



This document is a postprint version of an article published in Genomics© Elsevier after peer review. To access the final edited and published work see <https://doi.org/10.1016/j.ygeno.2019.12.005>

Document downloaded from:



1 **Discovery and annotation of novel microRNAs in the porcine genome by**
2 **using a semi-supervised transductive learning approach**

3 Emilio Mármol-Sánchez¹, Susanna Cirera², Raquel Quintanilla³, Albert Pla⁴, Marcel
4 Amills^{1,5}

5 ¹Centre for Research in Agricultural Genomics (CRAG), CSIC-IRTA-UAB-UB,
6 Universitat Autònoma de Barcelona, 08193 Bellaterra, Spain.

7 ²Department of Veterinary and Animal Sciences, Faculty of Health and Medical
8 Sciences, University of Copenhagen, Frederiksberg, Denmark.

9 ³Animal Breeding and Genetics Program, Institute for Research and Technology in
10 Food and Agriculture (IRTA), Torre Marimon, 08140 Caldes de Montbui, Spain.

11 ⁴Department of Medical Genetics, University of Oslo and Oslo University Hospital,
12 Oslo, Norway.

13 ⁵Departament de Ciència Animal i dels Aliments, Universitat Autònoma de Barcelona,
14 08193 Bellaterra, Barcelona, Spain.

15

16 Corresponding author: Emilio Mármol-Sánchez. Centre for Research in Agricultural
17 Genomics (CRAG), CSIC-IRTA-UAB-UB, Universitat Autònoma de Barcelona, 08193
18 Bellaterra, Spain. Phone: +34 935636600. Email: emilio.marmol@cragenomica.es

19

20

21

22

23 **Highlights**

- 24 • Motif search improved pre-miRNA reconstruction from mature microRNA
25 sequences.
- 26 • Semi-supervised methods outperformed canonical supervised classification
27 algorithms.
- 28 • The presence of multiple isomiRs in the porcine muscle miRNA repertoire was
29 uncovered.
- 30 • A total of 47 novel microRNA genes were identified in the porcine genome.
- 31 • RT-qPCR analyses allowed us to confirm the existence of three novel porcine
32 microRNAs.

33

34

35 **Abstract**

36 Despite the broad variety of available microRNA (miRNA) prediction tools, their
37 application to the discovery and annotation of novel miRNA genes in domestic species is
38 still limited. In this study we designed a comprehensive pipeline (eMIRNA) for miRNA
39 identification in the yet poorly annotated porcine genome and demonstrated the
40 usefulness of implementing a motif search positional refinement strategy for the accurate
41 determination of precursor miRNA boundaries. The small RNA fraction from *gluteus*
42 *medius* skeletal muscle of 48 Duroc gilts was sequenced and used for the prediction of
43 novel miRNA loci. Additionally, we selected the human miRNA annotation for a
44 homology-based search of porcine miRNAs with orthologous genes in the human
45 genome. A total of 20 novel expressed miRNAs were identified in the porcine muscle

46 transcriptome and 27 additional novel porcine miRNAs were also detected by homology-
47 based search using the human miRNA annotation. The existence of three selected novel
48 miRNAs (ssc-miR-483, ssc-miR484 and ssc-miR-200a) was further confirmed by reverse
49 transcription quantitative real-time PCR analyses in the muscle and liver tissues of
50 Göttingen minipigs. In summary, the eMIRNA pipeline presented in the current work
51 allowed us to expand the catalogue of porcine miRNAs and showed better performance
52 than other commonly used miRNA prediction approaches. More importantly, the
53 flexibility of our pipeline makes possible its application in other yet poorly annotated
54 non-model species.

55

56 **Keywords:** MicroRNA discovery; Motif search; Porcine skeletal muscle; Semi-
57 supervised learning; Small RNA-seq.

58

59

60 **Introduction**

61 The accurate annotation of a comprehensive set of miRNAs in different species has been
62 challenging since the first genome assemblies were published, although an ever-
63 increasing amount of knowledge about miRNA diversity across species has been
64 accumulating during the past years, being available in public databases [1-3]. Despite
65 these advances, many commonly studied domestic species still lack a complete and
66 reliable set of annotated miRNAs in their genomes [1].

67 The computational prediction of miRNAs in sequenced genomes initially relied on the
68 strong conservation of mature miRNA sequences across closely related species [4,5],
69 taking advantage of homology-based comparisons between well annotated genome

70 assemblies and other poorly annotated organisms [6-8]. Other approaches focused on
71 rule-based classification, integrating other sources of information such as sequencing data
72 or structural features to identify novel miRNAs [9-12]. More recently, several Machine
73 Learning (ML) approaches have been proposed for miRNA prediction. Different tools
74 have addressed the problem of correctly classifying miRNAs by training ML algorithms
75 with a set of positive (annotated miRNAs) and negative (other non-miRNA sequences)
76 data sets. [13-16]. Nevertheless, despite the broad array of available tools for novel
77 miRNA identification, their application to the discovery and annotation of novel miRNAs
78 in domestic species is still limited [17-25]. Moreover, the majority of miRNA surveys
79 carried out in domestic species do not generally take into account several issues regarding
80 miRNA genes prediction that have recently emerged. For instance, the set of positive
81 training annotated miRNAs often include misclassified sequences [26,27], whereas the
82 negative class is sometimes not clearly defined, i.e. different types of sequences have
83 been used as negative data sets (coding regions, pseudo-hairpins, non-coding hairpins or
84 artificial randomized miRNA sequences). Despite some efforts [28], obtaining a truly
85 representative negative class is still challenging and few approaches have critically
86 addressed this important issue [29-31]. Besides, miRNAs are thought to encompass a
87 small percentage of the total non-coding transcriptomic repertoire, with thousands of
88 other non-miRNA hairpin-like RNA molecules that represent a major fraction of it. This
89 circumstance contributes to create a high class-imbalance between positive and negative
90 sequences. Different approaches have dealt with such phenomenon [32], but recent
91 studies have shown that commonly used techniques for solving the high-class imbalance
92 problem in microRNA prediction may not be suited to a real-case classification scenario
93 [15].

94 In this study we present eMIRNA, a bioinformatics pipeline for miRNA discovery and
95 annotation in sequenced genomes. The proposed pipeline implements a semi-supervised
96 transductive learning approach to predict and annotate novel microRNAs in the porcine
97 genome, overcoming several of the drawbacks outlined above. In order to validate the
98 performance of our pipeline in a real-case scenario, we have applied it to the analysis of
99 a data set comprising the small RNA fraction of *gluteus medius* skeletal muscle from a
100 population of 48 Duroc gilts [33,34]. Furthermore, making use of the better annotated *H.*
101 *sapiens* miRNAome, an additional set of novel porcine miRNA genes were identified
102 based on a homology-based search approach. Finally, some of the identified novel porcine
103 miRNA candidates were independently validated in a Göttingen minipig population,
104 investigating their expression in skeletal muscle and liver tissues.

105

106

107 **Materials and methods**

108 A detailed flow chart depicting all steps described in the eMIRNA pipeline is shown
109 in [Figure 1](#). Additional instructions and modular scripts needed for the implementation of
110 eMIRNA are available at: <https://github.com/emarmolsanchez/eMIRNA/>.

111 **Positive and negative training data sets**

112 To define the corresponding positive (annotated miRNAs) data set required for novel
113 miRNA prediction, two approaches were considered:

114 1) The annotated pre-miRNA coordinates in Sscrofa11.1 genome assembly were obtained
115 from Ensembl repositories, release version 97
116 (<http://www.ensembl.org/info/data/ftp/index.html>), and the corresponding sequences
117 were extracted from the pig reference genome by using the BEDTools suite v2.27.0

118 software [35]. miRNA loci located in scaffolds were removed from further analyses,
119 resulting in a total of 484 annotated porcine miRNA genes. Sequence repeats from pre-
120 miRNA duplicated elements were removed from the retrieved positive data set by using
121 the CD-HIT Suite [36] with a 0.9 sequence identity cut-off value (i.e. sequences showing
122 a similarity $\geq 90\%$ to each other were removed and only unique representative pre-
123 miRNA candidates were retained). Moreover, to avoid the inclusion of miss-annotated
124 miRNA loci, an additional filtering based on secondary structure folding was applied. To
125 this end, the RNAfold tool from the ViennaRNA Package 2.0 [37] was used to select
126 sequences with canonical pre-miRNA hairpin secondary structures (stem-loop
127 conformation with one single terminal loop and two stems). Sequences that failed to
128 comply with required folding structure pre-requisites were removed.

129 2) In the second approach, the curated miRNA annotation for Sscrofa11.1 available in the
130 miRCarta database [2] was retrieved, and the same pre-filtering criteria based on sequence
131 identity and secondary structure employed in the analysis of the Ensembl data set were
132 applied. The miRCarta database [2] integrates one of the most comprehensive and curated
133 databases for miRNA annotation and functional activity, aiming to overcome the
134 limitations of other widely used miRNA databases such as miRBase [1].

135 Regarding the negative data set (other hairpin-like sequences), two different data sources
136 were used. First, the annotated non-coding transcripts in Ensembl repositories were
137 retrieved and non-miRNA sequences were retained. Analogously to what was
138 implemented for the positive data set, identity by sequence and secondary structure pre-
139 filters were applied, and non-miRNA non-coding hairpin-like unique sequences were
140 obtained. Only sequences ranging from 50 up to 150 nucleotides (nt) were retained, thus
141 removing hairpin-like long non-coding RNAs from the negative data set. Additionally, a
142 set of unlabeled sequences within the porcine reference genome (Sscrofa11.1) were

143 generated by extracting candidate pre-miRNA-like sequences from random blocks of 1
144 Megabase (Mb) in each of the chromosomes of the porcine assembly with
145 the *HextractoR* package [38], and the previously described pre-filters for the negative
146 class were subsequently applied.

147

148 **Obtaining putative miRNA candidate sequences from the porcine genome**

149 In order to test our method with pig transcriptomic data, a small RNA-seq data set was
150 generated by sequencing the muscle transcriptome of 48 gilts used in two previous studies
151 [33,34]. Upon collection, muscle samples were individually submerged in RNAlater and
152 snap-frozen in liquid nitrogen. Samples were pulverized and homogenized in 1 ml of TRI
153 Reagent (Thermo Fisher Scientific, Barcelona, Spain). Total RNA was isolated with the
154 RiboPure kit (Ambion, Austin, TX). A Nanodrop ND-100 spectrophotometer (Thermo
155 Fisher Scientific, Barcelona, Spain) was used to assess RNA concentration and quality.
156 RNA integrity expressed in RNA Integrity Number (RIN) units was measured with a
157 Bionalyzer-2100 equipment (Agilent Technologies Inc., Santa Clara, CA). High quality
158 RNA samples were then submitted to Sistemas Genómicos S.L.
159 (<https://www.sistemasgenomicos.com>) for small RNA sequencing. Library preparation
160 for each individual sample was carried out with the TruSeq Small RNA Sample
161 Preparation Kit (Illumina Inc., USA) and small RNA libraries were single-end sequenced
162 (1 × 50 bp) in a HiSeq 2500 platform (Illumina Inc., CA).

163 FASTQ sequence files were subjected to a quality control check as reported by Cardoso
164 et al. [33]. After preliminary quality-based filtering, sequencing adaptors were trimmed
165 with the Cutadapt software [39] and an acceptance sequence window of 15–30 nt per read
166 was established. Processed FASTQ files from all sequenced samples (N=48) were
167 pooled and collapsed to unique FASTA sequences with the FASTQ collapser tool from

168 FASTX-toolkit (http://hannonlab.cshl.edu/fastx_toolkit/). Unique FASTA sequences
169 represented by >10 reads-per-million (RPM) were considered to be significantly
170 expressed above the background noise [40], and thus selected for further analyses (File
171 S1). The CD-HIT Suite [36] was employed to build sequence clusters with >0.9 sequence
172 identity.

173 Furthermore, the human mature miRNA coordinates were obtained from Ensembl
174 repositories and the corresponding sequences were retrieved from the GRCh38.p12
175 assembly. Pre-filtering based on sequence identity was applied and a set of non-redundant
176 human mature miRNAs was generated for homology-based search in the Sscrofa11.1
177 porcine assembly (File S2).

