



MicroRNA discovery in the human parasite *Echinococcus multilocularis* from genome-wide data



L. Kamenetzky^{a,*}, G. Stegmayer^b, L. Maldonado^a, N. Macchiaroli^a, C. Yones^b, D.H. Milone^b

^a IMPAM-UBA-CONICET, Facultad de Medicina, Buenos Aires, Argentina

^b sinc(i)-FICH-UNL-CONICET, Ciudad Universitaria, Santa Fe, Argentina

ARTICLE INFO

Article history:

Received 15 December 2015

Received in revised form 6 April 2016

Accepted 18 April 2016

Available online 21 April 2016

ABSTRACT

The cestode parasite *Echinococcus multilocularis* is the aetiological agent of alveolar echinococcosis, responsible for considerable human morbidity and mortality. This disease is a worldwide zoonosis of major public health concern and is considered a neglected disease by the World Health Organization. The complete genome of *E. multilocularis* has been recently sequenced and assembled in a collaborative effort between the Wellcome Trust Sanger Institute and our group, with the main aim of analyzing protein-coding genes. These analyses suggested that approximately 10% of *E. multilocularis* genome is composed of protein-coding regions. This shows there is still a vast proportion of the genome that needs to be explored, including non-coding RNAs such as small RNAs (sRNAs). Within this class of small regulatory RNAs, microRNAs (miRNAs) can be found, which have been identified in many different organisms ranging from viruses to higher eukaryotes. MiRNAs are a key regulation mechanism of gene expression at post-transcriptional level and play important roles in biological processes such as development, proliferation, cell differentiation and metabolism in animals and plants. In spite of this, identification of miRNAs directly from genome-wide data only is still a very challenging task. There are many miRNAs that remain unidentified due to the lack of either sequence information of particular phylums or appropriate algorithms to identify novel miRNAs. The motivation for this work is the discovery of new miRNAs in *E. multilocularis* based on non-target genomic data only, in order to obtain useful information from the currently available unexplored data. In this work, we present the discovery of new pre-miRNAs in the *E. multilocularis* genome through a novel approach based on machine learning. We have extracted the most commonly used structural features from the folded sequences of the parasite genome: triplets, minimum free energy and sequence length. These features have been used to train a novel deep architecture of self-organizing maps (SOMs). This model can be trained with a high class imbalance and without the artificial definition of a negative class. We discovered 886 pre-miRNA candidates within the *E. multilocularis* genome-wide data. After that, experimental validation by small RNA-seq analysis clearly showed 23 pre-miRNA candidates with a pattern compatible with miRNA biogenesis, indicating them as high confidence miRNAs. We discovered new pre-miRNA candidates in *E. multilocularis* using non-target genomic data only. Predictions were meaningful using only sequence data, with no need of RNA-seq data or target analysis for prediction. Furthermore, the methodology employed can be easily adapted and applied on any draft genomes, which are actually the most interesting ones since most non-model organisms have this kind of status and carry real biological and sanitary relevance.

Availability

Web demo: <http://fich.unl.edu.ar/sinc/web-demo/mirna-som/>

Source code: <http://sourceforge.net/projects/sourcesinc/files/mirasom/>

© 2016 Elsevier Inc. All rights reserved.

1. Introduction

1.1. MicroRNAs in *Echinococcus* spp

Echinococcus multilocularis is a parasitic flatworm that causes human alveolar echinococcosis worldwide. It is among the world's most

dangerous zoonoses, developing tumor-like flatworm larvae growing in the body [37]. The metacestode of this parasite can grow in an aggressive manner budding exogenously, infiltrating and colonizing surrounding and distant tissues due to the metastatic nature of its germinative cells. The genome of *E. multilocularis* was recently sequenced and assembled in a collaborative effort between the Wellcome Trust Sanger Institute and our group [38]. Gene content analysis revealed that approximately 10% of the genome is protein-coding regions [6]. This shows that there is still a vast proportion of the genome that needs to be explored, including non-coding RNAs such as small RNAs (sRNAs).

* Corresponding author.

Within this class of small regulatory RNAs, microRNAs (miRNAs) have been identified in many different organisms. MiRNAs are endogenous ~22 nucleotide noncoding RNAs, which act as post-transcriptional regulators involved in the control of nearly all cellular pathways, from development to diseases in animals and plants [2]. MiRNAs act mainly silencing gene expression by binding to complementary sequences in the 3' untranslated regions (UTRs) of their target mRNAs. Animal miRNAs are processed in the nucleus from long primary RNA transcripts (pri-miRNAs) into ~70 nt long stem loop intermediates, known as miRNA precursors (pre-miRNAs), from which mature miRNAs are processed in the cytoplasm [4]. Pre-miRNAs (also known as hairpins) generated during biogenesis have well-known RNA secondary structures derived from primary structures that have allowed the development of computational algorithms for their identification. In a previous report, we experimentally found that miRNAs are expressed in *Echinococcus granulosus sensu lato* [5], a species closely related to *E. multilocularis*, suggesting that these small RNAs could be an essential mechanism of gene regulation in this genus. Profiling of miRNAs can be defined as the assessment of miRNA expression in a given cell type and condition [32]. Several methods are available to do this, and are preferentially used depending on a wide range of factors. The most important considerations tend to be related to the amount of biological material available, the experimental design and the final objectives of the study. As with model organisms, this kind of experiments is time-consuming and depends on the expression level of each biological stage. With the advent of new sequencing technologies, it is faster and easier to obtain genomic sequences from new organisms. However, only a few bioinformatics efforts are available to analyze this type of data, which, on the other hand, provide limited capabilities and low prediction performance for non-model organisms. To the best of our knowledge, no miRNA discovery studies from *E. multilocularis* genome wide data have been carried out to date. Thus, knowledge of the *E. multilocularis* miRNA repertoire needs to be explored.

