

Lawrence Berkeley National Laboratory

Environ Genomics & Systems Bio

Title

An integrated expression atlas of miRNAs and their promoters in human and mouse

Permalink

<https://escholarship.org/uc/item/6pr3p0r9>

Journal

Nature Biotechnology, 35(9)

ISSN

1087-0156

Authors

de Rie, Derek
Abugessaisa, Imad
Alam, Tanvir
et al.

Publication Date

2017-09-01

DOI

10.1038/nbt.3947

Peer reviewed



HHS Public Access

Author manuscript

Nat Biotechnol. Author manuscript; available in PMC 2018 March 01.

Published in final edited form as:

Nat Biotechnol. 2017 September ; 35(9): 872–878. doi:10.1038/nbt.3947.

An integrated expression atlas of miRNAs and their promoters in human and mouse

A full list of authors and affiliations appears at the end of the article.

Abstract

MicroRNAs (miRNAs) are short non-coding RNAs with key roles in cellular regulation. As part of the fifth edition of the Functional Annotation of Mammalian Genome (FANTOM5) project, we created an integrated expression atlas of miRNAs and their promoters by deep-sequencing 492 short RNA (sRNA) libraries, with matching Cap Analysis Gene Expression (CAGE) data, from 396 human and 47 mouse RNA samples. Promoters were identified for 1,357 human and 804 mouse miRNAs and showed strong sequence conservation between species. We also found that primary and mature miRNA expression levels were correlated, allowing us to use the primary miRNA measurements as a proxy for mature miRNA levels in a total of 1,829 human and 1,029 mouse CAGE libraries. We thus provide a broad atlas of miRNA expression and promoters in primary mammalian cells, establishing a foundation for detailed analysis of miRNA expression patterns and transcriptional control regions.

miRNAs¹ are a class of short (21–23 nt) non-coding RNAs with key roles in a wide range of biological processes including development and differentiation^{2,3}, immunity⁴, reproduction⁵, and longevity⁶. Dysregulation of miRNA expression has been implicated in numerous diseases⁷, including cancer^{8,9}. A detailed characterization of the expression profile of miRNAs across cell types and tissues is a fundamental requirement for understanding the function of miRNAs and their potential role in health and disease.

miRNAs inhibit specific mRNAs by binding to complementary sequences, usually located in the 3' UTR, leading to mRNA destabilization and a reduction in their translation output¹⁰. In the canonical miRNA biogenesis pathway^{1,11}, a primary miRNA transcript (pri-miRNA)

#Correspondence to Michiel de Hoon (michiel.dehoon@riken.jp) and Alistair Forrest (alistair.forrest@gmail.com).

§RIKEN Omics Science Center ceased to exist as of April 1st 2013 due to RIKEN reorganization.

Authors' contributions

P.A., G.Å., M.B., A.J.C., M.D., D.G., S.G., T.J.H., M.H., P.H., K.J.H., C.K., P.K., W.L., N.M., M.O., M.O.H., P.R., H.S., R.K.S., H.To., M.Y., N.Y., S.Z., P.G.Z., L.W., Y.Y., C.A.W., K.M.S., and A.R.R.F. provided RNA samples; E.A. and C.O.D. selected samples from the FANTOM5 time courses; Y.I., S.N., and H.Ta. produced the sRNA libraries; I.A., M.L., H.K., and T.K. managed the data; D.d.R., M.J.L.d.H., K.V.S., A.M.B., T.A., H.A., A.H., T.L., H.P., C.L., A.M., V.M., and M.R. carried out the bioinformatics analyses with the help of C.C.H., M.L., K.H., F.R., and J.S.; C.J.M. provided the cell ontology; K.M.S. created the Miru visualization; A.F., A.M., A.R.R.F., A.S., C.L., C.A.W., D.d.R., E.H., F.R., H.P., K.V.S., A.M.B., M.J.L.d.H., M.R., N.B., P.S., R.D., V.M., and Y.A.M. contributed to the manual miRNA promoter annotation; K.Y. and J.W.S. performed the expression validation experiments of known miRNAs; E.H. and C.A.W. performed the validation experiments of candidate miRNAs; C.G. and M.R. performed the RACE experiments; J.H. created the web visualization tool; D.d.R., A.R.R.F., and M.J.L.d.H. wrote the manuscript with the help of E.A., A.S., A.M.B., K.M.S., K.V.S., M.R., N.B., P.C., P.S., and C.A.W.; A.R.R.F. and M.J.L.d.H. designed the study; P.C. and Y.H. supervised the FANTOM5 project.

Competing financial interests

The authors declare no competing interests.

is cleaved by the endoRNase Drosha in the nucleus to excise the precursor miRNA (pre-miRNA), which is exported to the cytoplasm. The pre-miRNA has a characteristic hairpin secondary structure that is recognized and cleaved in the cytoplasm by the endoRNase Dicer, releasing the mature miRNA.