178

179 **Pre-miRNA reconstruction by sequence elongation and motif search**

180 Once putative mature miRNA candidate sequences from the small RNA-seq data set and
181 the human mature miRNA sequences were retrieved, they were aligned against the
182 porcine reference assembly (Sscrofa11.1) with the Bowtie aligner [41] and the following
183 specifications for short reads: 1) allowing 2 mismatches within the entire aligned
184 sequence with respect to the reference assembly, 2) removing reads with >50 putative
185 mapping sites and 3) reporting first single best stratum alignment (*bowtie -n 2 -l 25 -m 50*
186 *-k 1 --best --strata*). Reported alignment genome positions for successfully mapped
187 putative mature miRNAs were elongated upstream and downstream, thus ensuring an
188 adequate pre-miRNA reconstruction. As no prior knowledge about the 3p or 5p identity
189 of putative mature miRNA sequences was available for porcine small RNA-seq data, two
190 candidate pre-miRNA structures were generated for each expressed sequence. The same
191 procedure was applied to human mature miRNAs when 3p or 5p identity was not
192 specified. Candidate sequences that were aligned and extracted from overlapping regions

193 corresponding to other annotated non-miRNA non-coding loci were discarded from
194 further analyses.

195 Elongation patterns were based on previously reported pre-miRNA favored size, with a
196 stem length of $\sim 35 \pm 3$ nt and an apical loop ≥ 10 nt [42,43]. With these specifications, we
197 established two upstream and three downstream elongation pattern combinations: 1) from
198 the starting genome position of each aligned sequence, 15 and 30 nt were added upstream,
199 beginning from each mature miRNA sequence start position. 2) Additionally, 60, 70 and
200 80 nt were added from each miRNA end position, resulting in the following elongation
201 pattern combinations for each candidate sequence: 15/60, 30/60, 15/70, 30/70, 15/80 and
202 30/80 added nt (i.e. we generated a total of 12 putative elongated pre-miRNA candidates
203 per each aligned sequence). Besides, the presence of flanking microprocessor motifs was
204 assessed for positionally correcting the elongated pre-miRNA candidate sequences.
205 Downstream CNNC and upstream UG motifs were assessed within the 30/60, 30/70 and
206 30/80 elongated candidates for each sequence, as described in [44], whereas downstream
207 mismatched GHG and upstream CHC motifs were searched in 15/60, 15/70 and 15/80
208 candidates [42].

209 To determine the most prevalent positional range of flanking processing motifs
210 surrounding pre-miRNA sequences in the porcine genome, 30 and 15 nt were added at
211 the flanking positions of annotated porcine pre-miRNAs available at the curated
212 miRCarta database [2]. The presence of CNNC and UG motifs within flanking ± 30 nt, as
213 well as GHG and CHC motifs within ± 15 nt was hence assessed. According to positional
214 results (Figure 2A), the CNNC and UG flanking motifs appeared more prominently
215 located 18 nt after miRNA gene ending and 12 nt before miRNA starting points,
216 respectively. Therefore, when downstream CNNC or upstream UG motifs were found
217 within ± 30 nt flanking windows along pre-miRNA candidates, -18 and $+12$ nt positions

218 were added from CNNC and UG motifs location, respectively, so as to establish accurate
219 miRNA genes boundaries determined by the microprocessor machinery. In the event that
220 none of the aforementioned motifs within flanking upstream and/or downstream defined
221 regions were found, the original elongated pre-miRNA candidates with no motif-based
222 positional refinement were kept.

223

224 **Selecting putative pre-miRNA candidate sequences based on structural integrity**

225 To better assess the optimal elongation pattern for each candidate sequence, the structural
226 stability of the 12 pre-miRNA candidates per single sequence was determined based on
227 the randfold algorithm [45]. This approach assumes the estimated minimum free energy
228 (MFE) of the folded pre-miRNA hairpin to be consistently lower than that of other
229 random sequences resembling hairpin-like folded structures [45]. Based on this property
230 of pre-miRNA sequences, we implemented a Monte Carlo randomization test to select
231 the most stable hairpin, i.e. those having the least folding minimum free energy (MFE)
232 values among the 12 previously generated candidates during pre-miRNA elongation
233 reconstruction for each of the analyzed sequences. To this end, we generated a total of
234 100 randomized sequences per candidate by shuffling their nucleotide distribution while
235 maintaining k-let counts [46]. The corresponding MFE values for each shuffled and
236 original hairpin-folded sequences were calculated with the RNAfold tool [37] and the
237 structural integrity score (p) was defined as:

238

$$239 \quad p = \frac{R}{N + 1}$$

240

241 where R is the number of randomized sequences having an MFE value equal or smaller
242 than that of the MFE value of the original sequence and N is the number of generated
243 iterations (100 in this study).

244 Subsequently, the candidate sequence showing the higher structural integrity (i.e. the one
245 showing the smallest p score) among all 12 generated pre-miRNA candidates per
246 sequence was selected. The proportion of the most structurally stable sequences for each
247 elongation pattern is shown in Figure 2B. When two or more sequences had equal p scores
248 (i.e. they had equivalent structural stability irrespective of the elongation pattern) the
249 reconstructed candidates belonging to the motif-corrected (if available) and shortest
250 elongation pattern were retained. The proportion of each elongation pattern selected as
251 the most structurally stable among all 12 tested patterns from expression-based and
252 homology-based data is shown in Figure 2C and D, respectively.

253

254 **Candidates classification with semi-supervised transductive learning**

255 After defining training and candidate data sets, we selected a total of 100 features
256 representing structural and statistical properties from each pre-defined sequence. These
257 extracted features have been previously reported in other state-of-the-art methods and
258 thoroughly reviewed in [47]. A complete list of all used features is shown in Table 1.

259 For pre-miRNA classification, the *miRNAss* algorithm proposed by Yones et al. [31] was
260 applied. This method implements a semi-supervised transductive learning scheme by
261 using well defined labeled cases, either positives (annotated pre-miRNAs) or negatives
262 (comprising other annotated non-coding hairpin-like sequences and unlabeled cases with
263 unknown hairpins), in order to draw a graph-based representation of each sequence based
264 on input features. Each node in the graph represents a sequence, whereas the

265 corresponding edges account for the expected similarities among them. In order to
266 accurately represent the spatial distribution and connections of each node, the feature
267 importance is obtained by applying the Relief-F algorithm [48,49], where k-nearest
268 predictors are weighted based on conditional dependencies among all the considered
269 features and the response vector of labels. This algorithm penalizes those predictor
270 features giving different values to k-neighbors from the same label class and vice versa.
271 After graph construction, a prediction score is assigned to each sequence node [31].

272 Sscrofa11.1 pre-miRNA sequences from Ensembl and miRCarta databases were
273 evaluated and different imbalance ratios between positive (taken as reference) and
274 negative data sets were applied to assess the performance of the classification algorithm
275 for miRNA discovery in the porcine genome (i.e. 1:1, 1:2, 1:10, 1:20, 1:40, 1:60, 1:80,
276 1:100, 1:150 and 1:200 imbalance ratios were considered). Labeled sequences comprised
277 annotated pre-miRNAs (+1) as positive sequences, while other non-coding hairpin-like
278 transcripts (-1) were considered as negative. Genome-wide randomly extracted hairpins
279 were assigned as unlabeled cases (0) within the negative data set.

280 Testing subsets were randomly assigned from all proposed imbalanced training data set
281 combinations using a 0.25 ratio. The performance of the classification algorithm for
282 miRNA identification was assessed with a total of 100 random Monte Carlo iterations
283 and average performance measures based on Sensitivity (SE), Specificity (SP), Accuracy
284 (Acc), F-1 score (F1) and Adjusted Geometric-mean (Agm) [50] were estimated (Figure
285 3A). Furthermore, we evaluated the performance for each imbalance scenario by
286 computing the corresponding Receiver Operating Characteristics (ROC) curves and the
287 Precision-Recall (PR) curves. PR curves can be more informative than ROC curves for
288 highly imbalanced data sets [51]. ROC and PR curves as well as the corresponding Areas
289 under the curve (AUC) estimates are shown in Figure S1 and Table S1. The ability of the

290 algorithm to correctly classify the list of Ensembl and miRCarta annotated porcine
291 miRNAs was also assessed by incorporating the positive data set as unlabeled candidate
292 sequences during the classification process in each of the defined imbalance scenarios.
293 Results for annotated porcine miRNAs assignment are shown in Table S2.

294 Finally, the reconstructed expressed candidate sequences from the porcine small RNA-
295 seq data and *H. sapiens* homologous miRNAs detected in the porcine genome were used
296 for identifying putative novel miRNAs. For this purpose, annotated pre-miRNAs from
297 the Ensembl database were used as positive class and other hairpin-like sequences were
298 considered as either negative or unlabeled sequences. Candidates classification was
299 implemented with all previously proposed imbalance ratios. In order to reduce the false
300 positive rate (i.e. reducing the misclassification of non-miRNA short hairpins as true
301 miRNA candidates), the Ensembl miRNA data set was defined as the positive class, due
302 to its higher overall reported specificity (Figure 3A and B). Prediction of novel miRNA
303 candidates was carried independently with every defined imbalance ratio. Only
304 candidates consistently reported as putative miRNAs in all imbalance scenarios were kept
305 in order to minimize the number of false positive miRNA candidates, albeit probably at
306 the expense of increasing the false negative rate.

307 Besides, for homology-based predicted novel pre-miRNA candidates, we calculated the
308 proportion of shared neighboring genes (setting a 2 Mb window before and after each
309 annotated human miRNA detected in the porcine genome) present in both *S. scrofa* and
310 *H. sapiens* assemblies and expressed as a Neighborhood Score (N):

311

312

313

$$N = \frac{G_r \cap G_i}{G_r}$$

314

315 where Gr is the number of orthologous genes within the 4 Mb window in the model
316 species (*H. sapiens*) and Gi is the number of genes within the same window in the species
317 of interest (*S. scrofa*). Only homology-based novel pre-miRNA candidates with $N > 0.1$
318 were considered for further analyses, based on the assumption that microRNAs residing
319 in genomic regions with surrounding and/or host genes phylogenetically conserved across
320 species are more prone to be integrated in biologically relevant transcriptional networks
321 [52].

322 **Benchmarking for miRNA prediction performance**

323 One of the most cited and used prediction miRNA algorithms is miRDeep. This tool was
324 developed by Friedländer et al. [53], and further improvements were made in subsequent
325 updates [11,54]. This algorithm implements a series of heuristics to compute a score for
326 each miRNA candidate expressing the log-odds probability of a sequence being a true
327 miRNA gene against the probability of being a miRNA-like pseudo-hairpin [53]. In order
328 to benchmark the eMIRNA pipeline compared with the widely used miRDeep approach,
329 we used the miRDeep2 algorithm [54] to identify novel and annotated miRNAs by using
330 the same small RNA-seq data set employed for *de novo* miRNA identification with the
331 eMIRNA pipeline. To ensure a fair comparison, the arf alignment file needed for running
332 the miRDeep2 software was generated from the eMIRNA alignment pipeline using the
333 bowtie tool (*bowtie -n 2 -l 25 -m 50 -k 1 --best --strata*) on pre-filtered expressed small
334 RNA sequences generated in this study. After running the miRDeep2 algorithm, both
335 novel and already annotated pre-miRNA candidates were compared with those obtained
336 with the eMIRNA pipeline. The positional accuracy of the annotated pre-miRNA
337 candidates concurrently identified with both approaches was then determined using the

338 Ensembl annotation available for the Sscrofa11.1 assembly. To further determine which
339 of the two approaches provided a better positional annotation of predicted miRNAs, the
340 deviation rate (dr) of each miRNA gene commonly detected was calculated for both
341 eMIRNA and miRDeep2, expressed as the average number of upstream and downstream
342 overhanging nucleotides compared with the latest porcine miRNA Ensembl annotation
343 (v97). The differential deviation estimate (ΔD) was assessed separately for each predicted
344 pre-miRNA candidate, as follows:

345

$$346 \quad \Delta D = eMIRNA_{dr} - miRDeep2_{dr}$$

347

348 Additionally, the performance statistics of the semi-supervised transductive learning
349 method [31] implemented in the eMIRNA pipeline was compared with other canonical
350 widely used state-of-the-art supervised ML approaches for miRNA prediction, such as
351 Support Vector Machine (SVM), Random Forest (RF), K-nearest Neighbors (KNN),
352 Naïve Bayes (NB), Extreme Gradient Boosting Trees (XGB) and Light Gradient Boosting
353 Trees (LGBM). Only labeled positive and negative data sets were used for comparison
354 between semi-supervised and supervised algorithms. Training and testing subsets were
355 randomly generated with a 0.25 ratio for testing data and commonly used with all the
356 proposed methods. No imbalance correcting procedure was applied. The comparative
357 performance of these tools was assessed on the basis of SE, SP, F1-score, ROC and PR
358 curves obtained for each algorithm implementation. SVM, RF, KNN and NB algorithms
359 were trained allowing 10 iterations for parameter tuning and a 10-fold cross-validation
360 scheme, using built-in functions included in the *caret* R package [55]. The *xgboost* [56]
361 and *lightgbm* (<https://github.com/microsoft/LightGBM/tree/master/R-package>) R

362 packages with default parameters were employed for the training of XGB and LGBM
363 classifiers, respectively.

