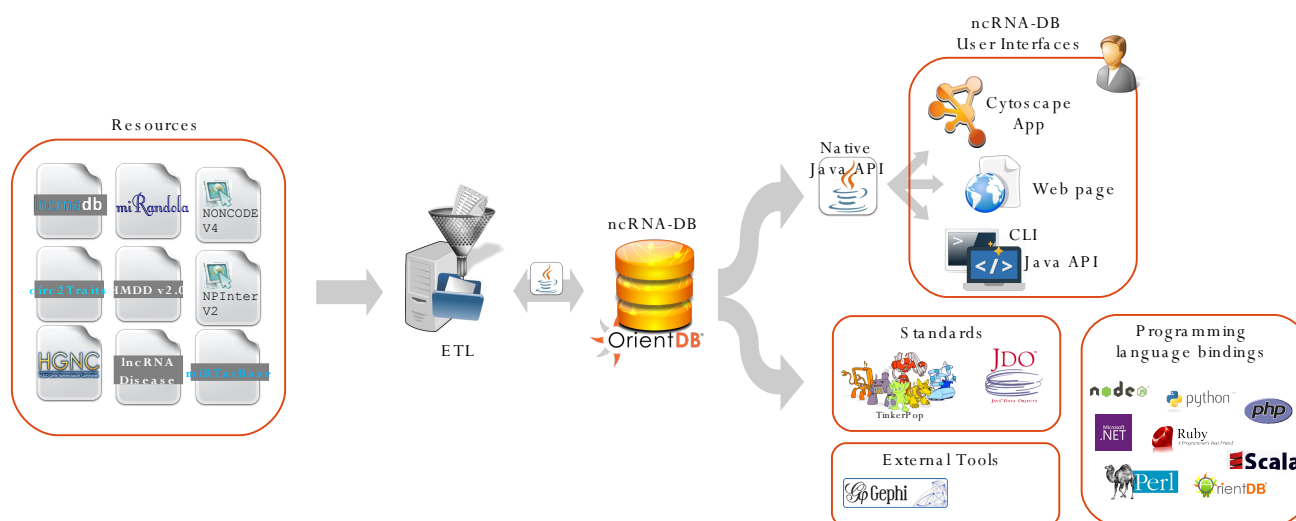


Figure 2. Complete system architecture and procedures.

276 identify format. The entity type can be one of *NCRNA*, *RNA*, *GENE* or *DISEASE* for fully qualified
 277 resource retrieving, or it can be a blank field, in this case the application will search for all the biological
 278 entity types. If the key *GROUP* is used as entity type, then the application associates all the matching
 279 aliases to the node. The user can decide whenever some biological entity types must be excluded from the
 280 retrieving by unflagging the corresponding check-boxes. Once those two columns have been specified, the
 281 application retrieves the matching biological entities, maps them to the CytoScape nodes and adds to the
 282 selected network the set of found relations between the retrieved nodes. If the aliases set of a CytoScape
 283 node matches with more than one ncRNA-DB entity, then the node is mapped to all of such entities. This
 284 behaviour allows the user to specify abstract nodes representing an entity group, for example, one can
 285 specify a group of aliases of different biological type and compact them in a single entity node. It also
 286 helps to eliminate ambiguity at this level of data representation.
 287 If the *Include neighbours* check-box is flagged, then the application retrieves all the ncRNA-DB
 288 neighbours of the matching entities and adds them to the selected CytoScape network, as well as relations
 289 between them and the other retrieved entities.
 290

291 The second function, *List entities*, allows the user to search for particular biological entities by
 292 specifying their aliases (fully or incomplete ncRNA-DB identifiers). Aliases can be put in the text form of
 293 the application or loaded from file. If aliases are specified by the text form, then a single CytoScape node
 294 is created, instead, if they are specified from file then a node is created for each file row. The file format
 295 expects a series of rows, each of which can contain one or more aliases. The application retrieves the
 296 matching entities and the user can select a list of them to be imported in the selected CytoScape network.
 297 For each entity the complete list of corresponding aliases, those taken in input plus those retrieved from
 298 the database, and its biological type are shown. Figure 4 shows the querying page of the proposed web
 299 interface.
 300