1.2. Tools for miRNA identification

MiRNAs can be identified either by bioinformatics approaches or by sequencing strategies, both of which need computational tools for the analysis of the sequences obtained [34]. Some of the oldest strategies for miRNAs discovery include RNA conformation based approaches using Mfold [47] and RNAfold [15,16,19] as core algorithms. Other approaches are based on homology methods using known miRNA and pre-miRNA sequences from several well-known model organisms. One potential drawback of these homology-based methods is their inability to identify completely novel miRNA sequences in non-model genomes, precisely due to the conservation criteria between related genomes on which they rely and that might not be true or known for brand-new recently sequenced genomes. More recently, machine-learning techniques for miRNA prediction have been proposed, based on properties and features of well-known miRNAs. Among them, mainly supervised machine-learning techniques have been employed, using sequence composition and structural conformation features to train a learning system capable of identifying miRNA candidates [36,49]. As opposed to homology based methods, this approach could be useful for species-specific miRNA discovery since it does not depend on evolutionary conservation. As mentioned above, many methods have been developed to predict pre-miRNA loci based on the genome sequence and structural properties of the candidate loci. The miRNA classifier methods use different features to evaluate, for example, the structural stability or sequence properties of the candidates, in order to produce a final prediction [18,22,23]. However, this is a non-trivial problem when addressing it in a purely computerized way, in particular with classical supervised learning because the artificial definition of a negative class is required [10]. Methods that use only positive samples to predict new miRNAs have been described [45]. However, it is well known that these methods fail when the negative class is complex

because this region of the feature space is not properly modeled. Actually, they do not model the negative class at all or they model it under very simplified assumptions. Furthermore, when the negative class is not artificially defined and genome-wide data wants to be used, a huge imbalance is often present between the positive class (a few known miRNAs) and the unlabeled data (hundreds of thousands of sequences). Since *E. multilocularis* genome was recently generated, mining this new genomic data will provide a deeper understanding of parasite miRNome. In this work, we identify candidate novel miRNA precursors in *E. multilocularis* through a novel approach based on self-organizing maps (SOM) [21,28].

2. Materials and methods

2.1. Biologically relevant data set and hairpin features extraction

The main pipeline used for the analysis of the genome-wide data is presented in Fig. 1. The complete *E. multilocularis* genome [38] was processed by Einverted software (EMBOSS package) as described by de Souza Gomes et al. [7]) with the following parameters: gap penalty 6, minimum score threshold 25, match score 3, mismatch score -3, maximum separation between the start and end of the inverted repeat 95. Then, the inverted repeats were folded into 491,532 sequences by RNAfold (Supp. file 1). The obtained sequences were then pre-processed. Sequences with minimum free energy (MFE) threshold of -20 and single-loop folded sequences were selected according to the miRNA biogenesis model [4]. The retained sequences were analyzed using BLAST algorithm [1] against an in-house database of CDS, tRNAs, rRNAs and long non coding RNAs flatworm sequences [6]. After this, 77,429 sequences were retained. Then, all *E. multilocularis* hairpin sequences were downloaded from miRBase v21, BLAST searches among the 77,429 sequences retained were performed and a total of 18 sequences were labeled as positive class. To represent the sequences, the 34 most commonly used features were extracted. We used the smallest and less costly to compute subset of features that are extensively used nowadays to identify novel pre-miRNAs: 32 triplets [42], sequence length and MFE [18]. These features were extracted with the web tool miRNAfe [43] recently developed by us. Then, the features extracted from 77,429 sequences were used to train the SOM classifier, which identified 886 sequences as the best pre-miRNA candidates.

2.2. Classifier

In this work, instead of training a classifier in a classical supervised manner, we identified miRNA precursors with a novel approach based on several nested SOMs. For SOM training, there is no need to define the negative miRNA class. Only some examples of positive class examples (well-known pre-miRNAs) are needed to identify the neurons that have the best miRNA candidates associated to them. In this context, each neuron in the SOM is a cluster of sequences. The SOM classifier is actually composed of several nested SOMs, which are hierarchically related. This deep architecture is shown at the top of Fig. 2, where a 10-layered ($h = 10$) example is provided. The training process of the hierarchical maps starts with the root SOM on the first layer (left), with the 77,429 sequences as input. This map undergoes standard training. After that, all the sequences grouped together in a neuron (cluster) having also well-known pre-miRNAs (painted in dark blue) are labeled as highly likely pre-miRNA candidates. These sequences are chosen as input to train the map in the following layer (indicated with black lines). This process is repeated several times, further refining the classifier level after level. With this approach, each internal map is trained with only a portion of the input data: the data mapped in the pre-miRNA clusters in the previous layer. At the bottom of Fig. 2, the number of candidates is shown for each level of the SOM. It can be clearly seen here that this method significantly reduces the number of possible pre-miRNA candidates, level after level, retaining at last the high-confidence pre-miRNAs.

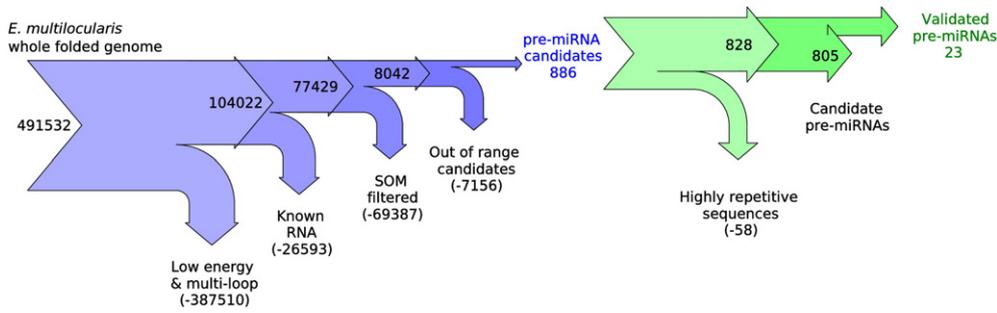


Fig. 1. Flow diagram of the pipeline proposed for miRNA discovery from *Echinococcus multilocularis* genome-wide data. The folded *E. multilocularis* genome (491,532 sequences) is used as input. Blue arrows indicate pre-processing and SOM analysis. Green arrows indicate pre-miRNA validation after RNA-seq data integration.