Currently, the miRBase reference database of miRNAs¹² lists 1,881 pre-miRNAs in human; around half (54%) are produced from intergenic non-coding pri-miRNA transcripts, while the remaining 46% are excised from the introns of protein-coding transcripts. A small proportion (6%) of human mature miRNAs annotated in miRBase are located in multiple pre-miRNAs encoded in different genomic loci.

Several high-throughput approaches are available to measure the expression levels of mature miRNAs, including high-throughput qPCR, microarray, and next-generation sequencing methods¹³. Profiling pri-miRNAs, which is more challenging owing to their transient character, has been accomplished by RNA-seq in cells expressing dominant-negative Drosha¹⁴. Additionally, since most pri-miRNAs are produced by RNA polymerase II and therefore have a 5' cap¹¹, they are amenable to Cap Analysis Gene Expression (CAGE) profiling^{15,16}, which identifies the pri-miRNA transcription start site and therefore the promoter region, while directly quantifying the pri-miRNA expression level.

Here, we analyze 492 sRNA sequencing libraries to evaluate the expression patterns of miRNAs in mammalian cells, with a particular emphasis on human primary cells. Each sRNA library was matched to a CAGE library produced from the same RNA sample, allowing us to create an integrated expression atlas of miRNAs and their promoters. The expression atlas can be accessed through a web interface at http://fantom.gsc.riken.jp/5/suppl/De_Rie_et_al_2017/. This work is part of the fifth edition of the Functional Annotation of Mammalian Genome project (FANTOM5)^{17,18}.

Results

Matched miRNA and CAGE expression profiles

In FANTOM5, a large collection of human and mouse primary cell types, cell lines, and tissues was profiled by CAGE to identify mRNA and long non-coding RNA transcription start sites and expression levels across a wide variety of biological states¹⁷. Here, we produced a complementary data set comprising 293 sRNA sequencing libraries using FANTOM5 RNA samples from human primary cells, and 87 sRNA libraries from RNA samples of six time courses of stimulated human cells¹⁸ (Table 1, S1 & S2). We also incorporated previously produced CAGE and sRNA sequencing libraries generated from human embryonic and induced pluripotent stem cells¹⁹ (Table 1 & S1) in our analysis. In total, our sRNA sequencing data set encompassed 121 distinct human cell types. In addition, we produced 6 sRNA sequencing libraries from human tissues, and 42 sRNA libraries from mouse samples (Table 1, S1 & S2). Most sRNA libraries were produced in biological triplicate. A matching CAGE library¹⁷⁻¹⁹ generated from the same RNA sample was available for 492 of the 500 sRNA libraries analyzed here (Table S3).

Establishing a robust set of miRNAs

Across the sRNA libraries, expression was confirmed for 98% (1842/1877) of human and 95% (1124/1186) of mouse pre-miRNAs annotated in release 21 of the miRBase database¹². To assess the confidence level of annotated miRNAs, the miRBase curators defined a set of five criteria evaluating their secondary structure and expression properties (Table 2), and used these criteria to mark 295 human pre-miRNAs as high-confidence annotations¹². Applying these criteria to the FANTOM5 sRNA data, we found that 571 human pre-miRNAs satisfied all five high-confidence criteria, 224 met four of them, and 1076 violated two or more criteria (Figure 1a). The 795 human and 502 mouse (Figure S1) pre-miRNAs satisfying at least 4 out of the 5 high-confidence criteria were defined as the FANTOM5 robust set, and the remaining 1076 human and 684 mouse pre-miRNAs as the permissive set (Table S4 and S5). The robust set encompasses 735 human and 438 mouse mature miRNAs, and covers more than 90% of the high-confidence pre-miRNAs in miRBase (Figures S2 and S3), 90% of miRNAs well-characterized in the scientific literature (Figure S4), as well as 91% (human) and 88% (mouse) of pre-miRNAs included in the manually curated MirGeneDB database²⁰ (Figure S5).

CAGE detects 3' cleavage products of Drosha

In zebrafish, the Drosha cleavage site at the 3' end of pre-miRNAs was recently found to be characterized by a distinctive CAGE peak²¹. We similarly observed a CAGE peak immediately downstream of the 3' end of human pre-miRNA loci in the ENCODE CAGE data²², and a slightly wider CAGE peak starting 1 nucleotide downstream in the FANTOM5 CAGE data^{17,18} (Figure 1b, c, human; Figure S6, mouse); the discrepancy between the ENCODE and FANTOM5 CAGE data was expected because of differences in the sequencer technologies employed (Figure S7). The ENCODE CAGE peak was found immediately downstream of the 3' end of the pre-miRNA locus (Figure S8) for 19 out of 25 pre-miRNAs with a full-length sequence in the FANTOM4 sRNA sequencing libraries²³, confirming that the CAGE peak marks the Drosha cleavage site. FANTOM5 and ENCODE CAGE tags at the peak were enriched in the nucleus (Figure S9), consistent with processing by Drosha. CAGE peaks were absent at the 3' end of pre-miRNA loci encoding mirtrons (Figure S10, human; Figure S11, mouse), which are excised by the spliceosomal machinery instead of by Drosha²⁴.