364

365 **Experimental confirmation of novel identified porcine miRNAs through the RT-** 366 **qPCR analysis of an independent Göttingen minipig population**

367 In order to investigate the existence of several of the novel putative predicted miRNAs in
368 the porcine genome, three well established orthologous novel miRNA candidates detected
369 by homology-based search and not previously annotated in the Sscrofa11.1 assembly
370 were selected (hsa-miR-483-3p, hsa-miR-484-5p and hsa-miR-200a-3p). The existence
371 of miRNA genes orthologous to hsa-miR-483-3p and hsa-miR-484-5p was supported by
372 the identification of the corresponding expressed mature miRNA sequences in our small
373 RNA-seq data set. Transcripts corresponding to hsa-miR-200a-3p were detected at very
374 low expression levels (RPM < 10) in the porcine skeletal muscle transcriptomic data, so
375 they were not considered as biologically relevant or functionally active in our
376 experimental conditions. *Longissimus dorsi* muscle and liver RNA samples were
377 collected from an independent Göttingen minipig population [57]. A total of 7 extracted
378 RNA samples from muscle and liver tissues were randomly selected and cDNA synthesis
379 was carried out as reported by Balcells et al. [58]. Primers for the qPCR amplification of
380 miRNAs were designed with the miRprimer software [59] according to described
381 protocols [60] and they are indicated in Table S3.

382 MiRspecific qPCR was performed on a MX3005P machine (Stratagene, USA). Briefly,
383 1 µl of cDNA diluted 8 fold, 5 µl of 2× QuantiFast SYBR Green PCR master mix (Qiagen,
384 Germany) and 250 nM of each primer (Table S3) were mixed in a final volume of 10 µl.
385 Cycling conditions were: 95 °C for 5 min followed by 40 cycles of 95 °C for 10 s and
386 60 °C for 30 s. Melting curve analyses (60 °C to 99 °C) were performed after completing

387 amplification reaction to ensure the specificity of the assays. Data were processed with
388 the MxPro qPCR associated software. Assays were considered successful when: 1) the
389 melting curve was specific (1 single peak) and 2) the samples had Cq values <33 cycles
390 (i.e. sufficiently expressed to be considered biologically functional). Finally, amplified
391 products for muscle and liver samples were visually inspected by electrophoresis in a 3%
392 agarose gel.

393

394

395 **Results**

396 **Motif-based positional refinement enhances structural stability of pre-miRNA** 397 **candidates**

398 We have evaluated the usefulness of previously reported flanking motifs that enhance
399 pre-miRNA processing [42,44] as possible novel determinants for pre-miRNA
400 reconstruction from mature sequences. The presence of UG and CHC motifs in upstream
401 flanking regions as well as of downstream CNNC and GHG motifs was assessed in the
402 curated porcine miRNA annotation available in the miRCarta database [2] (Figure 2A).
403 Consistent with data reported by Fang et al. [42] and Auyeung et al. [44], the most
404 common flanking upstream positions for UG and CHC motifs from the 5' start of the
405 porcine pre-miRNA genes were -13/-12 and -7/-5, respectively, whereas for
406 downstream CNNC and GHG motifs, the most common position from the 3' end of the
407 pre-miRNA genes were +18/+21 and +4/+6 (Figure 2A).

408 Moreover, we determined the percentage of annotated porcine miRNAs that were
409 surrounded by the aforementioned processing motifs, allowing ± 2 nt of positional
410 variation from their corresponding expected sites. From a total of 328 confidently

411 annotated porcine pre-miRNAs in the miRCarta database [2], CNNC, UG, GHG and
412 CHC flanking motifs were found in 53.05%, 42.68%, 30.79% and 33.54% of the
413 sequences, respectively. The high frequency of the CNNC motif agrees well with its key
414 role in the correct Drosha ribonuclease III (DROSHA) positioning through the
415 recruitment of Serine and Arginine rich splicing factor 3 (SRSF3) at the basal junction of
416 the processed pri-miRNA [61]. The proportion of the three other flanking motifs were
417 also consistent with previously reported surveys [42,44].

418 To further elucidate the contribution of each motif to better delineate the boundaries of
419 pri-miRNA processing, we compared the structural stability (i.e. the estimated p score of
420 the hairpin secondary structure with the randfold approach [45]) for every pre-miRNA
421 candidate in each of the 12 generated elongation patterns per sequence (15/60, 30/60,
422 15/70, 30/70, 15/80 and 30/80, with and without taking into account motif search
423 positional refinement). As depicted in Figure 2B, predictions of candidate miRNA
424 sequences based on positional information obtained through processing motif search
425 showed a consistently increased structural stability compared with non-positionally
426 corrected original sequences. This phenomenon was less evident for shorter elongation
427 patterns, where the structural stability of the positionally corrected hairpins resembled
428 that of non-corrected candidates (Figure 2B). In certain cases, both approaches resulted
429 in equally stable secondary structures. Furthermore, shorter elongation patterns appeared
430 to be more favored than their longer counterparts, showing higher overall structural
431 stability both in small RNA-seq and homology-based derived candidate sequences
432 (Figure 2C and D). This result highlights that the preferred length for pre-miRNA
433 processed transcripts would be approximately in the range of 80 to 90 nt, with few cases
434 showing longer stable hairpin structures. Interestingly, this favored pre-miRNA length
435 interval coincides with that reported by Roden et al. [43], who determined a preferred 2×

436 stem length of 35 nt and a terminal loop of ~10 nt, accounting for a total pre-miRNA
437 sequence length of ~80 nt. Indeed, the average length of annotated pre-miRNAs in the
438 porcine genome after filtering for secondary structure and sequence similarity was
439 84.63 nt, also in accordance with results obtained after selecting the most structurally
440 stable elongation pattern from all generated candidates per sequence.

441

442 **Classifier performance and feature importance**

443 For assessing the performance of transductive semi-supervised miRNA classification on
444 the porcine transcriptome, Ensembl and miRCarta positive pre-filtered porcine miRNA
445 data sets (415 Ensembl and 244 miRCarta non-redundant hairpin-like stable annotated
446 miRNAs) were tested against selected non-coding hairpin-like sequences (252 annotated
447 non-coding hairpin-like RNAs other than miRNAs) and different imbalance ratios were
448 applied by incorporating genome-wide randomly extracted hairpins (unlabeled). Overall,
449 SE and SP obtained with the Ensembl miRNA data set (Figure 3A) were slightly better
450 than those inferred for the miRCarta data set (Figure 3B). Ensembl average SE and SP
451 were 0.9199 and 0.9101 respectively, whereas results obtained with the miRCarta data
452 set were slightly worse (SE = 0.8975, SP = 0.9019). Optimal performance was achieved
453 by using a balanced ratio between positive and negative classes, with a slightly
454 descending trend in the classifier performance when increasing the imbalance ratio
455 (Figure 3A and B), a result that was also observed when analyzing the ROC and PR curves
456 (Figure S1). When we compared the performance of the semi-supervised approach vs that
457 of other supervised algorithms, the *miRNA*ss algorithm [31] implemented in the eMIRNA
458 pipeline outperformed the rest of supervised approaches, with the exception of IGBM,
459 which showed similar performance results (Table 2). SP, as well as AUROC and AUPR
460 estimates obtained with the *miRNA*ss method [31] showed its high ability to discard false

461 positives miRNA candidates, at the cost of a lower SE (Table 2). Additionally, after
462 evaluating the ability of the algorithm to correctly identify the annotated porcine miRNA
463 loci in all defined imbalance scenarios, a total of 399 (89.92%) and 213 (87.30%)
464 annotated miRNAs were consistently classified as miRNA sequences using Ensembl
465 (415) and miRCarta (244) positive databases, respectively.

466 The improved performance achieved with the Ensembl data set was expected because
467 Ensembl annotation includes a more diverse and complete miRNA catalogue (415) than
468 miRCarta (244). However, these differences are probably due to a more strict miRNA
469 annotation procedure in the case of miRCarta database [2], which only includes manually
470 curated bona fide miRNA genes. Nevertheless, the slight increase in overall performance
471 observed in the Ensembl miRNA data set evidenced that even when reducing the set of
472 positive sequences to a more stringent annotation, as that available in the miRCarta
473 database [2], the ability of the eMIRNA pipeline to accurately distinguish miRNA
474 sequences from other non-miRNA hairpins remained almost unaltered.

475 Besides, we determined the importance of the set of calculated features for classifying the
476 miRNA candidates with the relief-F algorithm [48,49]. The estimated importance of the
477 30 most discriminant features is depicted in Figure 3C. The estimated impact of each
478 feature on the accuracy of miRNA is shown in Table S4. Structural stability-related
479 features accounted for the most important variables for classifying miRNAs correctly
480 (MFEadj, EFEadj, MFE, EFE, MEAFE, MFEadj.GC and CFE). All of these parameters
481 represented different hairpin structure folding statistics and they were highly
482 intercorrelated (Figure 3D). The discriminant power of structural stability features is
483 better exemplified in Figure 3E, where Ensembl annotated pre-miRNA sequences had an
484 overall higher structural stability (i.e. lower MFEadj values) compared with that of other
485 non-coding hairpin-like RNA sequences. These results clearly show the utmost

486 importance of the structural folding configuration in order to discriminate true miRNA
487 candidates from other hairpin-like sequences, hence supporting the need of a careful
488 determination of pre-miRNA boundaries.

489

490 **Novel porcine miRNA identified in the muscle transcriptome and by homology-** 491 **based search**

492 After microRNA identification from the porcine small RNA-seq data set, a total of 1,403
493 reconstructed pre-miRNA candidates from expressed transcripts were successfully
494 identified as putative novel miRNAs in the porcine *gluteus medius* transcriptome, which
495 corresponded to 160 unique miRNA loci after assigning clustered isomiRs to consensus
496 single miRNA genes. Among these, 140 consensus candidates (87.5%) overlapped
497 already annotated miRNAs in the porcine genome, whereas the 20 remaining ones
498 (12.5%) were classified as novel miRNA candidates.

499 Regarding homology-based search miRNA discovery in the porcine assembly
500 (Sscrofa11.1), a total of 310 annotated human miRNAs had orthologous miRNA genes
501 in the porcine genome. The already annotated miRNAs in the porcine genome comprised
502 281 (90.64%) of the 310 homologous miRNAs detected with eMIRNA (File S3), and the
503 29 ($N > 0.1$) remaining candidates were classified as novel non-previously annotated
504 homologous miRNAs in the porcine assembly (Table 3). The miR-483 and miR-484
505 genes were also identified as novel expressed miRNA candidates in the *gluteus medius*
506 muscle transcriptome generated in our small RNA-seq experiment. A complete list of the
507 novel miRNA candidates obtained with *de novo* and homology-based approaches is
508 shown in Table 3. The full list of detected miRNAs that had been already annotated and
509 all isomiRs associated with novel miRNA sequences can be found in File S3. The
510 existence of multiple isoform candidates for single predicted miRNA loci, either

511 displaying polymorphisms within the mature miRNA sequence or corresponding to 5' or
512 3'-trimming variations (File S3), evidenced the wide variety of isomiR sequences
513 expressed at significant levels in our *gluteus medius* muscle transcriptomic data set.

514

515 **The eMIRNA pipeline accurately recalls miRNA loci**

516 The same *gluteus medius* skeletal muscle transcriptomic data from the small RNA-seq
517 experiment employed for de novo miRNA discovery with the eMIRNA pipeline was used
518 for running the miRDeep2 algorithm [54]. A total of 148 transcripts belonging to 134
519 unique annotated miRNA loci were identified with miRDeep2. These numbers were
520 slightly smaller than the 140 annotated porcine miRNAs recovered as expressed
521 transcripts by the eMIRNA pipeline. Among these, 126 annotated miRNAs (85.14%)
522 were consistently recovered with eMIRNA and miRDeep2, 14 (9.46%) were only
523 reported by eMIRNA, and 8 (5.41%) were exclusively predicted by miRDeep2 (Table
524 S5).