After four consecutive levels without changes in the number of data clustered into pre-miRNA neurons (8042 sequences), no more levels are added. These and the following levels are exactly the same since the map is trained with exactly the same data. Therefore, adding more levels does not cause over-training either. In the last level, each well-known pre-miRNA in the miRNA neurons (in blue) is grouped together with unlabeled sequences. Among them, the best bona fide candidates are selected (886) as those having feature values within ranges automatically defined by rules obtained according to the positive class (well-known miRNAs). This reduction was possible because each feature was evaluated individually with respect to its discriminative power for separating the positive class (well-known miRNAs) from the rest of the sequences. This was done iteratively, until all features were analyzed and all positive sequences were correctly classified. This way, several rules for the feature ranges were extracted, which were applied to the 8042 sequences in order to further reduce its number to 886.

2.3. Mature miRNA sequence extraction

The total number of candidate pre-miRNAs discovered by SOM analysis (886) was mapped to the complete *E. multilocularis* genome and sequences with more than 10 hits were removed (highly repetitive sequences, Fig. 1). Then, in order to extract mature miRNA sequences from pre-miRNAs retained in the previous step, 26.9 million clean mapped reads from small RNA-seq data of *E. multilocularis* metacystode stage retrieved from Cucher et al. [6]) were BLAST searched against the

pre-miRNAs sequences. BLAST algorithm was optimized for small sequences with word size set in 7, the filter for low complexity regions off, and an e-value set in 10. For each pre-miRNA with small RNAseq evidence in the stem region of the candidate pre-miRNA, the consensus mature sequence was extracted from alignments showing 100% of identity and 100% of coverage. This data was used for mature miRNA sequence determination and not for miRNA expression quantification. In order to extract additional mature miRNA sequences, all metazoan mature miRNA sequences from miRBase 21 and *Echinococcus* mature miRNAs reported in the literature that were not integrated in miRBase [3,25] were analyzed by BLAST and SSEARCH algorithms against candidate pre-miRNAs. Finally, for conservation analysis, all *E. multilocularis* mature sequences identified in previous steps were BLAST searched against related flatworm genomes: *E. granulosus*, *Echinococcus canadensis*, *Hymenolepis microstoma* and *Taenia solium*. The genomes were downloaded from <http://parasite.wormbase.org/index.html> and processed as previously described for *E. multilocularis* whole genome.

2.4. Further evaluation of the approach in a model organism

In order to further evaluate the proposal, a model organism has been used. *Caenorhabditis elegans* genome was processed in a similar way as previously described for *E. multilocularis*. The 1,739,460 sequences obtained were BLAST matched against miRBase v17 for pre-miRNA identification. A total of 200 well-known miRNAs of *C. elegans* included into miRBase v17 were labeled as positive class. All genome data (including the identified positive class) were used to train SOM until the level

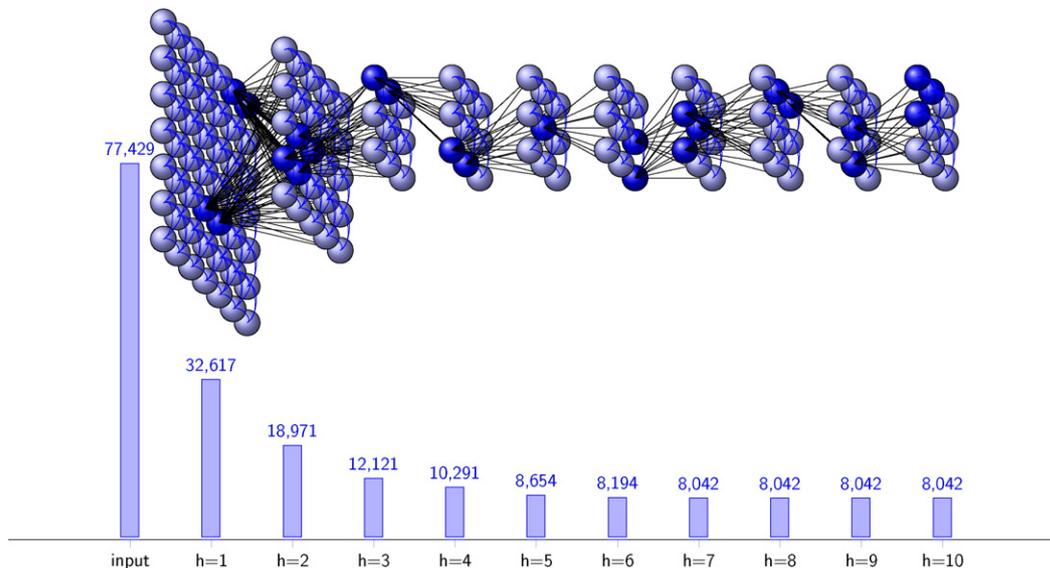


Fig. 2. Architecture developed to find pre-miRNA candidates in *E. multilocularis* genome. Top: Hierarchy of SOM classifier for 10 levels ($h = 10$). Dark blue neurons have highly likely pre-miRNA candidates, which are input to the next level SOM (black lines). Bottom: Number of pre-miRNA candidates in each level.

where the number of candidates did not change (as described previously for *E. multilocularis*). In order to evaluate the prediction performance of new miRNAs in a model organism, the miRNAs added to miRBase in its most recent version have been used as input test sequences. Therefore, the trained SOM was tested with 48 *C. elegans* pre-miRNA obtained from miRBase v19 to v21 (absent in miRBase v17). The trained model is available at <http://fich.unl.edu.ar/sinc/blog/web-demo/mirna-som-ce/>.