301 Finally, the *Import DB* application allows to import the whole ncRNA-DB database, or just a part of
 302 it by specifying a series of discriminating rulers. The user can select which of the biological entity type
 303 must be imported and which kind of relations. The retrieved elements are imported in a new Cytoscape
 304 network having the user specified name. We suggest to use this functionality properly, since hundreds of

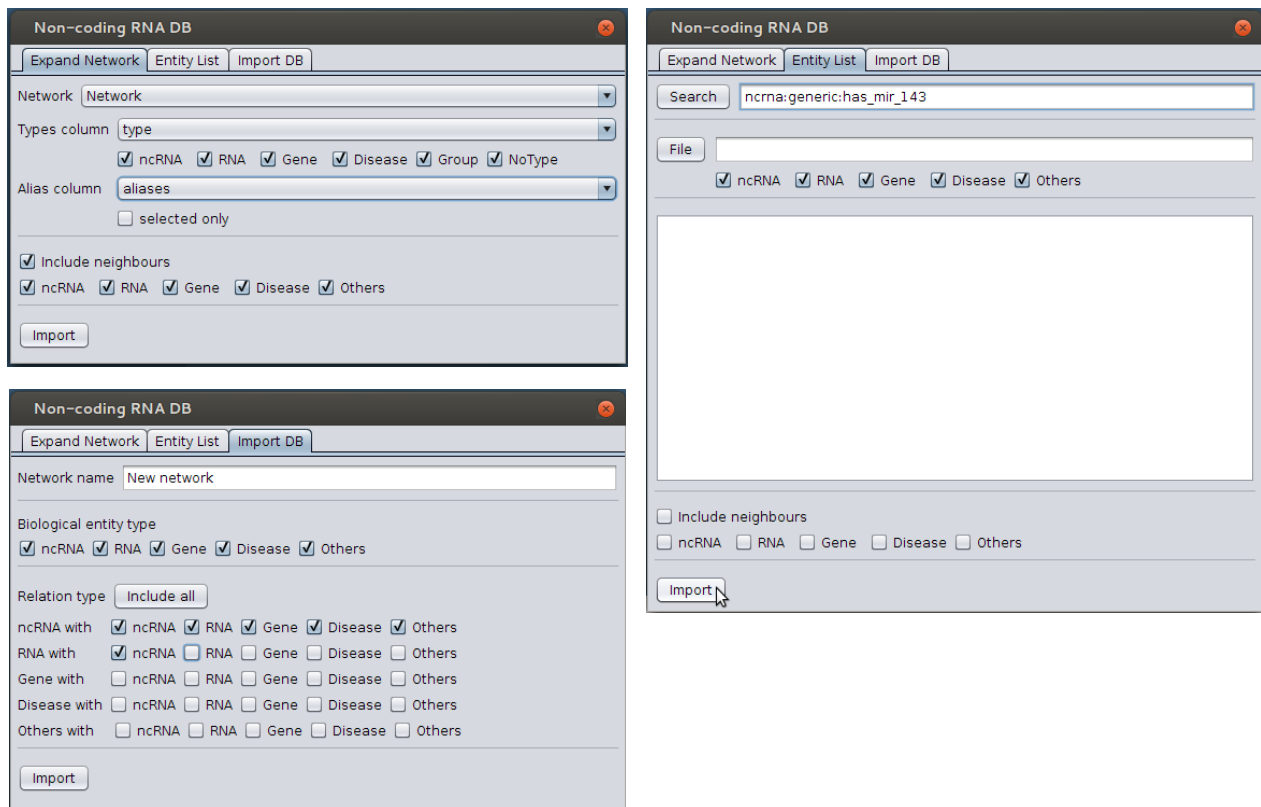


Figure 3.The Cytoscape plug-in.

Figure 4.The main page of the web interface.

305 thousands of items, both entities and relations, are stored in our databases, and the cause performance of
 306 minimal resource requirement issues when imported all together in the CytoScope environment.