525 Regarding novel candidates, miRDeep2 was able to recover a total of 11 putative novel
526 candidates belonging to 10 unique loci (Table S6). Seven of these candidates displayed
527 an estimated probability of being a true positive miRNA above 19% (miRDeep2
528 score ≥ 4 , Table S6). Noteworthy, two of the putatively true miRNAs detected by
529 miRDeep2 spanned other previously annotated non-coding RNAs in the porcine assembly
530 and were hence considered as miRNA-like false positives (Table S6). Among the 5
531 remaining candidates, 4 of them (miR-193a, miR-26a, miR-106b and miR-17) spanned
532 other already annotated miRNAs in the porcine assembly and were thus wrongly
533 classified as novel miRNAs by miRDeep2. The remaining candidate corresponded to
534 miR-483, which had already been identified with the eMIRNA pipeline (Table 3, Table
535 S6).

536 When comparing the accuracy of miRNA loci boundaries determined by the eMIRNA
537 pipeline and miRDeep2, the eMIRNA approach demonstrated an overall better capability
538 to accurately assign miRNA boundaries according to data from porcine miRNA loci
539 annotated in the Ensembl database. A total of 103 out of 126 (81.74%) annotated miRNA
540 genes detected by both eMIRNA and miRDeep2 showed reduced ΔD values (Table S7).
541 This result implies that genomic positions of miRNA precursors predicted with the
542 eMIRNA pipeline were more concordant with the annotation of the Sscrofa11.1 assembly
543 than those predicted with miRDeep2. This outcome illustrates the effectiveness of motif
544 search positional correction for reconstructing pre-miRNA candidates with a higher
545 reliability than the fixed elongation patterns strategy used by miRDeep2 [54]. Three of
546 the miRNA candidates showed no differences in positional accuracy between both
547 approaches, while the positions of the remaining sequences (15.87%) were more
548 accurately predicted with miRDeep2 (Table S7).

549

550 **Experimental confirmation of the existence of three novel miRNAs in the muscle** 551 **and liver tissues of Göttingen minipigs**

552 The RT-qPCR analyses allowed us to detect the expression of the novel ssc-miR-483,
553 ssc-miR-484 and ssc-miR-200a candidates in both *longissimus dorsi* skeletal muscle and
554 liver tissues (Figure S2A and B) retrieved from Göttingen minipigs. Both ssc-miR-483
555 and ssc-miR-484 were also detected as consistently expressed in the skeletal muscle of
556 Duroc gilts from our small RNA-seq experiment. The ssc-miR-200a was also detected in
557 our generated data set but at very low expression levels. Nevertheless, its expression was
558 further confirmed independently by RT-qPCR analyses. Amplification profiles and
559 melting curves for the three novel miRNA candidates detected by RT-qPCR are shown
560 in File S4.

561 **Discussion**

562 In the discovery of novel miRNA genes, one essential issue is the generation of pre-
563 miRNA sequence candidates, given that the majority of miRNA prediction tools are based
564 on feature extraction from the well-defined pre-miRNA hairpin structure [62]. At the
565 cellular level, the most abundant and stable miRNA transcripts are the mature miRNA
566 forms. Indeed, precursor stages, such as pri or pre-miRNAs, are much less abundant and
567 have shorter half-lives than mature miRNAs [63,64]. Therefore, the accurate definition
568 of pre-miRNA boundaries reconstructed from mature miRNAs is a crucial issue in order
569 to predict folding structure and minimum free energy (MFE) estimates in a robust manner.

570 Noteworthy, the majority of state-of-the-art methods for miRNA prediction are solely
571 focused on the miRNA classification of predefined candidate sequences. Moreover, many
572 of them do not contemplate the generation of such candidates for the identification of
573 unannotated miRNAs. On the contrary, they rely on well-known hairpins or on sets of
574 manually curated candidate sequences that are embedded in their prediction pipelines
575 [30,31,65-72].

576 Several other algorithms take advantage of the automated generation of hairpin
577 candidates, adopting fixed defined elongation patterns in order to reconstruct pre-miRNA
578 candidates from mature miRNA sequences [9,11,73,74]. However, fixed assumptions
579 about elongation patterns do not take into consideration the expected variable length of
580 pre-miRNA loci, and tend to generate candidate sequences that, despite harboring mature
581 miRNAs, might have unreliable boundaries. This may lead to inaccuracies in the folding
582 prediction and thus to an augmentation of the false negative rate. Even worse, non-
583 miRNA hairpin-like sequences strongly resembling pre-miRNAs may be generated
584 through the blind elongation of short sequences, which could result in the emergence of
585 false positive candidates. This situation is particularly critical when analyzing the

586 reliability of miRNA annotation in public databases [27,75,76]. Other approaches have
587 also adopted a multiple hairpin candidate search for each query sequence to further select
588 those showing a higher structural stability [77-79]. By using this strategy, we explored
589 the influence of flanking processing motifs on the accurate determination of the length
590 and boundaries of pre-miRNA candidates. By doing so, we have demonstrated that the
591 inclusion of processing motif search criteria for the estimation of pre-miRNA boundaries
592 resulted in an improved ability to better assess the optimal candidate sequences to be used
593 for miRNA prediction.

594 Compared with miRDeep2 [54], the eMIRNA pipeline showed an improved ability to
595 better assess the already annotated miRNA loci boundaries after pre-miRNA sequence
596 reconstruction. However, the presence of embedded processing motifs within the
597 boundaries of miRNA genes is not a universal feature, with a non-negligible amount of
598 miRNA loci lacking the well-known CNNC and UG motifs [44], as well as the CHC and
599 GHG mismatches [42] in their proximal surroundings. Additional work is needed to better
600 characterize other processing motifs or structural determinants that may also contribute
601 to miRNA maturation.

602 In contrast with pre-existing supervised methods for miRNA discovery, few semi-
603 supervised methods have been developed for such purpose [31,80]. From a biological
604 perspective, the scarce miRNA annotation typically found in non-model species poses a
605 great challenge when attempting to predict novel miRNA loci uniquely based on labeled
606 data. This happens because the amount of unknown non-miRNA sequences with hairpin-
607 like secondary structures is expected to be hundreds of times larger than the number of
608 confidently annotated miRNAs to be used for training supervised algorithms. Despite the
609 fact that good performance statistics may be obtained after classifier training, supervised
610 algorithms heavily depend on the existence of an extensive miRNA annotation. Indeed,

611 the ability of such classifiers to detect unannotated miRNA sequences is mainly driven
612 by the amount and diversity of positive and negative instances used for learning training.
613 On the contrary, semi-supervised transductive approaches [31] are able to overcome such
614 limitation by incorporating unlabeled cases to the training process, with the aim of
615 increasing the variability of the data used for target sequences classification. In fact,
616 allowing the classifier to check hundreds or thousands of unknown unlabeled sequences
617 has proven to increase the validity of microRNA prediction over other methods solely
618 based on labeled data [31], a result that was also verified when comparing the semi-
619 supervised approach used in this study with other broadly reported supervised methods
620 (Table 2). This strategy is particularly reliable when few positive data are available and
621 the annotated negative data set only represent a small proportion of the whole non-
622 miRNA class. Besides, in classification problems where the negative class is expected to
623 be dozens or hundreds of times larger than the positive class, the accurate identification
624 of false positives is crucial. Indeed, such scenario is completely applicable to miRNAs,
625 where thousands of non-miRNA sequences exist compared with the few hundreds of
626 reliably annotated miRNA genes, and the annotation of negative hairpin-like sequences
627 only represents a small proportion of the whole non-miRNA class.

628 After miRNA prediction, the detection of multiple isoforms for each single predicted
629 miRNA loci evidenced the existence of a broad array of isomiR sequences expressed at
630 significant levels in our *gluteus medius* muscle transcriptomic data set (File S3). Previous
631 studies have highlighted the importance of isomiRs in expanding the biological diversity
632 of miRNA function [81-84]. Like canonical miRNAs, isomiRs are also evolutionary
633 conserved [81]. Both 5' and 3' miRNA isoforms can be generated either from alternative
634 processing sites of DROSHA and Dicer [43,85] or from post-transcriptional

635 modifications, influencing miRNA half-lives as well as their interactions with RNA-
636 binding proteins (RBPs) [86,87].

637 More recently, other integrative approaches have addressed the detection of isomiRs and
638 the potential functional influence that subtle modifications in the 3' and 5' boundaries of
639 mature miRNA sequences might have on target recognition [88-91]. Other studies have
640 also reported 5' alternative processing events in a large number of miRNAs, contributing
641 to the expansion of their target repertoire at a higher rate than previously thought [92].
642 Despite these promising results, the biological implications of miRNA alternative
643 processing events leading to the generation of isomiRs are still poorly understood and
644 further research is needed in order to exclude potential biases in isomiR quantification
645 and functional validation, as variations in 3' or 5' ends of mature miRNAs can strongly
646 affect the reliability of stem-loop qPCR amplification protocols [93].

647 One potential limitation of our study is that 17 of the novel miRNAs predicted with
648 eMIRNA and based on muscle transcriptomic data have not been further investigated in
649 order to confirm their existence by RT-qPCR, so their experimental validation is still
650 pending. Indeed, we only investigated 3 out of 20 predicted novel porcine miRNAs.
651 Noteworthy, the three selected miRNAs were successfully confirmed as bona fide
652 miRNAs by RT-qPCR thus suggesting that eMIRNA predictions are accurate.

653 Among the three validated miRNAs, it is worth mentioning miR-483, which has been
654 functionally associated with cell growth regulation [94] as well as with insulin resistance
655 and metabolic syndrome susceptibility likely due to its strong implication in the
656 regulation of glucose metabolism [95,96]. Additionally, the expression of miR-483,
657 whose coding sequence maps to the second intron of the insulin growth factor 2 (*IGF2*)
658 gene, has been tightly associated with an enhancement of *IGF2* gene expression. This is
659 achieved through the binding of miR-483 to transcription factors in a positive feed-back

660 loop [97], although other authors have questioned such dependence [98]. Other relevant
661 successfully profiled miRNAs were ssc-miR-200a and ssc-miR-484. The miR-200a gene
662 has been mainly reported as a regulator of cell growth and differentiation through
663 targeting several protein-encoding transcripts like the growth factor receptor-bound 2
664 (*GRB2*), α -smooth muscle actin (*α -SMA*) or the fibroblast-specific protein-1 (*FSP-1*), thus
665 hampering the endothelial-mesenchymal transition [99]. Furthermore, miR-484 has been
666 associated with the inhibition of Fis1-mediated mitochondrial fission and apoptosis
667 signaling [100].

668

669

670 **Conclusions**

671 In this study we have implemented an end-to-end pipeline that may facilitate the
672 identification of novel miRNAs in the porcine genome. We have tested the eMIRNA
673 pipeline by following a homology-based approach making use of the well annotated
674 human microRNA transcriptome. Besides, we have analyzed the presence of non-
675 annotated miRNAs in the porcine genome using data from a small RNA-seq experiment
676 comprising muscle samples from 48 Duroc gilts. We have also taken into consideration
677 several issues that are critical to robustly predict miRNA genes, such as the accurate
678 reconstruction of candidate pre-miRNAs, the correct definition of negative training data
679 sets and the evaluation of the high class-imbalance phenomenon, which is not fully
680 addressed in many miRNA-prediction studies. In parallel, we have established hard-
681 threshold filtering steps to keep false positive predictions at a minimum. We have also
682 demonstrated the usefulness of positional refinement through flanking motif search to
683 better determine the boundaries of pre-miRNA hairpin-like candidate sequences. The

684 expression of several of the novel miRNAs described in this work was further confirmed
685 by RT-qPCR analyses. In the light of these results, we believe that the eMIRNA pipeline
686 will facilitate the discovery and annotation of novel miRNAs, thus broadening the
687 miRNA catalogue of non-model species with yet poorly annotated genome assemblies.

688

689

690 **Support**

691 The research presented in this publication was funded by grants AGL2013-48742-C2-1-
692 R and AGL2013-48742-C2-2-R awarded by the Spanish Ministry of Economy and
693 Competitiveness. E. Mármol-Sánchez was funded with a Ph.D. fellowship FPU15/01733
694 awarded by the Spanish Ministry of Education and Culture (MECD).