3. Results and discussion

In this work, we discovered 886 pre-miRNA candidates from *E. multilocularis* genome-wide data (Fig. 1). Although such quantity can be hard to validate experimentally, this must be interpreted as an important first step towards the discovery of new miRNAs in low explored genomes, such as the *E. multilocularis* one, where only few pre-miRNA sequences are available. Computationally identified miRNAs suggests that miRNA gene numbers are substantially higher than those currently known, as proposed by Piriyaopongsa et al. [31]. Most computational methods nowadays require expensive high-throughput RNA sequencing data as input [13,24]. However, we use NGS data only for validation after finding the pre-miRNA candidates, as in [35]. The few methods that have been proposed to identify miRNAs from a complete genome without such data obtain a very high number of initial candidates, hundreds of thousands or tens of thousands of sequences [26]. After that, a reduced list of the best candidates is obtained by manually applying ad hoc rules [27] in order to achieve a number of sequences that can be experimentally validated. However, for miRNA prediction most of the published approaches do not really deal with genome-wide data, instead they used data having positive and negative classes previously defined [8,11,12,14,17,20,33,41,42]. In these works, in order to train classifiers, and measure sensitivity and specificity in a cross-validation scheme, a reduced subset of negative examples must be artificially defined. Moreover, these unrealistic tests are performed over the genomes of model organisms, such as mammals or round worms, being only useful to

precisely measure the performance in cross-validation experiments, but they cannot be applied in real practical scenarios. In the proposed processing pipeline, only obvious non-miRNA sequences are filtered (according to loops, energy threshold and identity to known RNAs other than miRNAs). The remaining sequences from the original genome are all presented to the SOM for training and classification. The first advantage here is that the SOM does not require the artificial definition of negative class, thus it does not perform unrealistic tests. The second advantage is that it works directly on complete genome-wide data, which is being refined level after level, automatically discarding low-quality candidates. With this methodology, artificial examples to represent the negative class (which is actually unknown) must not be defined. The negative examples can be actually very hard to define, even for a model genome [39]. Thus, SOM is well suited to the analysis of genome data from novel non model organisms.

In order to classify each miRNA as conserved or novel, we analyzed the identity of all pre-miRNA candidates discovered by SOM with already reported metazoan miRNAs (miRBase v21) and *E. multilocularis* miRNAs [6]. This analysis allowed us to identify 13 pre-miRNAs previously described (Supplementary Table S1). Taking into account the 18 miRNAs used as positive class, the total of miRNAs found was 31 out of 37 miRNAs expected to be in *E. multilocularis* [6]. Since four miRNAs were absent in the genome input dataset because their folded structure did not match the filter criteria employed, the sensitivity of SOM reached 94% (31/33). Moreover, 10 new pre-miRNAs were also identified totaling 23 pre-miRNAs. The mature miRNA annotation, their clean mapped read counts and the biological function in other organisms are shown in Table 1. *E. multilocularis* RNA-seq clearly mapped to the hairpin stem region with a pattern compatible with miRNA biogenesis indicating them as high-confidence miRNAs. As an example, a schematic representation of the secondary structure from the conserved *E. multilocularis* premiRNA 36b is shown in Fig. 3.

These new pre-miRNAs represent, in the first place, flatworm-specific miRNAs since they were not detected in any other phyla. Also,

Table 1
Conserved and novel *Echinococcus multilocularis* microRNAs predicted from whole genome data.

MiRNA ID	Read counts ^a	Biological function ^{bc}	Reference ^b
emu-bantam-3p	1,184,581	Regulates the growth of dendrites in sensory neurons of <i>Drosophila melanogaster</i> epithelial cells. Present only in protostomes	Parrish et al. [30]
emu-miR-31-5p	88	Tumour suppressor in humans	O'Day and Lal [29]
emu-miR-36a-3p	617	Unknown, present only in protostomes	Macchiaroli et al. [25]
emu-miR-36b-3p	1075	Unknown, present only in protostomes	Cucher et al. [6]
emu-miR-61-3p	578,860	Promotes development in <i>Caenorhabditis elegans</i> . Present only in protostomes	Yoo and Greenwald [44]
emu-miR-281-3p	17,958	Enhance viral replication in <i>Aedes albopictus</i>	Zhou et al. [46]
emu-miR-307-3p	123,277	Unknown, present only in protostomes	Cucher et al. [6]
emu-miR-1992-3p	24	Unknown, present only in protostomes	Cucher et al. [6]
emu-miR-2162-3p	100,642	Unknown, present only in protostomes	Cucher et al. [6]
emu-miR-10,293-3p	4017	Unknown	Cucher et al. [6]
emu-miR-3479a-3p	56,603	Unknown	Cucher et al. [6]
emu-miR-3479b-3p	63,552	Unknown	Cucher et al. [6]
emu-miR-7b-5p	1070	Controls epidermal growth factor receptor signaling and promotes photoreceptor cell differentiation in <i>Drosophila</i>	Jiang et al. [48], Macchiaroli et al. [25] (egr-miR-7b-5p)
emu-miR-new1-5p	8	Unknown	This work and Bai et al. [3] (egr-new-48)
emu_miR-new2-3p	32	Unknown	This work
emu_miR-new3-5p	123	Unknown	This work
emu_miR-new4-5p	58	Unknown	This work and Bai et al. [3] (egr-new-12)
emu_miR-new5-3p	1	Unknown	This work and Bai et al. [3] (egr-new-25)
emu_miR-new6-5p	1	Unknown	This work and Bai et al. [3] (egr-new-114)
emu-miR-new7-5p	41	Unknown	This work and Bai et al. [3] (egr-new-7)
emu-miR-new8-3p	20	Unknown	This work and Bai et al. [3] (egr-new-24)
emu-miR-new9-3p	246	Unknown	This work
emu-miR-new10-5p	231	Unknown	This work and Bai et al. [3] (egr-new-29)
Total	2,133,125		

^a Number of clean mapped reads without normalization.

^b Described in model species.

^c Most relevant references for miRNA function in other organisms or studies on related *Echinococcus* species.

