

1 MirGeneDB2.0: the curated microRNA Gene Database

2 Non-coding RNAs (ncRNA), a significant part of the increasingly popular ‘*dark matter*’ of  
3 the human genome<sup>1</sup>, have gained substantial attention due to their involvement in animal  
4 development and human disorders such as cardiovascular diseases and cancer<sup>2</sup>. Although  
5 many different types of regulatory ncRNAs have been discovered over the last 25 years,  
6 microRNAs (miRNAs) are unique within these as they are the only class of ncRNAs with  
7 individual genes sequentially conserved across the animal kingdom<sup>3</sup>. Because of the  
8 conserved roles miRNAs play in establishing robustness of gene regulatory networks across  
9 Metazoa<sup>4</sup>, it is important that homologous miRNAs in different species are correctly  
10 identified, annotated, and named using consistent criteria<sup>5</sup> against the backdrop of numerous  
11 other types of coding and non-coding RNA fragments<sup>6</sup>.

12 Unlike miRBase<sup>7</sup>, which has developed organically through community-wide submissions,  
13 and thus does not use consistent annotation or nomenclature criteria<sup>6</sup>, MirGeneDB2.0  
14 (<http://mirgenedb.org>), a manually curated open source miRNA gene database, contains high  
15 quality annotations of 7,785 *bona fide* and consistently named miRNAs from 32 species  
16 representing major metazoan groups (including many invertebrate and vertebrate model  
17 organisms). The number of miRNAs conforming to the annotation criteria is almost four  
18 times higher than in miRBase (~2000 for the miRBase ‘high confidence’ set<sup>7</sup>), and can be  
19 considered free of *false positives*. For the expansion of the previous version, we used more  
20 than 250 publicly available sequencing datasets (for a total of 4.2 billion reads) derived from  
21 at least one representative dataset for each organism (such as whole organisms, organs,  
22 tissues or cell-types), which allowed for a consistent and uniform annotation of  
23 microRNAomes for each species (Supplementary File, “file\_info”; Supplementary  
24 Methods)<sup>8</sup>. Existing MirGeneDB.org miRNA complements for human, mouse, chicken and  
25 zebrafish were expanded from our initial effort by 65, 49, 28 and 100 genes, respectively

(Supplementary File, table), and annotation-accuracy was further improved using available Cap Analysis of Gene Expression (CAGE) data when available (Supplementary File, “CAGE”)<sup>9</sup>.

Because miRBase has become increasingly heterogeneous with respect to the number of *bona fide* miRNAs relative to other types of non-coding RNAs, it has considerable variation in the number of miRNAs for closely related groups (Supplementary File, graph miRBase). However, in MirGeneDB, congruent miRNA complements in terms of total miRNA genes and miRNA families were observed in related groups, such as the Vertebrates and arthropods<sup>3,10</sup> (Figure 1). Big differences between miRBase and MirGeneDB2.0 can be observed because miRBase has on the one hand a much larger number of annotated sequences for some of the 23 taxa shared with MirGeneDB2.0 including human, mouse, and chicken, accounting for 4,243 *false positives*, and on the other hand it lacks 22% of all MirGeneDB2.0 genes, accounting for 1,180 *false negatives* (Figure 2, Supplementary File, “overview”). Finally, 31% of the remaining 4,275 miRNAs are incompletely annotated in miRBase, whereas in MirGeneDB2.0 each miRNA has both arms annotated, with a clear distinction made between sequenced reads and predicted reads for each miRNA entry with predictions derived from both considerations of secondary structure and expressed orthologues in other taxa.

The expanded web-interface of MirGeneDB2.0 allows browsing, searching and downloading of miRNA-complements for each organism. Annotations are downloadable as fasta, gff, or bed-files containing distinct sub-annotations for all miRNA components such as precursor (pre), mature, loop, co-mature or star sequences. Unlike miRBase, seed sequences are also identified, and can be searched independently from the rest of the mature sequence. In addition, we included 30-nucleotide flanking regions on both arms for each precursor transcript to generate an extended precursor transcript, which again is downloadable.

51 MirGeneDB2.0 employs an internally consistent nomenclature system where genes of  
52 common descent are assigned the same miRNA family name, allowing for the easy  
53 recognition of both orthologues in other species, and paralogues within the same species.  
54 This nomenclature system allows for an accurate reconstruction of ancestral miRNA  
55 repertoires – both at the family level and at the gene level – that is now provided in  
56 MirGeneDB2.0 for all nodes leading to the 32 terminal taxa considered, which allows users  
57 to easily assess both gains and losses of miRNA genes through time. However, in order to not  
58 increase confusion about the naming of miRNA genes, we continue to provide commonly  
59 used miRBase names – if available – in our “*Browse*” section of MirGeneDB2.0 (i.e.  
60 <http://mirgenedb.org/browse/hsa>).