3.2 WEB INTERFACE

307 We developed a web interface for querying our databases. On the main page, the user specify a set of entity
 308 identifiers (as described in section 2.2) or an internal ncRNA-DB ORID, then the matching entities are
 309 printed on the result panel. For each entity, the set of its alias is shown, as well as its internal ncRNA-DB
 310 ORID, and a link for viewing the set of interacting entities.

3.3 COMMAND-LINE INTERFACE

311 We developed a command-line interface to our database for entity searching and relation retrieval. It is
 312 released as a Java package to be platform independent and it does not require any external dependency. It
 313 provides to different commands for accessing the database data. The *search* command takes as input
 314 a list of entity coordinates and returns the matching biological entities stored in the databases. This
 315 command helps to verify if an identifiers is included in the database and to retrieve all its alternative
 316 nomenclatures which were extracted during our data import phase. Moreover, it can be used to understand
 317 how a nomenclature has been catalogued. For example, a one level coordinate, e.g. a gene alias, search

318 helps to recognize different data sources using the same identifier or whenever the alias is used for distinct
319 biological entity types. The second command, *relations*, receives as input a list of entity coordinates and
320 outputs the relations between them stored in the database and their support information. If more than one
321 biological entities match the same input coordinate, then they are all included in the output.

4 CONCLUSION

DISCLOSURE/CONFLICT-OF-INTEREST STATEMENT

322 The authors declare that the research was conducted in the absence of any commercial or financial
323 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

324 The statement about the authors and contributors can be up to several sentences long, describing the tasks
325 of individual authors referred to by their initials and should be included at the end of the manuscript before
326 the References section.

ACKNOWLEDGEMENT

327 *Funding:* Francesco Russo has been supported by a fellowship sponsored by 'Progetto Istituto Toscano
328 Tumori Grant 2012 Prot.A00GRT'.

REFERENCES

- 329 Amaral, P. P., Clark, M. B., Gascoigne, D. K., Dinger, M. E., and Mattick, J. S. (2011), lncRNADB: a
330 reference database for long noncoding RNAs, *Nucleic acids research*, 39, suppl 1, D146–D151
- 331 Batista, P. J. and Chang, H. Y. (2013), Long noncoding RNAs: cellular address codes in development and
332 disease, *Cell*, 152, 6, 1298–1307
- 333 Betel, D., Wilson, M., Gabow, A., Marks, D. S., and Sander, C. (2008), The microRNA.org resource:
334 targets and expression, *Nucleic acids research*, 36, suppl 1, D149–D153
- 335 Bu, D., Yu, K., Sun, S., Xie, C., Skogerbø, G., Miao, R., et al. (2011), NONCODE v3. 0: integrative
336 annotation of long noncoding RNAs, *Nucleic acids research*, gkr1175
- 337 Cabili, M. N., Trapnell, C., Goff, L., Koziol, M., Tazon-Vega, B., Regev, A., et al. (2011), Integrative
338 annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses,
339 *Genes & development*, 25, 18, 1915–1927
- 340 Chen, G., Wang, Z., Wang, D., Qiu, C., Liu, M., Chen, X., et al. (2013), LncRNADisease: a database for
341 long-non-coding RNA-associated diseases, *Nucleic acids research*, 41, D1, D983–D986
- 342 Consortium, E. P. et al. (2012), An integrated encyclopedia of DNA elements in the human genome,
343 *Nature*, 489, 7414, 57–74
- 344 Consortium, U. et al. (2013), Update on activities at the Universal Protein Resource (UniProt) in 2013,
345 *Nucleic acids research*, 41, D1, D43–D47
- 346 Couzin, J. (2007), Erasing microRNAs reveals their powerful punch, *Science*, 316, 5824
- 347 Derrien, T., Johnson, R., Bussotti, G., Tanzer, A., Djebali, S., Tilgner, H., et al. (2012), The GENCODE
348 v7 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and expression,
349 *Genome research*, 22, 9, 1775–1789
- 350 Fatica, A. and Bozzoni, I. (2013), Long non-coding RNAs: new players in cell differentiation and
351 development, *Nature Reviews Genetics*