Mature emu-miR-36b-3p
5'- UCACCGGGUAGUUAUACGCCU-'3

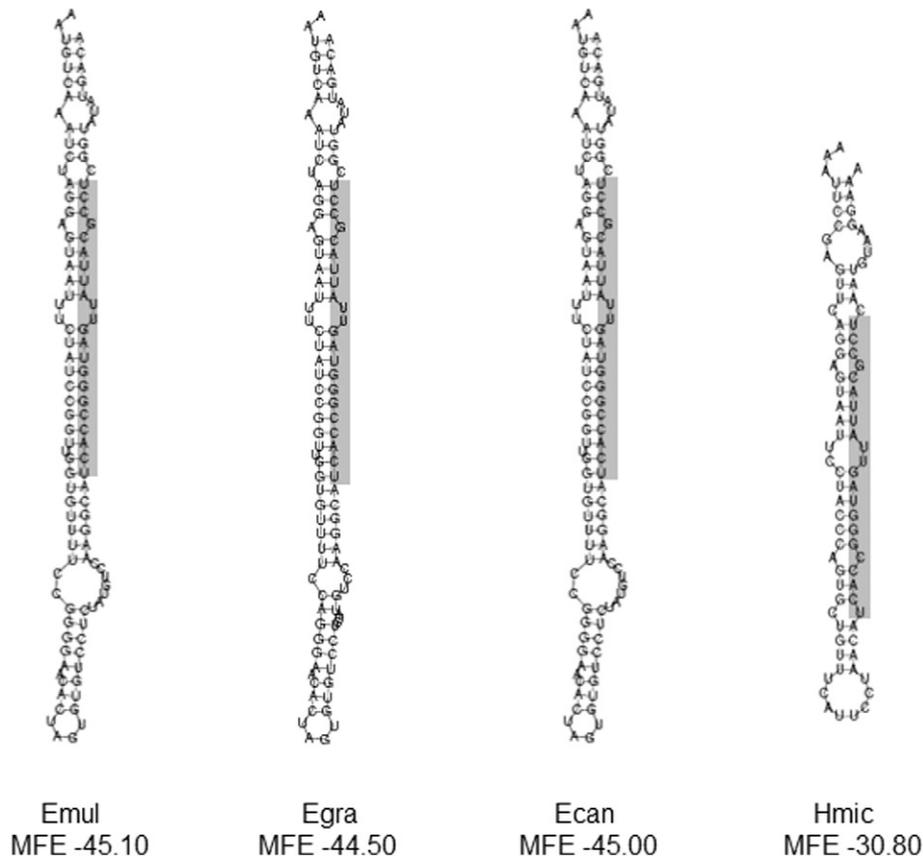


Fig. 3. Schematic representation of the secondary structure from the conserved pre-miRNA 36b discovered by the SOM. The secondary structure predictions for pre-miRNA-36b is shown for four species of flatworms. Emul: *E. multilocularis*; Egra: *E. granulosus*; Ecan: *E. canadensis*; Hmic: *H. microstoma*. Mature miRNA sequences are underlined. Minimum free energy (MFE) is expressed as kcal/mol.

some of them were recently reported in *E. granulosus* [3]. It can be noticed here the ability of the SOM to discover of new miRNAs, only with genomic data as input. Furthermore, the secondary structure from all new pre-miRNAs discovered by SOM analysis is shown in Fig. 4. Structural features such as MFE and mature miRNA sequences that mapped to them clearly showed that they were bona fide pre-miRNAs. All mature and pre-miRNA sequences and structures are available in Supplementary Table S1 and Fig. S1. Additionally, our method discovered miRNAs in *E. multilocularis* that were not identified by a recent bioinformatics approach [50] such as miR-36, miR-307, miR-1992, miR-3479, highlighting the potential of SOM analysis for miRNA discovery. Interestingly, this miRNAs were considered lost in *Echinococcus* [9] but SOM discovered them in coincidence with previously reports [6,25].

We have also searched for these 23 pre-miRNA sequences in closely related flatworm genomes. All of them were found in at least one of the four related flatworm species (Fig. 3, Supplementary Table S1). Several of the mature miRNAs found in this work are deeply conserved among bilateria such as emu-miR-281 and emu-miR-31, but others are found only in protostomia such as emu-bantam, emu-miR-36 and emu-miR-1992. So far, there is no information about the biological function of these miRNAs in *Echinococcus*. These results could be interpreted as a good indicator of the biological confidence of the predictions obtained with the pipeline proposed in this work, and indicate that the SOM could discover both conserved and novel miRNAs from *E. multilocularis* genome data. Although losses of conserved miRNAs have been

previously proposed in parasite flatworms [9,25], the presence of specific miRNAs is expected since novel miRNAs have been recently reported from small RNAseq data in other helminth parasites [3,40]. The new pre-miRNA sequences discovered in our work are good candidates to be flatworm-specific miRNAs since they have no identity with miRNAs from other phyla. These miRNA sequences are the most interesting ones because they could have a crucial role in the establishment and/or progression of human alveolar echinococcosis. As future work, it could be interesting to be able to determine the *E. multilocularis* life cycle stage where the new miRNAs discovered in this work are expressed which could be done following approaches previously published by us [25]. The knowledge of the complete repertoire of miRNAs, conserved and specific ones, is key to understand the development of the parasite and the progression and control of this neglected disease.

The validation of the proposed methodology in a non-model organism has proved its effectiveness. However, benchmarking it in a well-known reference genome can provide evidence of its utility in a wide number of organisms. Thus, we have performed a benchmarking test of the proposed SOM approach with a well-known reference genome. The SOM was trained with the complete genome data plus a total of 200 *C. elegans* well-known pre-miRNA sequences present in miRBase v17. Then, the trained SOM has been tested with 48 pre-miRNAs more recently added to miRBase v18–21 and absent in v17. In this test, 44 out of 48 pre-miRNA have been identified as positive class, resulting in a SOM sensitivity of 92%.