61 *Gene-pages* for each miRNA gene contain names, orthologues & paralogues, downloadable  
62 sequences, structure, and a range of other previously available information including genomic  
63 coordinates (i.e. <http://mirgenedb.org/show/hsa/Let-7-P1>). New features in MirGeneDB2.0  
64 include accurate information on 3’ non-templated uridylation, which characterize an  
65 important sub-group of miRNAs; information of the presence or absence of the recently  
66 discovered sequential motifs (UG, UGUG, CNNC); and the visualization of at least one  
67 expression dataset for each gene in each organism. Further, *read-pages* are also provided for  
68 each gene (i.e. <http://mirgenedb.org/static/graph/hsa/results/Hsa-Let-7-P1.html>), which show  
69 an overview of read-stacks on the corresponding extended precursor sequence of each *gene-*  
70 *page*. They contain detailed representation of templated and non-templated reads for  
71 individual datasets for each gene including reports on miRNA isoforms and downloadable  
72 read-mappings.

73 The establishment of this carefully curated data base of miRNA genes, supplementing  
74 existing databases including miRBase, allows for a stable and robust foundation for miRNA  
75 studies, in particular studies that rely on cross-species comparisons to explore the roles

76 miRNAs play in development and disease, as well as the evolution of miRNAs (and animals)  
77 themselves.

78 Note: Supplementary Methods and Supplementary files are available in the online version of  
79 the paper.

## 80 ACKNOWLEDGMENTS

81 We thank Victor Ambros, David Bartel, Marc Friedländer, Marc Halushka, Andreas Keller,  
82 Gianvito Urgese for discussions, Georgios Magklaras for IT support. B.F. has been supported  
83 by the South-Eastern Norway Regional Health Authority (Grant No. 2014041). AM has been  
84 supported by the Norwegian Research Council, Helse Sør-Øst, and the University of Oslo  
85 through the Centre for Molecular Medicine Norway (NCMM), which is part of the Nordic  
86 European Molecular Biology Laboratory partnership for Molecular Medicine. K.J.P. was  
87 supported by NASA-Ames.

## 88 AUTHOR CONTRIBUTIONS

89 BF and KJP conceived MirGeneDB2.0, compiled miRNA complements for all organisms.  
90 DD created read-pages and heatmaps. MJ set up the framework and database. MH processed  
91 sRNA sequencing data, AM processed and analyzed CAGE data. EHøyve created scripts for  
92 mature/star annotation. EHovig and KF provided infrastructure and all authors read and  
93 commented on the manuscript.

## 94 COMPETING FINANCIAL INTERESTS

95 The authors declare no competing financial interests.

96 Bastian Fromm 1, Diana Domanska 2, Michael Hackenberg 3, Anthony Mathelier 4,5, Eirik  
97 Høyve 1, Morten Johansen 6, Eivind Hovig 1,2,6, Kjersti Flatmark 1,7,8 & Kevin J. Peterson  
98 9.

99 1 - Department of Tumor Biology, Institute for Cancer Research, The Norwegian Radium  
100 Hospital, Oslo University Hospital, Nydalen, N-0424 Oslo, Norway, 2 - Department of  
101 Informatics, University of Oslo, Blindern, N-0318 Oslo, Norway, 3 - Department of Genetics,  
102 Faculty of Sciences, University of Granada, Granada, 1s8071, Spain, 4 - Centre for Molecular

103 Medicine Norway (NCMM), Nordic EMBL Partnership, University of Oslo, N-0318 Oslo,  
 104 Norway, 5 - Department of Cancer Genetics, Institute for Cancer Research, The Norwegian  
 105 Radium Hospital, Oslo University Hospital Radiumhospitalet, N-0424 Oslo, Norway, 6 -  
 106 Institute for Medical Informatics, The Norwegian Radium Hospital, Oslo University  
 107 Hospital, N-0424 Oslo, Norway, 7 - Department of Gastroenterological Surgery, The  
 108 Norwegian Radium Hospital, Oslo University Hospital, Nydalen, N-0424 Oslo, Norway, 8 -  
 109 Institute of Clinical Medicine, University of Oslo, Blindern, N-0318 Oslo, Norway, 9 -  
 110 Department of Biological Sciences, Dartmouth College, Hanover, New Hampshire 03755,  
 111 U.S..  
 112 e-mail: BastianFromm@gmail.com or Kevin.J.Peterson@dartmouth.edu