- 352 Ferro, A., Giugno, R., Laganà, A., Mongiovì, M., Pigola, G., Pulvirenti, A., et al. (2009),
353 miRScape: A Cytoscape Plugin to Annotate Biological Networks with microRNAs, in Network
354 Tools and Applications in Biology (NETTAB), Focused on Technologies, Tools and Applications for
355 Collaborative and Social Bioinformatics Research and Development
- 356 Ghosal, S., Das, S., Sen, R., Basak, P., and Chakrabarti, J. (2013), Circ2Traits: a comprehensive database
357 for circular RNA potentially associated with disease and traits, *Frontiers in genetics*, 4
- 358 Gray, K. A., Daugherty, L. C., Gordon, S. M., Seal, R. L., Wright, M. W., and Bruford, E. A. (2012),
359 Genenames.org: the HGNC resources in 2013, *Nucleic acids research*, gks1066
- 360 Gupta, S. K., Bang, C., and Thum, T. (2010), Circulating microRNAs as biomarkers and potential
361 paracrine mediators of cardiovascular disease, *Circulation: Cardiovascular Genetics*, 3, 5, 484–488
- 362 Guttman, M. and Rinn, J. L. (2012), Modular regulatory principles of large non-coding RNAs, *Nature*,
363 482, 7385, 339–346
- 364 Hanke, M., Hoefig, K., Merz, H., Feller, A. C., Kausch, I., Jocham, D., et al. (2010), A robust
365 methodology to study urine microRNA as tumor marker: microRNA-126 and microRNA-182 are
366 related to urinary bladder cancer, in *Urologic Oncology: Seminars and Original Investigations*,
367 volume 28 (Elsevier), volume 28, 655–661
- 368 Huang, G., Athanassiou, C., and Benos, P. (2011), mirConnX: condition-specific mRNA-microRNA
369 network integrator, *Nucleic Acids Research (Web server issue)*, 39
- 370 Jeggari, A., Marks, D. S., and Larsson, E. (2012), miRcode: a map of putative microRNA target sites in
371 the long non-coding transcriptome, *Bioinformatics*, 28, 15, 2062–2063
- 372 Jiang, Q., Wang, Y., Hao, Y., Juan, L., Teng, M., Zhang, X., et al. (2009), miR2Disease: a manually
373 curated database for microRNA deregulation in human disease, *Nucleic acids research*, 37, suppl 1,
374 D98–D104
- 375 Kertesz, M., Iovino, N., Unnerstall, U., Gaul, U., and Segal, E. (2007), The role of site accessibility in
376 microRNA target recognition, *Nature genetics*, 39, 10, 1278–1284
- 377 Khalil, A. M., Guttman, M., Huarte, M., Garber, M., Raj, A., Morales, D. R., et al. (2009), Many
378 human large intergenic noncoding RNAs associate with chromatin-modifying complexes and affect
379 gene expression, *Proceedings of the National Academy of Sciences*, 106, 28, 11667–11672
- 380 Kozomara, A. and Griffiths-Jones, S. (2013), miRBase: annotating high confidence microRNAs using
381 deep sequencing data, *Nucleic acids research*, gkt1181
- 382 Krek, A., Grün, D., Poy, M. N., Wolf, R., Rosenberg, L., Epstein, E. J., et al. (2005), Combinatorial
383 microRNA target predictions, *Nature genetics*, 37, 5, 495–500
- 384 Kutmon, M., Kelde, T., Mandaviya, P., Evelo, C. T., and Coort, S. L. (2013), CyTargetLinker: A
385 Cytoscape App to Integrate Regulatory Interactions in Network Analysis, *PLoS One*
- 386 Laczny, C., Leidinger, P., Haas, J., Ludwig, N., Backes, C., Gerasch, A., et al. (2012), miRTrail - a
387 comprehensive webserver for analyzing gene and miRNA patterns to enhance the understanding of
388 regulatory mechanisms in diseases, *BMC Bioinformatics*, 13, 36
- 389 Laganà, A., Forte, S., Giudice, A., Arena, M., Puglisi, P., Giugno, R., et al. (2009), miR: a miRNA
390 knowledge base, *DATABASE*
- 391 Lewis, B. P., Burge, C. B., and Bartel, D. P. (2005), Conserved seed pairing, often flanked by adenosines,
392 indicates that thousands of human genes are microRNA targets, *cell*, 120, 1, 15–20
- 393 Li, Y., Qiu, C., Tu, J., Geng, B., Yang, J., Jiang, T., et al. (2013), HMDD v2. 0: a database for
394 experimentally supported human microRNA and disease associations, *Nucleic acids research*, gkt1023
- 395 Loher, P. and Rigoutsos, I. (2012), Interactive exploration of RNA22 microRNA target predictions,
396 *Bioinformatics*, 28, 24, 3322–3323
- 397 Lu, J., Getz, G., Miska, E. A., Alvarez-Saavedra, E., Lamb, J., Peck, D., et al. (2005), MicroRNA
398 expression profiles classify human cancers, *nature*, 435, 7043, 834–838
- 399 Mattick, J. S. (2001), Non-coding RNAs: the architects of eukaryotic complexity, *EMBO reports*, 2, 11,
400 986–991
- 401 Memczak, S., Jens, M., Elefsinioti, A., Torti, F., Krueger, J., Rybak, A., et al. (2013), Circular RNAs are
402 a large class of animal RNAs with regulatory potency, *Nature*, 495, 7441, 333–338