695

696 **Conflict of interest**

697 The authors declare no conflict of interest.

698

699 **Acknowledgements**

700 The authors would like to thank the Department of Veterinary Animal Sciences in the
701 Faculty of Health and Medical Sciences of the University of Copenhagen for providing
702 their facilities and materials for RT-qPCR experiments. We would also like to express
703 our gratitude to Dr. Caroline M. Junker Mentzel for kindly providing RNA samples for
704 qPCR analyses. We also acknowledge Selección Batallé S.A. for providing animal
705 material and the Spanish Ministry of Economy and Competitiveness for the Center of
706 Excellence Severo Ochoa 2016-2019 (SEV-2015-0533) grant awarded to the Center for

707 Research in Agricultural Genomics (CRAG). Thanks also to the CERCA Programme of
708 the Generalitat de Catalunya for their support.

709

710

711 **References**

712 [1] A. Kozomara, M. Birgaoanu, S. Griffiths-Jones, miRBase: from microRNA sequences
713 to function, *Nucleic Acids Res.* 47 (2019) D155–D162.

714 [2] C. Backes, T. Fehlmann, F. Kern, T. Kehl, H.-P. Lenhof, E. Meese, A. Keller,
715 miRCarta: a central repository for collecting miRNA candidates, *Nucleic Acids Res.* 46
716 (2018) D160–D167.

717 [3] B. From, D. Domanska, L. Høy, V. Ovchinnikov, W. Kang, E. Aparicio-Puerta, M.
718 Johansen, K. Flatmark, A. Mathelier, E. Hovig, M. Hackenberg, M.R. Friedländer, K.J.
719 Peterson, MirGeneDB2.0: the metazoan microRNA complement, *Nucleic Acids Res.*
720 (2019) gkz885.

721 [4] J. Meunier, F. Lemoine, M. Soumillon, A. Liechti, M. Weier, K. Guschanski, H. Hu,
722 P. Khaitovich, H. Kaessmann, Birth and expression evolution of mammalian microRNA
723 genes, *Genome Res.* 23 (2013) 34–45.

724 [5] M. Warnefors, A. Liechti, J. Halbert, D. Valloton, H. Kaessmann, Conserved
725 microRNA editing in mammalian evolution, development and disease, *Genome Biol.* 15
726 (2014) R83.

727 [6] L.P. Lim, N.C. Lau, E.G. Weinstein, A. Abdelhakim, S. Yekta, M.W. Rhoades, C.B.
728 Burge, D.P. Bartel, The microRNAs of *Caenorhabditis elegans*, *Genes Dev.* 17 (2003)
729 991–1008.

- 730 [7] E.C. Lai, P. Tomancak, R.W. Williams, G.M. Rubin, Computational identification of
731 *Drosophila* microRNA genes, *Genome Biol.* 4 (2003) R42.
- 732 [8] X. Wang, J. Zhang, F. Li, J. Gu, T. He, X. Zhang, Y. Li, MicroRNA identification
733 based on sequence and structure alignment, *Bioinformatics.* 21 (2005) 3610–3614.
- 734 [9] A. Mathelier, A. Carbone, MIRENA: finding microRNAs with high accuracy and no
735 learning at genome scale and from deep sequencing data, *Bioinformatics.* 26 (2010)
736 2226–2234.
- 737 [10] K. Qian, E. Auvinen, D. Greco, P. Auvinen, miRSeqNovel: An R based workflow
738 for analyzing miRNA sequencing data, *Mol. Cell. Probes* 26 (2012) 208–211.
- 739 [11] J. An, J. Lai, M.L. Lehman, C.C. Nelson, MiRDeep*: An integrated application tool
740 for miRNA identification from RNA sequencing data, *Nucleic Acids Res.* 41 (2013) 727–
741 737.
- 742 [12] T.B. Hansen, M.T. Venø, J. Kjems, C.K. Damgaard, miRIdentify: high stringency
743 miRNA predictor identifies several novel animal miRNAs, *Nucleic Acids Res.* 42 (2014)
744 e124.
- 745 [13] D. Kleftogiannis, A. Korfiati, K. Theofilatos, S. Likothanassis, A. Tsakalidis, S.
746 Mavroudi, Where we stand, where we are moving: surveying computational techniques
747 for identifying miRNA genes and uncovering their regulatory role, *J. Biomed. Inform.* 46
748 (2013) 563–573.
- 749 [14] M. Bortolomeazzi, E. Gaffo, S. Bortoluzzi, A survey of software tools for microRNA
750 discovery and characterization using RNA-seq, *Brief. Bioinform.* 20 (2017) 918–930.
- 751 [15] G. Stegmayer, L.E. Di Persia, M. Rubiolo, M. Gerard, M. Pividori, C. Yones, L.A.
752 Bugnon, T. Rodriguez, J. Raad, D.H. Milone, Predicting novel microRNA: a

753 comprehensive comparison of machine learning approaches, *Brief. Bioinform.* (2018)
754 bby037.

755 [16] A. Rajendiran, A. Chatterjee, A. Pan, Computational approaches and related tools to
756 identify microRNAs in a species: a bird's eye view, *Interdiscip. Sci. Comput. Life Sci.* 10
757 (2018) 616–635.

758 [17] J.-E. Long, H.-X. Chen, Identification and characteristics of cattle microRNAs by
759 homology searching and small RNA cloning, *Biochem. Genet.* 47 (2009) 329–343.

760 [18] Z. Wang, K. He, Q. Wang, Y. Yang, Y. Pan, The prediction of the porcine
761 pre-microRNAs in genome-wide based on support vector machine (SVM) and homology
762 searching, *BMC Genomics* 13 (2012) 729.

763 [19] X. Hou, Z. Tang, H. Liu, N. Wang, H. Ju, K. Li, Discovery of microRNAs associated
764 with myogenesis by deep sequencing of serial developmental skeletal muscles in pigs,
765 *PLoS One* 7 (2012) e52123.

766 [20] C. Yuan, X. Wang, R. Geng, X. He, L. Qu, Y. Chen, Discovery of cashmere goat
767 (*Capra hircus*) microRNAs in skin and hair follicles by Solexa sequencing, *BMC*
768 *Genomics* 14 (2013) 511.

769 [21] J. Sun, M. Li, Z. Li, J. Xue, X. Lan, C. Zhang, C. Lei, H. Chen, Identification and
770 profiling of conserved and novel microRNAs from Chinese Qinchuan bovine longissimus
771 thoracis, *BMC Genomics* 14 (2013) 42.

772 [22] T. Buza, M. Arick, H. Wang, D.G. Peterson, Computational prediction of disease
773 microRNAs in domestic animals, *BMC Res. Notes.* 7 (2014) 403.

774 [23] B. Sadeghi, H. Ahmadi, S. Azimzadeh-Jamalkandi, M.R. Nassiri, A. Masoudi-
775 Nejad, BosFinder: a novel pre-microRNA gene prediction algorithm in *Bos taurus*, *Anim.*
776 *Genet.* 45 (2014) 479–484.

777 [24] J. Wu, H. Zhu, W. Song, M. Li, C. Liu, N. Li, F. Tang, H. Mu, M. Liao, X. Li, W.
778 Guan, X. Li, J. Hua, Identification of conservative microRNAs in Saanen dairy goat testis
779 through deep sequencing, *Reprod. Domest. Anim.* 49 (2014) 32–40.

780 [25] Z. Li, H. Wang, L. Chen, L. Wang, X. Liu, C. Ru, A. Song, Identification and
781 characterization of novel and differentially expressed microRNAs in peripheral blood
782 from healthy and mastitis Holstein cattle by deep sequencing, *Anim. Genet.* 45 (2014)
783 20–27.

784 [26] D.M.D. Saçar, H. Hamzeiy, J. Allmer, Can miRBase provide positive data for
785 machine learning for the detection of miRNA hairpins? *J. Integr. Bioinform.* 10 (2013)
786 1–11.

787 [27] N. Ludwig, M. Becker, T. Schumann, T. Speer, T. Fehlmann, A. Keller, E. Meese,
788 Bias in recent miRBase annotations potentially associated with RNA quality issues, *Sci.*
789 *Rep.* 7 (2017) 5162.

790 [28] L. Wei, M. Liao, Y. Gao, R. Ji, Z. He, Q. Zou, Improved and promising identification
791 of human microRNAs by incorporating a high-quality negative set, *IEEE/ACM Trans.*
792 *Comput. Biol. Bioinforma.* 11 (2014) 192–201.

793 [29] M. Yousef, J. Allmer, W. Khalifa, Accurate plant microRNA prediction can be
794 achieved using sequence motif features, *J. Intell. Learn. Syst. Appl.* 8 (2016) 9–22.

795 [30] G. Stegmayer, C. Yones, L. Kamenetzky, D.H. Milone, High class-imbalance in
796 premiRNA prediction: a novel approach based on deepSOM, *IEEE/ACM Trans. Comput.*
797 *Biol. Bioinforma.* 14 (2017) 1316–1326.

798 [31] C. Yones, G. Stegmayer, D.H. Milone, C. Sahinalp, Genome-wide pre-miRNA
799 discovery from few labeled examples, *Bioinformatics.* 34 (2018) 541–549.

- 800 [32] Y. Wang, X. Li, B. Tao, Improving classification of mature microRNA by solving
801 class imbalance problem, *Sci. Rep.* 6 (2016) 25941.
- 802 [33] T.F. Cardoso, R. Quintanilla, J. Tibau, M. Gil, E. Mármol-Sánchez, O.
803 GonzálezRodríguez, R. González-Prendes, M. Amills, Nutrient supply affects the mRNA
804 expression profile of the porcine skeletal muscle, *BMC Genomics* 18 (2017) 603.
- 805 [34] M. Ballester, M. Amills, O. González-Rodríguez, T.F. Cardoso, M. Pascual, R.
806 González-Prendes, N. Panella-Riera, I. Díaz, J. Tibau, R. Quintanilla, Role of AMPK
807 signalling pathway during compensatory growth in pigs, *BMC Genomics* 19 (2018) 682.
- 808 [35] A.R. Quinlan, I.M. Hall, BEDTools: a flexible suite of utilities for comparing
809 genomic features, *Bioinformatics.* 26 (2010) 841–842.
- 810 [36] Y. Huang, B. Niu, Y. Gao, L. Fu, W. Li, CD-HIT suite: a web server for clustering
811 and comparing biological sequences, *Bioinformatics.* 26 (2010) 680–682.
- 812 [37] R. Lorenz, S.H. Bernhart, C. Höner zu Siederdisen, H. Tafer, C. Flamm, P.F.
813 Stadler, I.L. Hofacker, ViennaRNA Package 2.0, *Algorithms Mol. Biol.* 6 (2011) 26.
- 814 [38] C. Yones, HextractorR: Integrated tool for hairpin extraction of RNA sequences, R
815 Package Version 1.3, 2018 <https://cran.r-project.org/package=HextractorR>.
- 816 [39] M. Martin, Cutadapt removes adapter sequences from high-throughput sequencing
817 reads, *EMBnet.journal.* 17 (2011) 10.
- 818 [40] Y. Lu, A.S. Baras, M.K. Halushka, miRge 2.0 for comprehensive analysis of
819 microRNA sequencing data, *BMC Bioinforma.* 19 (2018) 275.
- 820 [41] B. Langmead, C. Trapnell, M. Pop, S. Salzberg, Ultrafast and memory-efficient
821 alignment of short DNA sequences to the human genome, *Genome Biol.* 10 (2009) R25.

822 [42] W. Fang, D.P. Bartel, The menu of features that define primary microRNAs and
823 enable de novo design of microRNA genes, *Mol. Cell* 60 (2015) 131–145.

824 [43] C. Roden, J. Gaillard, S. Kanoria, W. Rennie, S. Barish, J. Cheng, W. Pan, J. Liu, C.
825 Cotsapas, Y. Ding, J. Lu, Novel determinants of mammalian primary microRNA
826 processing revealed by systematic evaluation of hairpin-containing transcripts and human
827 genetic variation, *Genome Res.* 27 (2017) 374–384.

828 [44] V.C. Auyeung, I. Ulitsky, S.E. McGeary, D.P. Bartel, Beyond secondary structure:
829 primary-sequence determinants license pri-miRNA hairpins for processing, *Cell.* 152
830 (2013) 844–858.

831 [45] E. Bonnet, J. Wuyts, P. Rouze, Y. Van de Peer, Evidence that microRNA precursors,
832 unlike other non-coding RNAs, have lower folding free energies than random sequences,
833 *Bioinformatics.* 20 (2004) 2911–2917.