## 113 References

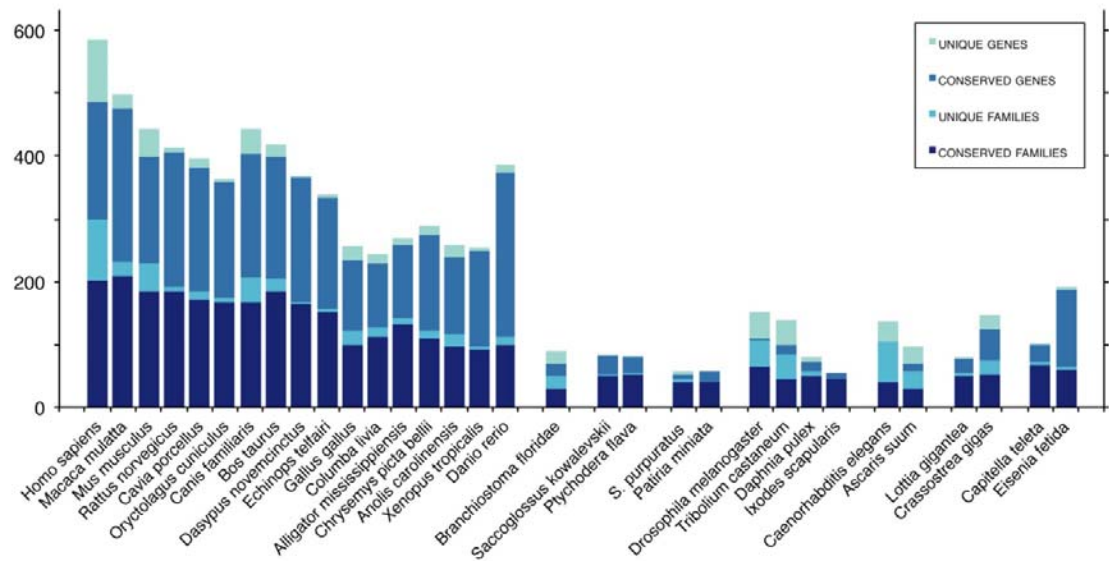
- 114 1 Blaxter, M. Genetics. Revealing the dark matter of the genome. *Science* **330**, 1758-1759,  
115 doi:10.1126/science.1200700 (2010).
- 116 2 Esteller, M. Non-coding RNAs in human disease. *Nature reviews. Genetics* **12**, 861-874,  
117 doi:10.1038/nrg3074 (2011).
- 118 3 Wheeler, B. *et al.* The deep evolution of metazoan microRNAs. *Evolution & development* **11**,  
119 50 - 68 (2009).
- 120 4 Ebert, M. S. & Sharp, P. A. Roles for microRNAs in conferring robustness to biological  
121 processes. *Cell* **149**, 515-524, doi:10.1016/j.cell.2012.04.005 (2012).
- 122 5 Ambros, V. A uniform system for microRNA annotation. *Rna* **9**, 277-279,  
123 doi:10.1261/rna.2183803 (2003).
- 124 6 Tosar, J. P., Rovira, C. & Cayota, A. Non-coding RNA fragments account for the majority of  
125 annotated piRNAs expressed in somatic non-gonadal tissues. *Communications Biology* **1**, 2,  
126 doi:10.1038/s42003-017-0001-7 (2018).
- 127 7 Kozomara, A. & Griffiths-Jones, S. miRBase: annotating high confidence microRNAs using  
128 deep sequencing data. *Nucleic acids research* **42**, D68-73, doi:10.1093/nar/gkt1181 (2014).
- 129 8 Fromm, B. *et al.* A Uniform System for the Annotation of Vertebrate microRNA Genes and  
130 the Evolution of the Human microRNAome. *Annual review of genetics* **49**, 213-242,  
131 doi:10.1146/annurev-genet-120213-092023 (2015).
- 132 9 de Rie, D. *et al.* An integrated expression atlas of miRNAs and their promoters in human and  
133 mouse. *Nature biotechnology* **35**, 872-878, doi:10.1038/nbt.3947 (2017).
- 134 10 Tarver, J. E. *et al.* miRNAs: small genes with big potential in metazoan phylogenetics.  
135 *Molecular biology and evolution* **30**, 2369-2382, doi:10.1093/molbev/mst133 (2013).