- 403 Mercer, T. R., Dinger, M. E., Sunkin, S. M., Mehler, M. F., and Mattick, J. S. (2008), Specific expression
404 of long noncoding RNAs in the mouse brain, *Proceedings of the National Academy of Sciences*, 105, 2,
405 716–721
- 406 Mitchell, P. S., Parkin, R. K., Kroh, E. M., Fritz, B. R., Wyman, S. K., Pogosova-Agadjanyan, E. L., et al.
407 (2008), Circulating microRNAs as stable blood-based markers for cancer detection, *Proceedings of the*
408 *National Academy of Sciences*, 105, 30, 10513–10518
- 409 Pauli, A., Valen, E., Lin, M. F., Garber, M., Vastenhouw, N. L., Levin, J. Z., et al. (2012), Systematic
410 identification of long noncoding RNAs expressed during zebrafish embryogenesis, *Genome Research*,
411 22, 3, 577–591
- 412 Rinn, J. L. and Chang, H. Y. (2012), Genome regulation by long noncoding RNAs, *Annual review of*
413 *biochemistry*, 81
- 414 Russo, F., Di Bella, S., Bonnici, V., Laganà, A., Rainaldi, G., Pellegrini, M., et al. (2014), A knowledge
415 base for the discovery of function, diagnostic potential and drug effects on cellular and extracellular
416 miRNAs, *BMC Genomics*, 15, 3, 1–7
- 417 Russo, F., Di Bella, S., Nigita, G., Macca, V., Lagana, A., Giugno, R., et al. (2012), miRandola:
418 extracellular circulating microRNAs database, *PLoS One*, 7, 10, e47786
- 419 Sales, G., Coppe, A., Bisognin, A., Biasiolo, M., Bortoluzzi, S., and Romualdi, C. (2010), MAGIA, a
420 web-based tool for miRNA and Genes Integrated Analysis, *Nucleic Acids Res*
- 421 Shannon, P., Markiel, A., Ozier, O., Baliga, N., Wang, J., Ramage, D., et al. (2003), Cytoscape: a
422 software environment for integrated models of biomolecular interaction networks, *Genome Research*,
423 13, 11, 2498–504
- 424 Wu, T., Wang, J., Liu, C., Zhang, Y., Shi, B., Zhu, X., et al. (2006), NPInter: the noncoding rnas and
425 protein related biomacromolecules interaction database, *Nucleic acids research*, 34, suppl 1, D150–
426 D152
- 427 Yang, J.-H., Li, J.-H., Shao, P., Zhou, H., Chen, Y.-Q., and Qu, L.-H. (2011), starBase: a database for
428 exploring microRNA–mRNA interaction maps from Argonaute CLIP-Seq and Degradome-Seq data,
429 *Nucleic acids research*, 39, suppl 1, D202–D209

FIGURES

- 430 **Figure 1.** Enter the caption for your figure here. Repeat as necessary for each of your figures.