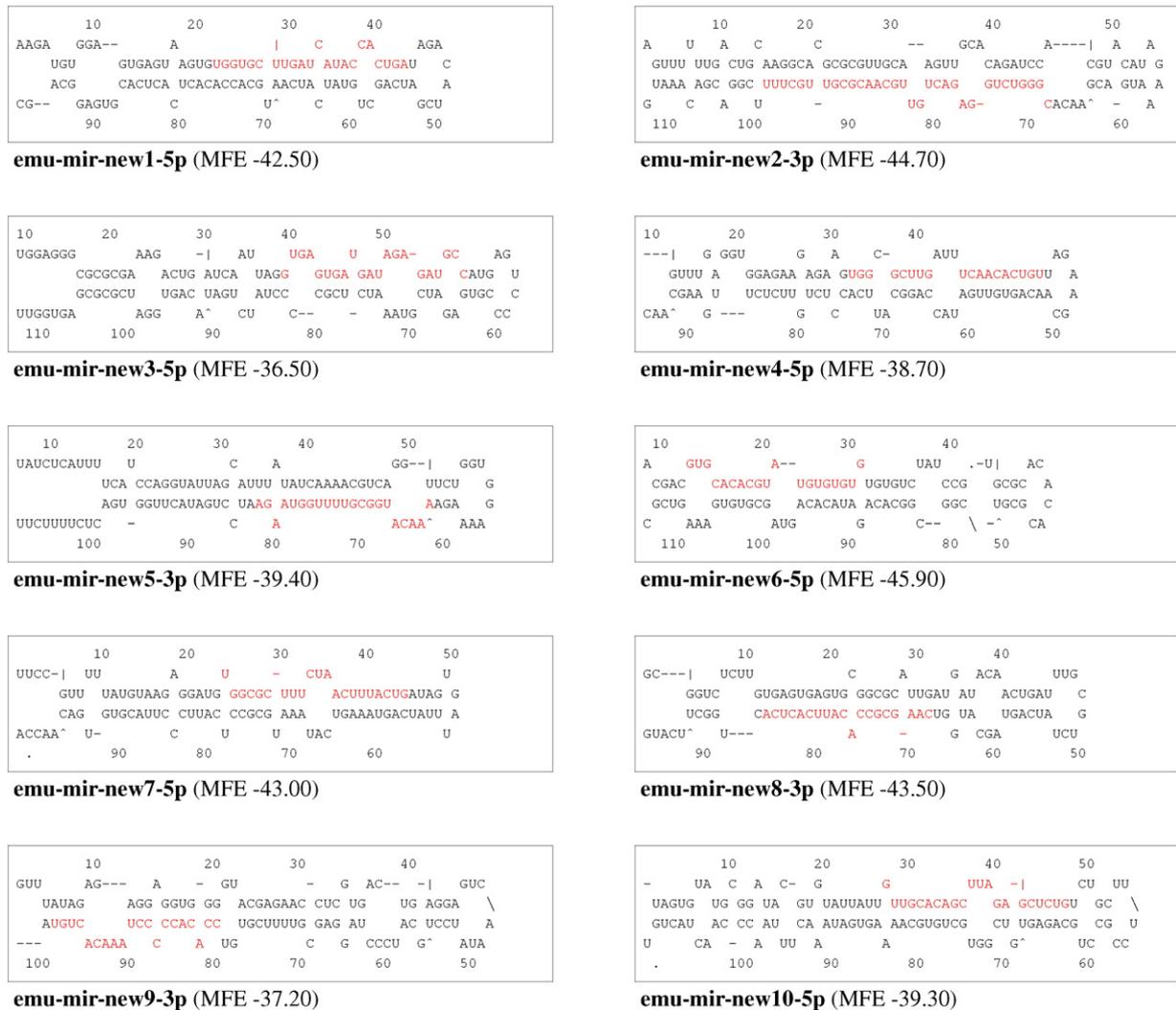


Fig. 4. The secondary structure predictions of all new miRNAs from *E. multilocularis* discovered by SOM analysis. Mature miRNAs are indicated in red. Minimum free energy (MFE) is expressed as kcal/mol.

4. Conclusions

We applied SOM analysis for *E. multilocularis* miRNA prediction and demonstrated its effectiveness and usefulness. Although using purely computational methods for de novo miRNA prediction was a real challenge and a difficult problem to address, this analysis allow us to discover good candidates from *E. multilocularis* genome sequencing data. Most pre-miRNA prediction methods based on supervised machine learning methods, which need to artificially define the negative class, cannot handle the class imbalance existing in such genome-wide data. However, the proposed method addressed the problem effectively without requiring the artificial definition of a negative class dataset. With this approach, complete genomes containing thousands of hairpins sequences could be analyzed and only highly likely hairpin sequences can be further selected for biological validation. We found novel *E. multilocularis* pre-miRNAs from non target genomic data without the need of RNA-seq data and all of them conserved in at least one related flatworm species. These results clearly indicate that there are still several genomic sequences to be classified and ready to be analyzed deeply. We found expression of mature miRNAs derived from pre-miRNA candidates adding confidence to the predictions obtained by SOM analysis. The data obtained in this work will be useful to search for new mature miRNAs expressed in the human parasite *E. multilocularis* resulting in new tools for the diagnosis, prevention

and developmental regulation of alveolar echinococcosis neglected disease.

Supplementary data to this article can be found online at <http://dx.doi.org/10.1016/j.ygeno.2016.04.002>.

Authors' contributions

LK, GS and DHM wrote the manuscript and designed the experiments. GS and DHM designed and implemented the SOM deep architecture and training scripts. CY developed the scripts for feature extraction and data pre-processing. LK, NM and LM analyzed data from high-throughput experiments. All authors read and approved the manuscript.

Acknowledgements

This work was supported by Agencia Nacional de Promoción Científica y Tecnológica (ANPCyT), Argentina, project PICT-CABBIO 2012 No 3044 and PICT 2014 No 2627, Universidad Nacional del Litoral (UNL) CAI + D 2011 548, and by Consejo Nacional de Investigaciones Científicas y Tecnológicas (CONICET) project PIP 114 2011 and PIP 117 2013. High-throughput analysis was performed in a local server (ID 924) at Instituto de Investigaciones en Microbiología y Parasitología Médicas (IMPAM) which is part of Sistema Nacional de Computación de Alto Desempeño (SNCAD) of Ministerio de Ciencia, Tecnología e

Innovación Productiva (MINCYT). Thanks to Dr. Marcela Cucher for making raw data of *E. multilocularis* available.