136

137

138    Figures

139

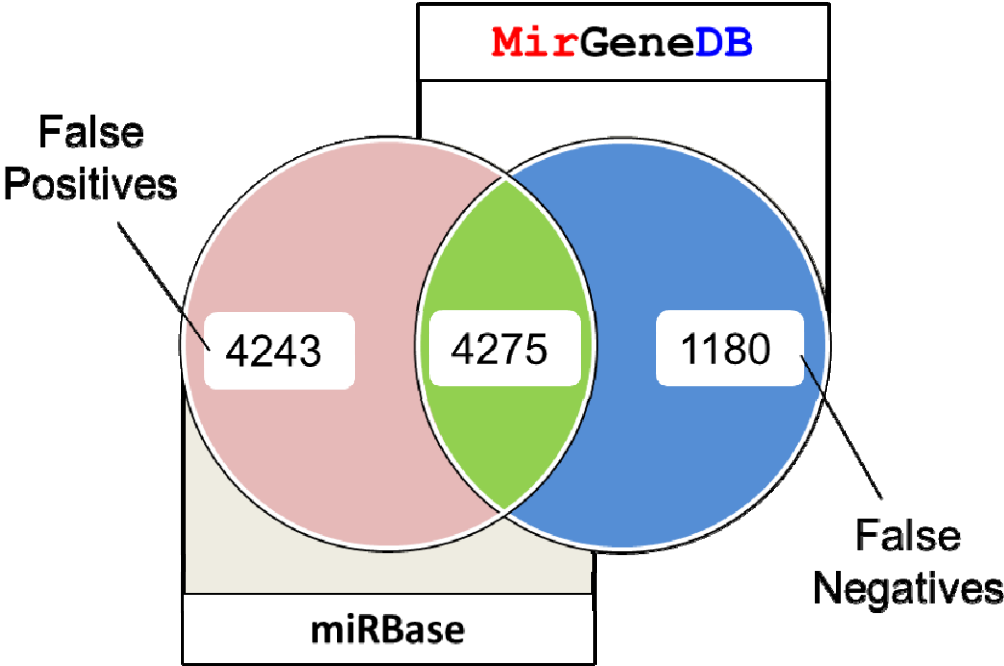


140

141    Figure 1: High consistency of conserved miRNA gene and family numbers in closely related groups in  
142    MirGeneDB2.0 can be observed for groups with more than two representatives. High variation in  
143    gene-numbers for *Danio* and *Eisenia* (double asterisks) are explainable by genome-duplication events  
144    within that particular monophyletic group (vertebrates and annelids, respectively), while high  
145    numbers of unique /novel genes and families in *Homo*, *Mus*, *Canis*, *Drosophila*, *Tribolium* and  
146    *Caenorhabditis* might be explainable by the significantly higher number of studies and/or the  
147    relatively higher number of absolute small RNA reads on these organisms (single asterisks).



148



149

150 Figure 2: High number of incorrect and missing miRNA annotations in miRBase as compared to MirGeneDB.  
151 A comparison of the microRNA complements of 23 organisms shared between miRBase and MirGeneDB  
152 revealed that only 4,275 of the 8,531 entries in miRBase are shared with MirGeneDB (green). An additional  
153 4,243 miRBase entries represent false positives (red), miRNAs found in miRBase that do not satisfy standard  
154 annotation criteria, whereas 1,180 MirGeneDB entries represent false negatives (blue), miRNAs that are present  
155 in these taxa that are not currently annotated in miRBase.

## Supplementary Methods

### *Features of miRNAs*

In the last two decades the small non-coding RNA field has significantly expanded beyond snRNAs & snoRNAs<sup>1</sup> to include piRNAs<sup>2</sup>, siRNAs<sup>3</sup>, novel small RNAs derived from known non-coding RNAs including tRNAs<sup>4</sup> and rRNAs<sup>5</sup> and, of course, microRNAs (miRNAs)<sup>6-9</sup>.

Each of these types of smallRNAs is characterized by a distinctive suite of characteristics and a unique evolutionary history. MiRNAs can be distinguished from other small genomically encoded RNA families by a set of unique features as described earlier<sup>10,11</sup>: The presence of two 20-26nt long reads that are expressed from each of the two arms derived from a stable hairpin precursor is essential to assess whether or not Drosha and Dicer were involved in the processing. Since the ends of canonical miRNA reads are generated enzymatically, the 5' ends of the reads are homogeneous (>90%). The hairpin precursor shows imperfect complementarity and base pairs in at least 16 of the ~ 22 nucleotides. The 5p and 3p reads are offset by 2 nucleotides on both ends due to the sequential processing of the miRNA transcript by Drosha and Dicer to generate the mature ~22 nucleotide read(s). In some cases, the Drosha offset is only offset by 1 templated nucleotide, but in these cases the 3' end of the 3p arm is monouridylated<sup>12,13</sup>. The length of the loop is at least 8 nucleotides long; there is no apparent maximum in loop length, even in organisms possessing only a single Dicer gene, contra our earlier statement<sup>11</sup>, even though most taxa like vertebrates with single Dicer genes never show loop lengths greater than ~40 nucleotides.