834 [46] M. Jiang, J. Anderson, J. Gillespie, M. Mayne, uShuffle: a useful tool for shuffling
835 biological sequences while preserving the k-let counts, *BMC Bioinforma.* 9 (2008) 192.

836 [47] I. Lopes, A. Schliep, A.C.L.F. de Carvalho, The discriminant power of RNA features
837 for pre-miRNA recognition, *BMC Bioinforma.* 15 (2014) 124.

838 [48] I. Kononenko, E. Šimec, M. Robnik-Šikonja, Overcoming the myopia of inductive
839 learning algorithms with RELIEFF, *Appl. Intell.* 7 (1997) 39–55.

840 [49] M. Robnik-Šikonja, I. Kononenko, Theoretical and empirical analysis of ReliefF and
841 RReliefF, *Mach. Learn.* 53 (2003) 23–69, <https://doi.org/10.1023/A:1025667309714>.

842 [50] R. Batuwita, V. Palade, Adjusted geometric-mean: a novel performance measure for
843 imbalanced bioinformatics data sets learning, *J. Bioinforma. Comput. Biol.* 10 (2012)
844 1250003.

845 [51] J. Davis, M. Goadrich, The relationship between precision-recall and ROC curves,
846 ACM Int. Conf. Proceeding Ser. (2006) 233–240.

847 [52] G.S. França, M.D. Vibranovski, P.A.F. Galante, Host gene constraints and genomic
848 context impact the expression and evolution of human microRNAs, Nat. Commun. 7
849 (2016) 11438.

850 [53] M.R. Friedländer, W. Chen, C. Adamidi, J. Maaskola, R. Einspanier, S. Knespel, N.
851 Rajewsky, Discovering microRNAs from deep sequencing data using miRDeep, Nat.
852 Biotechnol. 26 (2008) 407–415.

853 [54] M.R. Friedländer, S.D. MacKowiak, N. Li, W. Chen, N. Rajewsky, MiRDeep2
854 accurately identifies known and hundreds of novel microRNA genes in seven animal
855 clades, Nucleic Acids Res. 40 (2012) 37–52.

856 [55] M. Kuhn, Building predictive models in R using the caret package, J. Stat. Softw. 28
857 (2008) 1–26.

858 [56] T. Chen, C. Guestrin, XGBoost: A scalable tree boosting system, Proc. ACM
859 SIGKDD Int. Conf. Knowl. Discov. Data Min., 2016, pp. 785–794.

860 [57] C.M.J. Mentzel, C. Anthon, M.J. Jacobsen, P. Karlskov-Mortensen, C.S. Bruun, C.B.
861 Jørgensen, J. Gorodkin, S. Cirera, M. Fredholm, Gender and obesity specific microRNA
862 expression in adipose tissue from lean and obese pigs, PLoS One 10 (2015) e0131650.

863 [58] I. Balcells, S. Cirera, P.K. Busk, Specific and sensitive quantitative RT-PCR of
864 miRNAs with DNA primers, BMC Biotechnol. 11 (2011) 70.

865 [59] P.K. Busk, A tool for design of primers for microRNA-specific quantitative RT-
866 qPCR, BMC Bioinforma. 15 (2014) 29.

867 [60] S. Cirera, P.K. Busk, Quantification of miRNAs by a simple and specific qPCR
868 method, *Methods Mol. Biol.* (2014) 73–81.

869 [61] K. Kim, T. Duc Nguyen, S. Li, T. Anh Nguyen, SRSF3 recruits DROSHA to the
870 basal junction of primary microRNAs, *RNA*. 24 (2018) 892–898.

871 [62] D.P. Bartel, Metazoan microRNAs, *Cell*. 173 (2018) 20–51.

872 [63] L. Gan, B. Denecke, Profiling pre-microRNA and mature microRNA expressions
873 using a single microarray and avoiding separate sample preparation, *Microarrays*. 2
874 (2013) 24–33.

875 [64] Y. Guo, J. Liu, S.J. Elfenbein, Y. Ma, M. Zhong, C. Qiu, Y. Ding, J. Lu,
876 Characterization of the mammalian miRNA turnover landscape, *Nucleic Acids Res.* 43
877 (2015) 2326–2341.

878 [65] C. Xue, F. Li, T. He, G.-P. Liu, Y. Li, X. Zhang, Classification of real and pseudo
879 microRNA precursors using local structure-sequence features and support vector
880 machine, *BMC Bioinforma.* 6 (2005) 310, <https://doi.org/10.1186/1471-2105-6-310>.

881 [66] P. Jiang, H. Wu, W. Wang, W. Ma, X. Sun, Z. Lu, MiPred: classification of real and
882 pseudo microRNA precursors using random forest prediction model with combined
883 features, *Nucleic Acids Res.* 35 (2007) W339–W344.

884 [67] R. Batuwita, V. Palade, microPred: effective classification of pre-miRNAs for
885 human miRNA gene prediction, *Bioinformatics*. 25 (2009) 989–995.

886 [68] Y. Wu, B. Wei, H. Liu, T. Li, S. Rayner, MiRPara: a SVM-based software tool for
887 prediction of most probable microRNA coding regions in genome scale sequences, *BMC*
888 *Bioinforma.* 12 (2011) 107.

889 [69] A. Gudyś, M.W. Szcześniak, M. Sikora, I. Makałowska, HuntMi: an efficient and
890 taxon-specific approach in pre-miRNA identification, *BMC Bioinforma.* 14 (2013) 83.

891 [70] Q. Zou, Y. Mao, L. Hu, Y. Wu, Z. Ji, miRClassify: an advanced web server for
892 miRNA family classification and annotation, *Comput. Biol. Med.* 45 (2014) 157–160.

893 [71] D. Klefogiannis, K. Theofilatos, S. Likothanassis, S. Mavroudi, YamiPred: a novel
894 evolutionary method for predicting pre-miRNAs and selecting relevant features,
895 *IEEE/ACM Trans. Comput. Biol. Bioinforma.* 12 (2015) 1183–1192.

896 [72] D.M.D. Saçar, J. Baumbach, J. Allmer, On the performance of pre-microRNA
897 detection algorithms, *Nat. Commun.* 8 (2017) 330.

898 [73] D.M. Vitsios, E. Kentepozidou, L. Quintais, E. Benito-Gutiérrez, S. van Dongen,
899 M.P. Davis, A.J. Enright, Mirnovo: genome-free prediction of microRNAs from small
900 RNA sequencing data and single-cells using decision forests, *Nucleic Acids Res.* 45
901 (2017) e177.

902 [74] R.J. Peace, M. Sheikh Hassani, J.R. Green, miPIE: NGS-based prediction of miRNA
903 using integrated evidence, *Sci. Rep.* 9 (2019) 1548.

904 [75] M.J. Axtell, B.C. Meyers, Revisiting criteria for plant microRNA annotation in the
905 era of big data, *Plant Cell* 30 (2018) 272–284.

906 [76] J. Alles, T. Fehlmann, U. Fischer, C. Backes, V. Galata, M. Minet, M. Hart, M.
907 AbuHalima, F.A. Grässer, H.-P. Lenhof, A. Keller, E. Meese, An estimate of the total
908 number of true human miRNAs, *Nucleic Acids Res.* 47 (2019) 3353–3364.

909 [77] J. Lei, Y. Sun, miR-PREFeR: an accurate, fast and easy-to-use plant miRNA
910 prediction tool using small RNA-seq data, *Bioinformatics.* 30 (2014) 2837–2839.

911 [78] M. Evers, M. Huttner, A. Dueck, G. Meister, J.C. Engelmann, miRA: adaptable
912 novel miRNA identification in plants using small RNA sequencing data, *BMC*
913 *Bioinforma.* 16 (2015) 370, <https://doi.org/10.1186/s12859-015-0798-3>.

914 [79] C. Paicu, I. Mohorianu, M. Stocks, P. Xu, A. Coince, M. Billmeier, T. Dalmay, V.
915 Moulton, S. Moxon, miRCat2: accurate prediction of plant and animal microRNAs from
916 next-generation sequencing data sets, *Bioinformatics.* 33 (2017) 2446–2454.

917 [80] M. Sheikh Hassani, J.R. Green, Multi-view co-training for microRNA prediction,
918 *Sci. Rep.* 9 (2019) 10931.

919 [81] G.C. Tan, E. Chan, A. Molnar, R. Sarkar, D. Alexieva, I.M. Isa, S. Robinson, S.
920 Zhang, P. Ellis, C.F. Langford, P.V. Guillot, A. Chandrashekan, N.M. Fisk, L.
921 Castellano, G. Meister, R.M. Winston, W. Cui, D. Baulcombe, N.J. Dibb, 5' isomiR
922 variation is of functional and evolutionary importance, *Nucleic Acids Res.* 42 (2014)
923 9424–9435.

924 [82] A.G. Telonis, P. Loher, Y. Jing, E. Londin, I. Rigoutsos, Beyond the one-locus-one-
925 miRNA paradigm: microRNA isoforms enable deeper insights into breast cancer
926 heterogeneity, *Nucleic Acids Res.* 43 (2015) 9158–9175.

927 [83] F. Yu, K.A. Pillman, C.T. Neilsen, J. Toubia, D.M. Lawrence, A. Tsykin, M.P.
928 Gantier, D.F. Callen, G.J. Goodall, C.P. Bracken, Naturally existing isoforms of miR-222
929 have distinct functions, *Nucleic Acids Res.* 45 (2017) 11371–11385.

930 [84] P. Sheng, C. Fields, K. Aadland, T. Wei, O. Kolaczowski, T. Gu, B. Kolaczowski,
931 M. Xie, Dicer cleaves 5'-extended microRNA precursors originating from RNA
932 polymerase II transcription start sites, *Nucleic Acids Res.* 46 (2018) 5737–5752.

933 [85] B. Kim, K. Jeong, V.N. Kim, Genome-wide mapping of DROSHA cleavage sites on
934 primary microRNAs and noncanonical substrates, *Mol. Cell* 66 (2017) 258–269.e5.

935 [86] C.T. Neilsen, G.J. Goodall, C.P. Bracken, IsomiRs – the overlooked repertoire in the
936 dynamic microRNAome, *Trends Genet.* 28 (2012) 544–549.

937 [87] X. Bofill-De Ros, A. Yang, S. Gu, IsomiRs: expanding the miRNA repression
938 toolbox beyond the seed, *Biochim. Biophys. Acta - Gene Regul. Mech.* (2019) 194373.

939 [88] G. Urgese, G. Paciello, A. Acquaviva, E. Ficarra, isomiR-SEA: an RNA-seq analysis
940 tool for miRNAs/isomiRs expression level profiling and miRNA-mRNA interaction sites
941 evaluation, *BMC Bioinforma.* 17 (2016) 148.

942 [89] Y. Zhang, Q. Zang, B. Xu, W. Zheng, R. Ban, H. Zhang, Y. Yang, Q. Hao, F. Iqbal,
943 A. Li, Q. Shi, IsomiR Bank: a research resource for tracking IsomiRs, *Bioinformatics.* 32
944 (2016) 2069–2071.

945 [90] X. Bofill-De Ros, K. Chen, S. Chen, N. Tesic, D. Randjelovic, N. Skundric, S. Nestic,
946 V. Varjacic, E.H. Williams, R. Malhotra, M. Jiang, S. Gu, QuagmiR: a cloud-based
947 application for isomiR big data analytics, *Bioinformatics.* 35 (2019) 1576–1578.

948 [91] X. Bofill-De Ros, W.K. Kasprzak, Y. Bhandari, L. Fan, Q. Cavanaugh, M. Jiang, L.
949 Dai, A. Yang, T.-J. Shao, B.A. Shapiro, Y.-X. Wang, S. Gu, Structural differences
950 between pri-miRNA paralogs promote alternative Drosha cleavage and expand target
951 repertoires, *Cell Rep.* 26 (2019) 447–459.e4.

952 [92] H. Kim, J. Kim, K. Kim, H. Chang, K. You, V.N. Kim, Bias-minimized
953 quantification of microRNA reveals widespread alternative processing and 3' end
954 modification, *Nucleic Acids Res.* 47 (2019) 2630–2640.

955 [93] A. Schamberger, T.I. Orbán, 3' IsomiR species and DNA contamination influence
956 reliable quantification of microRNAs by stem-loop quantitative PCR, *PLoS One* 9 (2014)
957 e106315.