References

- [1] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, D.J. Lipman, Basic local alignment search tool, *J. Mol. Biol.* 215 (1990) 403–410.
- [2] S. Ameres, P. Zamore, Diversifying microRNA sequence and function, *Nat. Rev. Mol. Cell Biol.* 14 (8) (2013) 475–488.
- [3] Y. Bai, Z. Zhang, L. Jin, H. Kang, Y. Zhu, L. Zhang, L. Xia, F. Ma, L. Zhao, B. Shi, J. Li, D. McManus, W. Zhang, S. Wang, Genome-wide sequencing of small RNAs reveals a tissue-specific loss of conserved microRNA families in *Echinococcus granulosus*, *BMC Genomics* 1 (15) (2014) 736.
- [4] D. Bartel, MicroRNAs: genomics, biogenesis, mechanism, and function, *Cell* 116 (2004) 281–297.
- [5] M. Cucher, L. Prada, G. Mourglia-Ettlin, S. Dematteis, F. Camicia, S. Asurmendi, M. Rosenzvit, Identification of *Echinococcus granulosus* microRNAs and their expression in different life cycle stages and parasite genotypes, *Int. J. Parasitol.* 41 (3–4) (2011) 439–448.
- [6] M. Cucher, N. Macchiaroli, L. Kamenetzky, L. Maldonado, K. Brehm, M.C. Rosenzvit, High-throughput characterization of *Echinococcus* spp. metacestode miRNomes, *Int. J. Parasitol.* 45 (4) (2015) 253–267.
- [7] M. de Souza Gomes, M.K. Muniyappa, S.G. Carvalho, R. Guerra-S, C. Spillane, Genome-wide identification of novel microRNAs and their target genes in the human parasite *Schistosoma mansoni*, *Genomics* 98 (2) (2011) 96–111.
- [8] J. Ding, S. Zhou, J. Guan, MiRenSVM: towards better prediction of microRNA precursors using an ensemble SVM classifier with multi-loop features, *BMC Bioinform.* 11 (11) (2010) S11.
- [9] B. Fromm, M. Worren, C. Hahn, E. Hovig, L. Bachmann, Substantial loss of conserved and gain of novel MicroRNA families in flatworms, *Mol. Biol. Evol.* 30 (12) (2013) 2619–2628.
- [10] C.P.C. Gomes, J.-H. Cho, L. Hood, O.L. Franco, R.W. Pereira, K.A. Wang, Review of computational tools in microRNA discovery, *Front. Genet.* 4 (2013) 81.
- [11] K. Gkirtzou, I. Tsamardinos, P. Tsakalides, P. Poirazi, MatureBayes: a probabilistic algorithm for identifying the mature miRNA within novel precursors, *PLoS One* 5 (8) (2010) e11843.
- [12] A. Gudy, M. Szczeniak, M. Sikora, I. Makalowska, HuntMi: an efficient and taxon-specific approach in pre-miRNA identification, *BMC Bioinform.* 14 (1) (2013) 83+.
- [13] M. Hackenberg, N. Rodriguez-Ezpeleta, A. Aransay, miRanalyzer: an update on the detection and analysis of microRNAs in high-throughput sequencing experiments, *Nucleic Acids Res.* 39 (Suppl. 2) (2011) W132–W138.
- [14] J. Hertel, P.F. Stadler, Hairpins in a Haystack: recognizing microRNA precursors in comparative genomics data, *Bioinformatics* 22 (14) (2006) e197–e202.
- [15] I. Hofacker, W. Fontana, P. Stadler, S. Bonhoeffer, M. Tacker, P. Schuster, Fast folding and comparison of RNA secondary structures, *Monatsh. Chem./Chem. Mon.* 125 (1994) 167–188.
- [16] I.L. Hofacker, The Vienna RNA secondary structure server, *Nucleic Acids Res.* 31 (2003) 3429–3431.
- [17] T.H. Huang, B. Fan, M. Rothschild, Z.L. Hu, K. Li, S.H. Zhao, MiRFinder: an improved approach and software implementation for genome-wide fast microRNA precursor scans, *BMC Bioinform.* 8 (1) (2007) 341.
- [18] de O.N. Lopes, A. Schliep, A. de Carvalho, The discriminant power of RNA features for pre-miRNA recognition, *BMC Bioinform.* 15 (1) (2014) 124+.
- [19] A. Jacobson, M. Zuker, Structural analysis by energy dot plot of a large mRNA, *J. Mol. Biol.* 233 (2) (1993) 261–269.
- [20] P. Jiang, H. Wu, W. Wang, W. Ma, X. Sun, Z. Lu, MiPred: classification of real and pseudo microRNA precursors using random forest prediction model with combined features, *Nucleic Acids Res.* 35 (1) (2007) W339–W344.
- [21] T. Kohonen, M.R. Schroeder, T.S. Huang, *Self-Organizing Maps*, Springer-Verlag New York, Inc., 2005.
- [22] L. Li, J. Xu, D. Yang, X. Tan, H. Wang, Computational approaches for microRNA studies: a review, *Mamm. Genome* 21 (1) (2010) 1–12.
- [23] B. Liu, J. Li, M. Cairns, Identifying mirnas, targets and functions, *Brief. Bioinform.* 15 (1) (2014) 1–19.
- [24] M. Friedlander, S. Mackowiak, N. Li, W. Chen, N. Rajewsky, miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades, *Nucleic Acids Res.* 40 (1) (2012) 37–52.
- [25] N. Macchiaroli, M. Cucher, M. Zarowiecki, L. Maldonado, L. Kamenetzky, M.C. Rosenzvit, microRNA profiling in the zoonotic parasite *Echinococcus canadensis* using a high-throughput approach, *Parasites Vectors* 8 (1) (2015) 83.
- [26] N. Mendes, A. Freitas, A. Vasconcelos, M.-F. Sagot, Combination of measures distinguishes pre-mirnas from other stem-loops in the genome of the newly sequenced *Anopheles darlingi*, *BMC Genomics* 11 (1) (2010) 529.