There are other features of miRNAs, in particular structural and evolutionary signatures that allow them to be further distinguished from other small RNAs. The mature miRNA sequence usually starts with A or U, and is often mismatched with the complementary arm, which seems to facilitate arm selection by Argonaute (at least in mammals)<sup>11,14,15</sup>. Nucleotide

positions 2 through 8 of the mature sequence (the "seed") are strongly conserved through evolution, as are positions 13-16 (the 3' complementary region)<sup>11,16</sup>. Recently it was demonstrated that processing motifs are often (but not always) present in the primary miRNA transcript including a UG motif 14 nucleotides upstream of the 5p arm, a UGU motif at the 3' end of the 5p arm, and a CNNC motif 17 nucleotides downstream of the 3p arm<sup>17-19</sup>.

Recognition and utilization of clear and, for the most part mechanistically well understood, criteria for the annotation of miRNAs allows the delineation of *bona fide* miRNAs from the myriad small RNAs generated in eukaryotic cells, providing deeper and more significant insights into their function, possible mis-regulation, and evolution.

### *Data processing*

Publicly available smallRNA sequencing data of whole organisms, healthy organs, tissue or cell-isolates was downloaded from European Genome-phenome Archive (EGA), the Sequence Read Archive (SRA) and the the Gene Expression Omnibus (GEO) respectively (see Supplementary File, "file\_info"). For download and processing we used the latest version of sRNAbench<sup>20</sup>. Corresponding files were automatically downloaded and converted into fastq files. All datasets were consistently processed with the following parameters: 3' adapter sequences were automatically identified and trimmed using sRNAbench (detection of at least 10nt of the adapter allowing 1 mismatch)<sup>20</sup>; reads within length of 18 and 27 nts were retained and collapsed for mapping employing fastx-toolkit, and custom perl scripts. Collapsed reads were mapped to miRBase complements and MirGeneDB<sup>11</sup> using bowtie1.2<sup>21</sup>, requiring an 18 nucleotide seed sequence of zero mismatches to avoid cross-mapping. All mappings were transformed to bam-files using SAMtools<sup>22</sup>.

## 203 *Confident miRNA complement annotations for 32 metazoan taxa*

204 Similar to our previous efforts for the four vertebrates of MirGeneDB1.0<sup>11</sup>, we analyzed all  
 205 miRBase “miRNA” entries for the 23 available taxa. We included, however, at least one  
 206 smallRNAseq dataset to assess the status of the unique miRNA features for each miRNA  
 207 individually. For this task, we used an improved version of MirMiner to visualize read-  
 208 mappings and structures, which also allowed us to predict previously missing genes and the  
 209 miRNA complements for the nine organisms which are not found in miRBase<sup>16</sup>.

210 Because these criteria can only be applied in miRNAs processed by the canonical pathway  
 211 (i.e. they are processed by Drosha and Dicer respectively), and literally no non-canonical  
 212 miRNA is conserved beyond very close relatives, we have only considered canonical  
 213 miRNAs. One exception we made was the non-canonical erythroid miRNA Mir-451 that is a  
 214 very important regulator of erythroid development and highly conserved in all vertebrates<sup>23</sup>  
 215 (<http://mirgenedb.org/browse/ALL?family=MIR-451>).

## 216 *Mature arm annotations*

217 Mature and star, and Co-mature status of miRNAs was assigned by assessing the expression  
 218 of 5’ and 3’ arms overall available datasets, respectively. Only if one arm was expressed  
 219 more than twofold higher as the other mature status was assigned, else Co-mature status was  
 220 given (Supplementary markdown). In the few cases where arms were not both expressed we  
 221 used information from orthologous genes in related organisms and assigned predicted status  
 222 of mature /star based on the expression ratios in the corresponding datasets of the related  
 223 species.