958 [94] T.H. Vu, N.V. Chuyen, T. Li, A.R. Hoffman, M. Blick, F. Fornari, N. Zanesi, H.
959 Alder, G. D'Elia, L. Gramantieri, L. Bolondi, G. Lanza, P. Querzoli, A. Angioni, C.M.
960 Croce, M. Negrini, Loss of imprinting of IGF2 sense and antisense transcripts in Wilms'
961 tumor, *Cancer Res.* 63 (2003) 1900–1905.

962 [95] D. Ferland-McCollough, D.S. Fernandez-Twinn, I.G. Cannell, H. David, M. Warner,
963 A.A. Vaag, J. Bork-Jensen, C. Brøns, T.W. Gant, A.E. Willis, K. Siddle, M. Bushell, S.E.
964 Ozanne, Programming of adipose tissue miR-483-3p and GDF-3 expression by maternal
965 diet in type 2 diabetes, *Cell Death Differ.* 19 (2012) 1003–1012.

966 [96] F. Pepe, S. Pagotto, S. Soliman, C. Rossi, P. Lanuti, C. Braconi, R.
967 MarianiCostantini, R. Visone, A. Veronese, Regulation of miR-483-3p by the O-linked
968 N-acetylglucosamine transferase links chemosensitivity to glucose metabolism in liver
969 cancer cells, *Oncogenesis.* 6 (2017) e328.

970 [97] M. Liu, A. Roth, M. Yu, R. Morris, F. Bersani, M.N. Rivera, J. Lu, T. Shioda, S.
971 Vasudevan, S. Ramaswamy, S. Maheswaran, S. Diederichs, D.A. Haber, The IGF2
972 intronic miR-483 selectively enhances transcription from IGF2 fetal promoters and
973 enhances tumorigenesis, *Genes Dev.* 27 (2013) 2543–2548.

974 [98] A. Veronese, L. Lupini, J. Consiglio, R. Visone, M. Ferracin, F. Fornari, N. Zanesi,
975 H. Alder, G. D'Elia, L. Gramantieri, L. Bolondi, G. Lanza, P. Querzoli, A. Angioni, C.M.
976 Croce, M. Negrini, Oncogenic role of miR-483-3p at the IGF2/483 locus, *Cancer Res.* 70
977 (2010) 3140–3149.

978 [99] H. Zhang, J. Hu, L. Liu, MiR-200a modulates TGF- β 1-induced endothelial-to-
979 mesenchymal shift via suppression of GRB2 in HAECs, *Biomed. Pharmacother.* 95
980 (2017) 215–222.

981 [100] K. Wang, B. Long, J.-Q. Jiao, J.-X. Wang, J.-P. Liu, Q. Li, P.-F. Li, miR-484
982 regulates mitochondrial network through targeting Fis1, Nat. Commun. 3 (2012) 781.

983

984

985

986

987

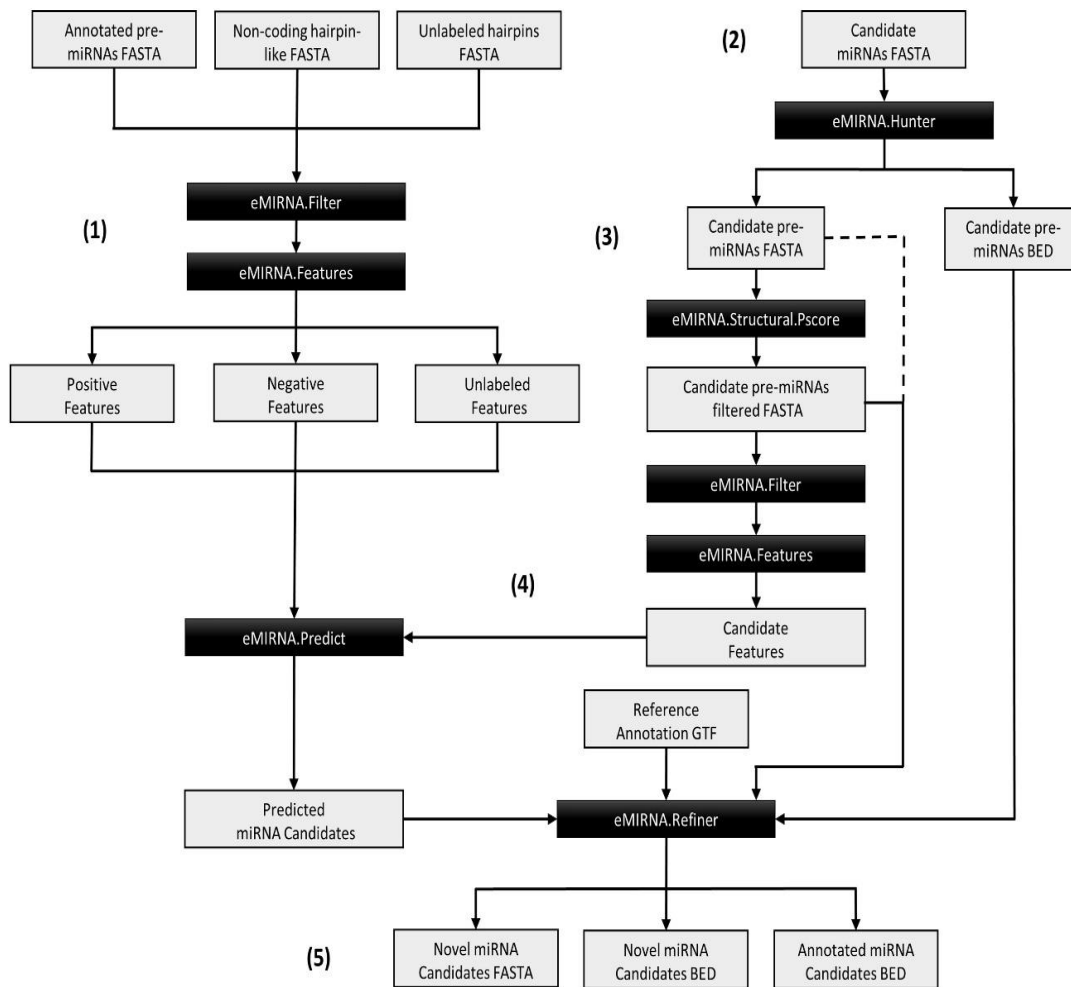
988

989

990

991 **Figures**

992



993

994

995 **Figure 1:** eMIRNA pipeline scheme for homology-based miRNA prediction using data
 996 from closely related species and *de novo* miRNA prediction from small RNA-seq data.

997 (1) Positive, negative and unlabeled data are filtered based on size and secondary folding

998 structure and a set of features is extracted for each sequence. (2) Mature miRNA

999 sequences from small RNA-seq data or related model species are mapped against the

1000 selected genome assembly and elongated to reconstruct putative pre-miRNA candidates.

1001 (3) Candidate precursors are filtered based on size and secondary folding structure and a

1002 set of features is extracted for each candidate sequence. Optionally, sequences showing

1003 unstable secondary structure are removed. (4) Candidate sequences are embedded in the

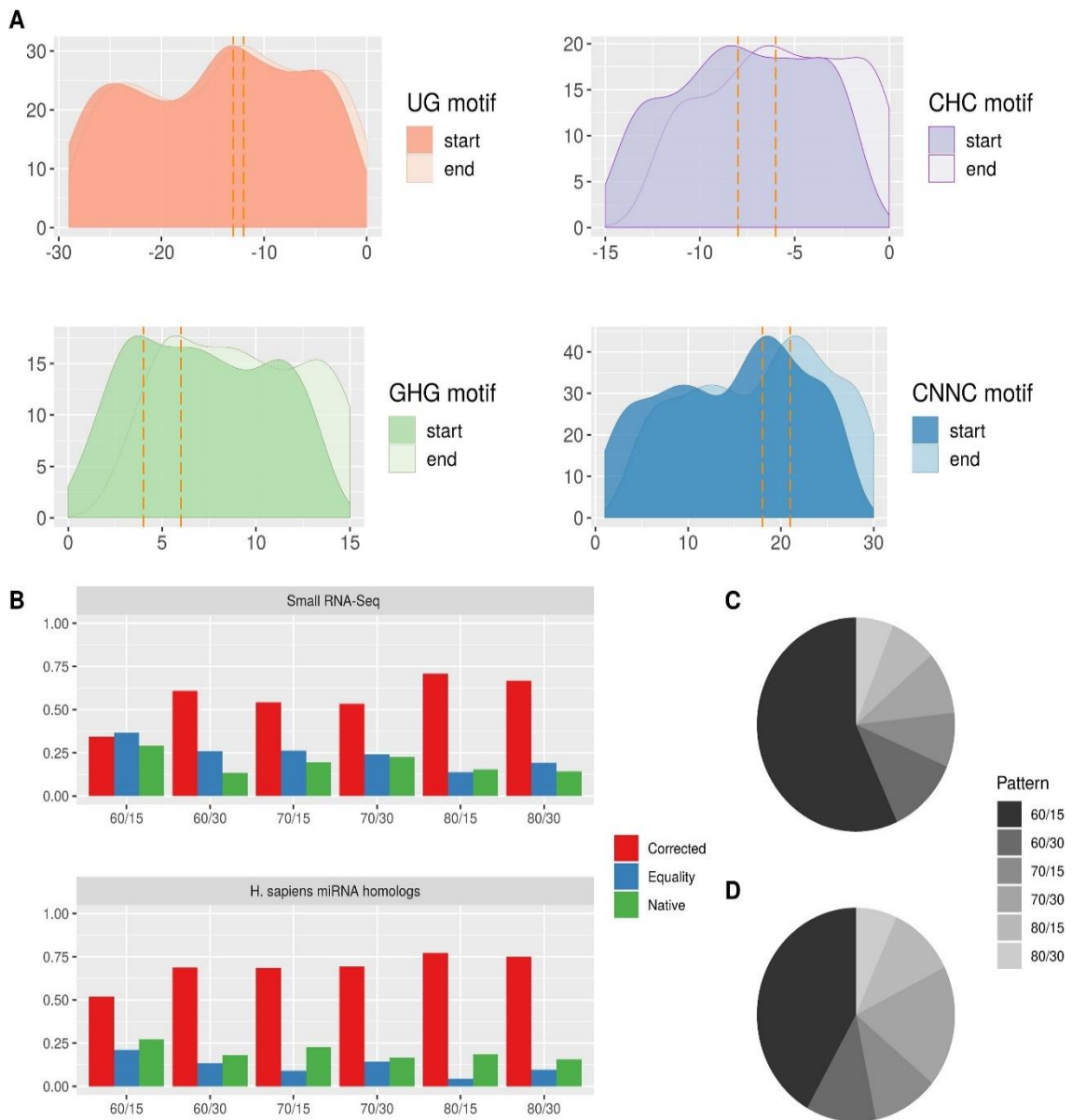
1004 semi-supervised transductive classifier and a list of putative miRNAs is predicted. (5)

1005 Predicted miRNAs are either assigned to already annotated miRNA loci in the provided
 1006 reference assembly or classified as putative novel miRNAs genes.