- [27] N.D. Mendes, S. Heyne, A.T. Freitas, M.-F. Sagot, R. Backofen, Navigating the unexplored seascape of pre-miRNA candidates in single-genome approaches, *Bioinformatics* 28 (23) (2012) 3034–3041.
- [28] D. Milone, G. Stegmayer, L. Kamenetzky, M. López, J. Lee, J. Giovannoni, F. Carrari, omeSOM: a software for clustering and visualization of transcriptional and metabolite data mined from interspecific crosses of crop plants, *BMC Bioinform.* 11 (2010) 438–447.
- [29] E. O'Day, A. Lal, MicroRNAs and their target gene networks in breast cancer, *Breast Cancer Res.* 12 (2) (2010) 201.
- [30] J. Parrish, P. Xu, C. Kim, L. Jan, Y. Jan, The microRNA bantam functions in epithelial cells to regulate scaling growth of dendrite arbors in drosophila sensory neurons, *Neuron* 63 (6) (2009) 788–802.
- [31] J. Priyapongsa, L. Ramírez, K. Jordan, Origin and evolution of human microRNAs from transposable elements, *Genetics* 176 (2) (2007).
- [32] C. Pritchard, H. Cheng, M. Tewari, MicroRNA profiling: approaches and considerations, *Nat. Rev. Genet.* 13 (5) (2012) 358–369.
- [33] M.E. Rahman, R. Islam, S. Islam, S.I. Mondal, M.R. Amin, MiRANN: a reliable approach for improved classification of precursor microRNA using Artificial Neural Network model, *Genomics* 99 (4) (2012) 189–194.
- [34] Rosenzvit, M., Cucher, M., Kamenetzky, L., Macchiaroli, N., Prada, L., and Camicia, F. (2013). *MicroRNAs in Endoparasites*. Nova Science Publishers. Book, *MicroRNA and Non-Coding RNA: Technology, Developments and Applications*. Series: Genetics - Research and Issues. Pages: 7x10 - (NBC-R) Editors: James C. Johnson Editorial: Nova Science Publishers, Inc. Nueva York, USA. 978-1-62618-443-5.
- [35] M.D. Saçar, C. Bağcı, J. Allmer, Computational prediction of microRNAs from *Toxoplasma gondii* potentially regulating the hosts' gene expression, *Genomics Proteomics Bioinformatics* 12 (5) (2014) 228–238.
- [36] P. Sætrom, O. Snøve, Robust machine learning algorithms predict MicroRNA genes and targets, *Methods Enzymol.* 427C (2007) 25–49.
- [37] P. Torgerson, K. Keller, M. Magnotta, N. Ragland, The global burden of alveolar echinococcosis, *PLoS Negl. Trop. Dis.* 4 (6) (2010) e722.
- [38] I. Tsai, M. Zarowiecki, N. Holroyd, et al., The genomes of four tapeworm species reveal adaptations to parasitism, *Nature* 496 (2013) 57–63.
- [39] L. Wei, M. Liao, Y. Gao, R. Ji, Z. He, Q. Zou, Improved and promising identification of human MicroRNAs by incorporating a high-quality negative set, *IEEE/ACM Trans. Comput. Biol. Bioinform.* 11 (1) (2014) 192–201.
- [40] A. Winter, W. Wei, M. Hunt, M. Berriman, J. Gilleard, E. Devaney, C. Britton, Diversity in parasitic nematode genomes: the microRNAs of *Brugia pahangi* and *Haemonchus contortus* are largely novel, *BMC Genomics* 13 (1) (2012) 4.
- [41] Y. Xu, X. Zhou, W. Zhang, MicroRNA prediction with a novel ranking algorithm based on random walks, *Bioinformatics* 24 (1) (2008) i50–i58.
- [42] C. Xue, F. Li, T. He, G.-P. Liu, Y. Li, X. Zhang, Classification of real and pseudo microRNA precursors using local structure-sequence features and support vector machine, *BMC Bioinformatics* 6 (1) (2005) 310.
- [43] C.A. Yones, G. Stegmayer, L. Kamenetzky, D.H. Milone, miRNafe: a comprehensive tool for feature extraction in microRNA prediction, *Biosystems* 138 (2015) 1–5.
- [44] A. Yoo, I. Greenwald, LIN-12/Notch activation leads to microRNA-mediated down-regulation of Vav in *C. elegans*, *Science* 310 (5752) (2005) 1330–1333.
- [45] M. Yousef, S. Jung, L. Showe, M. Showe, Learning from positive examples when the negative class is undetermined-microRNA gene identification, *Algorithm. Mol. Biol.* 3 (2008) 2.
- [46] Y. Zhou, Y. Liu, H. Yan, Y. Li, H. Zhang, J. Xu, S. Puthiyakunnon, X. Chen, miR-281, an abundant midgut-specific miRNA of the vector mosquito *Aedes albopictus* enhances dengue virus replication, *Parasites Vectors* 1 (7) (2014) 488.
- [47] M. Zuker, Mfold web server for nucleic acid folding and hybridization prediction, *Nucleic Acids Res.* 31 (13) (2003) 3406–3415.
- [48] L. Jiang, X. Liu, Z. Chen, Y. Jin, C.E. Heidbreder, A. Kolokythas, A. Wang, Y. Dai, X. Zhou, MicroRNA-7 targets IGF1R (insulin-like growth factor 1 receptor) in tongue squamous cell carcinoma cells, *Biochem. J.* 432 (1) (2010) 199–205.
- [49] W.C. Chan, M.R. Ho, S.C. Li, K.W. Tsai, C.H. Lai, C.N. Hsu, W.C. Lin, MetaMirClust: discovery of miRNA cluster patterns using a data-mining approach, *Genomics* 100 (3) (2012) 141–148.
- [50] X. Jin, L. Lu, H. Su, Z. Lou, F. Wang, Y. Zheng, G.T. Xu, Comparative analysis of known miRNAs across platyhelminths, *FEBS J.* 280 (16) (2013) 3944–3951.