## 224 *Refinement of pre-miRNA 3’end annotation with CAGE data*

## 225 *Human annotation*

We downloaded the hg38 bigwig files associated to all ENCODE CAGE experiments from the ENCODE data portal (see [https://www.encodeproject.org/metadata?type=Experiment&assay\\_slims=Transcription&assay\\_title=CAGE&assembly=GRCh38&files.file\\_type=bigWig/metadata.tsv](https://www.encodeproject.org/metadata?type=Experiment&assay_slims=Transcription&assay_title=CAGE&assembly=GRCh38&files.file_type=bigWig/metadata.tsv)). We merged the data from all the experiments and converted the files in the BED format. Computation and plotting of the distribution of CAGE tags around the 3' end of pre-miRNAs annotated in MirGeneDB were performed using the deepTools v. 2.4.0<sup>24,25</sup> (Supplementary Figure 1a). As described in previous studies<sup>26,27</sup> we observed a peak for CAGE tags 1 nt downstream (i.e. the +1 nt) of the pre-miRNA 3' ends (Supplementary Figure 1). We considered for manual curation the pre-miRNAs showing a higher number of CAGE tags at positions 0 or +2 with respect to the annotated pre-miRNA 3' end, which could correspond to a 1 nt off misannotation (see [http://fantom.gsc.riken.jp/zenbu/gLyphs/#config=ufw7Z\\_rvFF5juG\\_FZhbOD](http://fantom.gsc.riken.jp/zenbu/gLyphs/#config=ufw7Z_rvFF5juG_FZhbOD) for an example). After manual curation through the Zenbu genome browser<sup>28</sup>, we corrected the 3' end position for pre-miRNAs Hsa-Mir-145, Hsa-Let-7-P12, Hsa-Let-7-P7 (Supplementary File, "CAGE").

#### Zebrafish annotation

We applied the same methodology to the CAGE data obtained from 12 developmental stages of embryogenesis in zebrafish<sup>27</sup>. Bigwig files of CAGE tags mapping were retrieved using the CAGEr R package<sup>29</sup>. Data from all developmental stages were merged to analyze the distribution of CAGE tags around pre-miRNA 3' ends (Supplementary Figure 1b). After manual curation, we updated the 3' end position of the pre-miRNAs Dre-Mir-153-P1a and Dre-Let-7-P6 (Supplementary File, "CAGE").

#### Water flea annotation

249 CAGE data for *Daphnia pulex* derived from three developmental states were retrieved from  
250 GEO (GSE80141)<sup>30</sup>. We followed the same steps as described above for human and zebrafish  
251 CAGE data but did not find any 3' end annotation of pre-miRNAs to update.