1007

1008

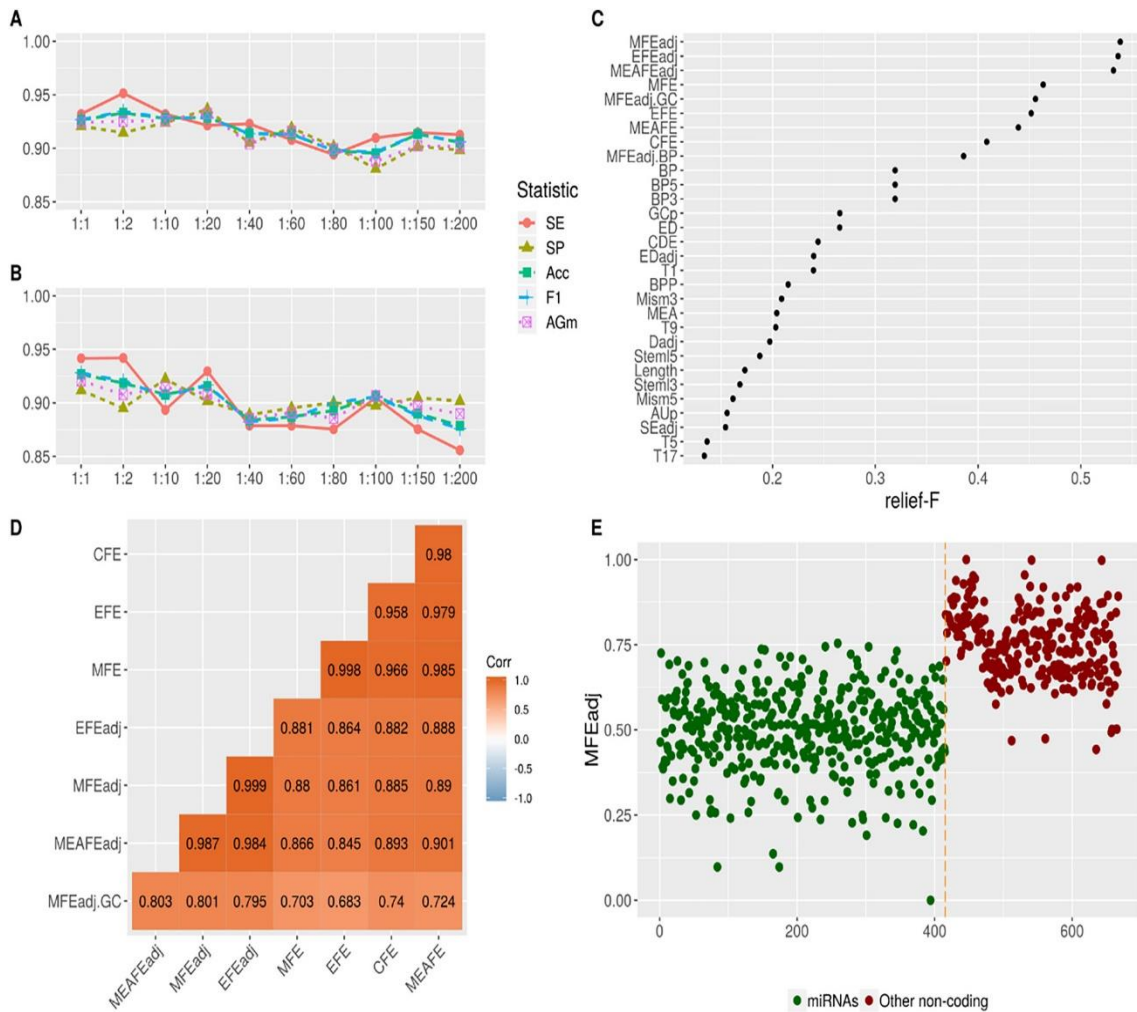
1009



1010

1011

1012 **Figure 2:** Processing motifs distribution and structural stability metrics. **(A)** Positional
1013 distribution of upstream and downstream motifs across annotated pre-miRNA boundaries
1014 in the porcine genome. **(B)** Proportion of candidate sequences for each elongation pattern
1015 showing the most stable folding structure according to randfold p score. The proportion
1016 of sequences for which the structural stability was higher in motif corrected candidates
1017 or, conversely, in non-corrected (native) candidates are shown as red and green bars,
1018 respectively. The proportion of sequences for which the structural stability was equivalent
1019 between motif corrected and native candidates were labeled as equally stable (blue). **(C)**
1020 Proportion of selected pre-miRNA candidates detected in the porcine *gluteus medius*
1021 muscle small RNA-seq data and **(D)** Proportion of selected pre-miRNA candidates
1022 detected through a *H. sapiens* homology-based miRNA search strategy, according to the
1023 most structurally stable elongation pattern tested. If two or more pre-miRNA sequences
1024 showed equivalent stability, the shortest motif-corrected candidate was selected.
1025



1026

1027

1028 **Figure 3:** Classification performance and feature importance statistics. Performance
 1029 metrics for Sensitivity (SE), Specificity (SP), Accuracy (Acc), F1-score (F1) and
 1030 Adjusted Geometric-mean (Agm) across incremental imbalance-ratios by using positive
 1031 miRNAs from (A) Ensembl and (B) miRCarta databases. (C) Thirty most discriminant
 1032 features according to the relief-F algorithm. (D) Pearson's correlation coefficient among
 1033 the seven most discriminant features associated with secondary structure stability metrics.
 1034 (E) Comparison of the folding structure stability between annotated miRNAs and other
 1035 hairpin-like non-coding RNA sequences present in the porcine genome. Stability is
 1036 expressed as the scaled Minimum Free Energy of the folded hairpins adjusted by sequence
 1037 length (MFEadj).

1038 **Tables**1039 **Table 1:** List of calculated features extracted from candidate hairpins.

Sequence Features	Symbol	Number of variables
Triplet Elements by SVM-Triplet	T1 ... T32	32
Sequence Length	Length	1
G+C/Length	GC	1
A+U/G+C	AU.GCr	1
A, U, G, C/Length	Ar, Ur, Gr, Cr	4
Dinucleotide/Length	Aar, GGr, CCr ...	16
Secondary Structure metrics	Symbol	Number of variables
Hairpin loop Length	HI	1
5' and 3' Stems Length	Stem5, Stem3	2
Basepairs in Secondary Structure	BP	1
Matches in 5' and 3' Stems	BP5, BP3	2
Mismatches in 5' and 3' Stems	Mism5, Mism3	2
Bulges in 5' and 3' Stems	B5, B3	2
Bulges in 5' and 3' Stems of types 1 to 7 mismatch	BN1.5, BN1.3 ...	14
A-U, G-C and G-U basepairs	Aup, GCp, Gup	3
Structural Statistics	Symbol	Number of variables
Minimum Free Energy	MFE	1
Ensemble and Centroid Free Energy	EFE, CFE	2
Centroid Distance to Ensemble	CDE	1
Maximum Expected Accuracy	MEA, MEAFE	2
BP/Length	BPP	1
MFE Ensemble Frequency	Efreq	1

Ensemble Diversity	ED	1
MFE/Length, EFE/Length and CDE/Length	MFEadj, EFEadj, Dadj	3
Shannon Entropy/Length	Seadj	1
MFE-EFE/Length	DiffMFE.EFE	1
MFEadj/GC and MFEadj/BP	MFEadj.GC, MFEadj.BP	2
MEAFE/Length and ED/Length	MEAFEadj, Edadj	2

1040

1041

1042 **Table 2:** Comparative benchmarking between the semi-supervised transductive learning
1043 approach employed by the *miRNAss* algorithm and other state-of-the-art supervised
1044 algorithms (i.e. SVM: Support Vector Machine, RF: Random Forest, KNN: k-Nearest
1045 Neighbors, NB: Naïve Bayes, XGB: Extreme Gradient Boosting and IGBM: light
1046 Gradient Boosting Tree) for miRNA classification. Only labeled positive and negative
1047 data sets were used for training.

1048 SE: Sensitivity; SP: Specificity; F-1: F-score measure of the harmonic mean of the
1049 precision and recall; AUROC: Area under the Receiver Operating Characteristics (ROC)
1050 curve; AUPR: Area under the Precision-Recall curve.

Statistic	SVM	RF	KNN	NB	XGB	IGBM	miRNAss
SE	0.932	0.932	0.9223	0.9126	0.9515	0.9223	0.8835
SP	0.8413	0.9524	0.9524	0.9683	0.9365	0.9048	0.9683
F-1	0.9187	0.9505	0.9453	0.9447	0.9561	0.9314	0.9226
AUROC	0.6428	0.7246	0.5757	0.4291	0.7063	0.9781	0.9783
AUPR	0.7222	0.8489	0.6751	0.5818	0.8509	0.9873	0.987

1051

1052

1053 **Table 3:** Novel porcine miRNA genes predicted through a homology-based comparison
 1054 with human miRNA annotation and on the basis of data generated by sequencing small
 1055 RNAs expressed in the *gluteus medius* muscle of Duroc pigs.

1056 Chr: Chromosome; N: Neighborhood score.

Chr	Start	End	Strand	ID	N
1	191218572	191218651	+	miR-3529	0.33
1	268816970	268817050	+	miR-219b	0.92
2	32718	32792	+	miR-6743	0.82
2	1473428	1473495	-	miR-483	0.84
2	1474436	1474513	-	3229-4643	-
2	40104336	40104403	-	1325-14520	-
2	134660802	134660897	-	1323-14559	-
3	7180536	7180603	-	miR-484	0.1
3	40421320	40421409	+	427-63874	-
3	40772345	40772445	+	176-178526	-
4	22195784	22195880	+	2340-6855	-
5	3397056	3397130	-	1111-18619	-
5	17410008	17410122	+	1794-9841	-
5	95548384	95548458	+	miR-3059	1
6	56426941	564267012	-	miR-520e	0.3
6	63490755	63490822	+	miR-200a	0.6
8	1205684	1205760	-	miR-4800	0.85
9	52087075	52087155	+	1864-9314	-
9	114528009	114528076	+	miR-3120	0.7
10	27079413	27079489	-	miR-24-1	0.79
11	1824995	1825062	+	504-51258	-
11	49808356	49808431	-	miR-3665	0.86
12	1538011	1538119	+	337-84973	-

12	1601453	1601506	-	miR-3065	0.82
12	18989584	18989651	+	399-69074	-
12	45088806	45088863	+	miR-451b	0.78
12	45597382	45597459	+	miR-4523	0.81
12	46211527	46211594	-	miR-3184	0.61
12	48162620	48162704	-	miR-132	0.84
12	56201226	56201300	-	518-49963	-
13	30242047	30242114	+	772-29980	-
13	33152284	33152383	+	miR-4787	0.83
13	197168804	197168901	+	miR-6501	0.97
14	87673881	87673954	+	3552-4147	-
14	109233945	109234032	-	miR-3085	0.95
14	122706280	122706361	+	miR-6715a	0.96
14	122706285	122706353	-	miR-6715b	0.96
14	127016706	127016794	-	miR-9851	0.83
14	140979533	140979627	+	3525-4198	-
15	128165751	128165827	-	miR-5702	0.86
17	61915309	61915376	+	1544-12001	-
X	41793240	41793315	+	451-58980	-
X	43716471	43716538	+	miR-502	0.73
X	59551153	59551220	+	miR-374c	0.8
X	94122543	94122610	+	miR-1264	0.83
X	96979691	96979765	+	miR-1277	0.68
X	124724889	124724956	-	miR-718	0.89

1057

1058

1059

1060

1061 **Supplementary Materials**

1062 **Figure S1:** (A) Receiver Operating Characteristics (ROC) and (B) Precision-Recall (PR)
1063 curves computed for each pre-defined imbalance scenario using porcine Ensembl
1064 annotation for positive (miRNAs) and negative (other hairpin-like non-coding RNAs)
1065 data sets.

1066 **Figure S2:** RT-qPCR results of selected novel miRNAs. Successfully profiled novel
1067 miRNAs in (A) the *longissimus dorsi* skeletal muscle and (B) liver tissues from 7
1068 Göttingen minipigs.

1069 **File S1:** FASTA file of collapsed expressed sequences (RPM > 10) used in the *de novo*
1070 discovery of miRNAs expressed in the porcine *gluteus medius* skeletal muscle.

1071 **File S2:** Non-redundant annotated mature miRNA sequences obtained from the *H.*
1072 *sapiens* GRCh38.p12 genome assembly used as a reference in the homology-based search
1073 of novel miRNAs in the current release of the porcine genome (Sscrofa11.1).

1074 **File S3:** List of already annotated miRNAs and all isomiRs detected as expressed (RPM
1075 > 10) in the porcine *gluteus medius* skeletal muscle.

1076 **File S4:** Amplification profiles and melting curves for the three novel miRNA candidates
1077 subjected to confirmation by RT-qPCR analyses.

1078 **Table S1:** Area under the curve (AUC) computed for each pre-defined imbalance
1079 scenario using Ensembl annotation for positive and negative data sets.

1080 **Table S2:** True positive ratio of porcine miRNA loci annotated in the Ensembl and
1081 miRCarta databases and identified by the eMIRNA pipeline in all considered imbalance
1082 scenarios.

1083 **Table S3:** Mature miRNAs and primers used for RT-qPCR confirmation of selected
1084 novel miRNA candidates.

1085 **Table S4:** Feature importance according to the relief-F algorithm.

1086 **Table S5:** Previously annotated miRNAs genes that are correctly classified as miRNAs
1087 by eMIRNA and miRDeep2.

1088 **Table S6:** miRDeep2 algorithm results for miRNA prediction using the *gluteus medius*
1089 muscle small RNA-seq data generated in the present study.

1090 **Table S7:** Deviation rates (dr) and Differential deviation (ΔD) estimates for miRNA
1091 genomic positional prediction with eMIRNA and miRDeep2.

1092