## 252     *Supplementary References*

- 253     1     Matera, A. G., Terns, R. M. & Terns, M. P. Non-coding RNAs: lessons from the small nuclear  
254           and small nucleolar RNAs. *Nature reviews. Molecular cell biology* **8**, 209-220,  
255           doi:10.1038/nrm2124 (2007).
- 256     2     Lau, N. C. *et al.* Characterization of the piRNA complex from rat testes. *Science* **313**, 363-367,  
257           doi:10.1126/science.1130164 (2006).
- 258     3     Hamilton, A. J. & Baulcombe, D. C. A species of small antisense RNA in posttranscriptional  
259           gene silencing in plants. *Science* **286**, 950-952 (1999).
- 260     4     Goodarzi, H. *et al.* Endogenous tRNA-Derived Fragments Suppress Breast Cancer Progression  
261           via YBX1 Displacement. *Cell* **161**, 790-802, doi:10.1016/j.cell.2015.02.053 (2015).
- 262     5     Chak, L. L., Mohammed, J., Lai, E. C., Tucker-Kellogg, G. & Okamura, K. A deeply conserved,  
263           noncanonical miRNA hosted by ribosomal DNA. *Rna* **21**, 375-384,  
264           doi:10.1261/rna.049098.114 (2015).
- 265     6     Lee, R. C. & Ambros, V. An extensive class of small RNAs in *Caenorhabditis elegans*. *Science*  
266           **294**, 862-864 (2001).
- 267     7     Lee, R. C., Feinbaum, R. L. & Ambros, V. The *C. elegans* heterochronic gene *lin-4* encodes  
268           small RNAs with antisense complementarity to *lin-14*. *Cell* **75**, 843-854 (1993).
- 269     8     Lau, N. C., Lim, L. P., Weinstein, E. G. & Bartel, D. P. An abundant class of tiny RNAs with  
270           probable regulatory roles in *Caenorhabditis elegans*. *Science* **294**, 858-862 (2001).
- 271     9     Lagos-Quintana, M., Rauhut, R., Lendeckel, W. & Tuschl, T. Identification of novel genes  
272           coding for small expressed RNAs. *Science* **294**, 853-858 (2001).
- 273     10    Ambros, V. A uniform system for microRNA annotation. *Rna* **9**, 277-279,  
274           doi:10.1261/rna.2183803 (2003).
- 275     11    Fromm, B. *et al.* A Uniform System for the Annotation of Vertebrate microRNA Genes and  
276           the Evolution of the Human microRNAome. *Annual review of genetics* **49**, 213-242,  
277           doi:10.1146/annurev-genet-120213-092023 (2015).
- 278     12    Kim, B. *et al.* TUT7 controls the fate of precursor microRNAs by using three different  
279           uridylation mechanisms. *The EMBO journal* **34**, 1801-1815, doi:10.15252/embj.201590931  
280           (2015).
- 281     13    Kim, Y. K., Kim, B. & Kim, V. N. Re-evaluation of the roles of DROSHA, Exportin 5, and DICER  
282           in microRNA biogenesis. *Proceedings of the National Academy of Sciences of the United*  
283           *States of America* **113**, E1881-1889, doi:10.1073/pnas.1602532113 (2016).
- 284     14    Suzuki, H. I. *et al.* Small-RNA asymmetry is directly driven by mammalian Argonautes. *Nature*  
285           *structural & molecular biology* **22**, 512-521, doi:10.1038/nsmb.3050 (2015).
- 286     15    Schirle, N. T., Sheu-Gruttadauria, J. & MacRae, I. J. Structural basis for microRNA targeting.  
287           *Science* **346**, 608-613, doi:10.1126/science.1258040 (2014).
- 288     16    Wheeler, B. *et al.* The deep evolution of metazoan microRNAs. *Evolution & development* **11**,  
289           50 - 68 (2009).
- 290     17    Nguyen, T. A. *et al.* Functional Anatomy of the Human Microprocessor. *Cell* **161**, 1374-1387,  
291           doi:10.1016/j.cell.2015.05.010 (2015).
- 292     18    Fang, W. & Bartel, D. P. The Menu of Features that Define Primary MicroRNAs and Enable De  
293           Novo Design of MicroRNA Genes. *Molecular cell* **60**, 131-145,  
294           doi:10.1016/j.molcel.2015.08.015 (2015).
- 295     19    Auyeung, V. C., Ulitsky, I., McGeary, S. E. & Bartel, D. P. Beyond secondary structure:  
296           primary-sequence determinants license pri-miRNA hairpins for processing. *Cell* **152**, 844-858,  
297           doi:10.1016/j.cell.2013.01.031 (2013).
- 298     20    Rueda, A. *et al.* sRNAtoolbox: an integrated collection of small RNA research tools. *Nucleic*  
299           *Acids Res* **43**, W467-473, doi:10.1093/nar/gkv555 (2015).

300 21 Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient  
301 alignment of short DNA sequences to the human genome. *Genome biology* **10**, R25,  
302 doi:10.1186/gb-2009-10-3-r25 (2009).

303 22 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-  
304 2079, doi:10.1093/bioinformatics/btp352 (2009).

305 23 Jee, D. *et al.* Dual Strategies for Argonaute2-Mediated Biogenesis of Erythroid miRNAs  
306 Underlie Conserved Requirements for Slicing in Mammals. *Molecular cell* **69**, 265-278 e266,  
307 doi:10.1016/j.molcel.2017.12.027 (2018).

308 24 Ramírez, F., Dündar, F., Diehl, S., Grüning, B. A. & Manke, T. deepTools: a flexible platform  
309 for exploring deep-sequencing data. *Nucleic Acids Res.* **42**, W187-191,  
310 doi:10.1093/nar/gku365 (2014).

311 25 Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data  
312 analysis. *Nucleic Acids Res.* **44**, W160-165, doi:10.1093/nar/gkw257 (2016).

313 26 de Rie, D. *et al.* An integrated expression atlas of miRNAs and their promoters in human and  
314 mouse. *Nat. Biotechnol.* **35**, 872-878, doi:10.1038/nbt.3947 (2017).

315 27 Nepal, C. *et al.* Transcriptional, post-transcriptional and chromatin-associated regulation of  
316 pri-miRNAs, pre-miRNAs and moRNAs. *Nucleic Acids Res.* **44**, 3070-3081,  
317 doi:10.1093/nar/gkv1354 (2016).

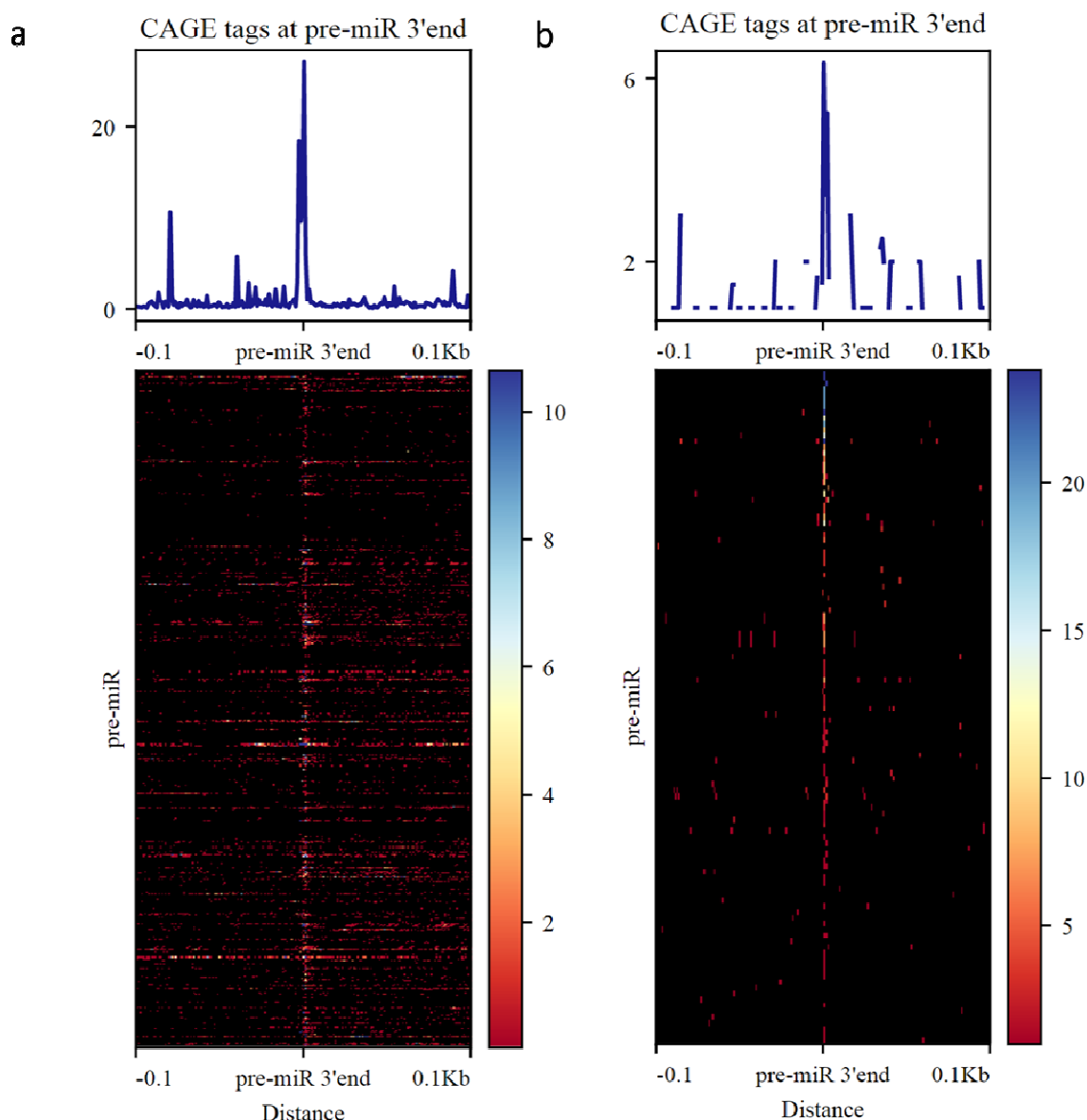
318 28 Severin, J. *et al.* Interactive visualization and analysis of large-scale sequencing datasets  
319 using ZENBU. *Nat. Biotechnol.* **32**, 217-219, doi:10.1038/nbt.2840 (2014).

320 29 Haberle, V., Forrest, A. R. R., Hayashizaki, Y., Carninci, P. & Lenhard, B. CAGEr: precise TSS  
321 data retrieval and high-resolution promoterome mining for integrative analyses. *Nucleic*  
322 *Acids Res.* **43**, e51, doi:10.1093/nar/gkv054 (2015).

323 30 Raborn, R. T., Spitze, K., Brendel, V. P. & Lynch, M. Promoter Architecture and Sex-Specific  
324 Gene Expression in *Daphnia pulex*. *Genetics* **204**, 593-612, doi:10.1534/genetics.116.193334  
325 (2016).

326





Supplementary Figure 1: The distribution of CAGE tags around the 3' end of pre-miRNAs annotated in MirGeneDB for a) human and b) zebrafish shows a clear peak for CAGE tags 1 nt downstream (i.e. the +1 nt) of the pre-miRNA 3' ends as described before<sup>26,27</sup>.