To rule out the possibility that these CAGE tags originated from an independent transcript, we analyzed the first nucleotide of the CAGE tags at the Drosha cleavage site. Most CAGE tags originating from a transcription start site have an additional guanine as their first nucleotide, as the 7-methylguanosine cap at the 5' end of transcripts produced by RNA polymerase II can be recognized as a guanine nucleotide during reverse transcription (Figure S7). No such additional guanine nucleotides were found at the Drosha CAGE peak (Figure S12), confirming that the detected RNAs were not due to an independent transcription initiation event. The lack of guanine nucleotide enrichment also suggested that the 3' Drosha cleavage products were uncapped RNAs that were nonetheless observed to some extent in the CAGE library due to their cellular abundance. Alternatively, these RNAs may have a hypermethylated cap, as previously found for small nucleolar RNAs (snoRNAs) produced by excision from a host gene transcript²⁵; no additional guanines were found as

the first nucleotide of CAGE tags mapping to the 5' end of snoRNAs (Figure S12), since hypermethylation of the cap prevents base-pairing during reverse transcription.

Excluding mirtrons, about half of the robust pre-miRNAs had a significant ($P < 0.05$) Drosha CAGE peak (52%, human, Figure 1a; 64%, mouse, Figure S1). This percentage decreased from 56% for human pre-miRNAs satisfying all five of the miRBase high-confidence criteria to 37% if one of the criteria was violated, while only 7% of miRNAs in the permissive set had a Drosha CAGE peak (Figure 1a). Similar results were obtained for mouse (Figure S1). The analysis of Drosha CAGE peaks thus provided independent support for the stringency of the selection criteria used to define the FANTOM5 robust and permissive set of miRNAs.

Discovery of candidate novel miRNAs

To discover potential miRNAs that had not been described previously, the miRDeep2 software²⁶ was applied on all unannotated sRNAs. In total, 6,543 candidate miRNAs in human (Table S6, S7) and 1,444 in mouse (Table S8, S9) were identified. Most of the candidate miRNAs were expressed at low levels, with fewer than 5% of them having sufficient tag counts on both arms of the pre-miRNA to enable a full evaluation of the high-confidence criteria (Table 2). The 282 human and 34 mouse candidate miRNAs meeting at least 4 of the 5 high-confidence criteria formed the robust candidate set, while the permissive candidate set consisted of the remaining candidate miRNAs (Table S4 and Figure S13, human; Table S5 and Figure S14, mouse). The robust candidate set comprised 279 (human) and 33 (mouse) unique mature sequences, whereas the permissive candidate set provided an additional 5,826 (human) and 1,354 (mouse) mature sequences. Nearly 11% of robust and 5% of permissive human candidate miRNAs had a significant ($P < 0.05$) Drosha CAGE peak (Figure S13; Figure S14 for mouse). Validation by qPCR of a selection of robust candidate miRNAs identified in monocyte and macrophage libraries confirmed their expression in these cell types in multiple donors (Figure S15, Table S10).

The robust candidate set showed good concordance (127/282 or 45%) with the 3,524 putative miRNAs identified recently in a study of tissue- and primate-specific miRNAs²⁷, whereas the permissive candidate set yielded a smaller overlap (352/6,261 or 6%). Few of these putative miRNAs²⁷ had a significant Drosha CAGE peak (258/3,524 or 7%), which may be due to their low expression levels in the samples surveyed in FANTOM5.

We conclude that the vast majority of canonical, highly expressed miRNAs had already been annotated. However, our analysis also provides evidence of extensive transcription of short RNAs expressed at low levels from specific genomic loci.

Expression variability of miRNAs in human primary cells

The cell type dependence of expression of individual miRNAs was evaluated by analyzing the distribution of miRNA abundance across the FANTOM5 primary cells and tissues. First, we assessed various expression normalization strategies, and found that a counts per million (cpm) normalization (i.e., dividing the tag count of each miRNA by the total number of tags mapping to miRNA loci, and multiplying by 1,000,000) yielded the best reproducibility between different donors for the same cell type, while maintaining the distinction in

expression profile between different cell types (Figure S16). We then created miRNA expression tables across the FANTOM5 samples for human (Table S11) and mouse (Table S12), using cpm normalization in our further analysis. Figure 2a shows a graphical overview of the human primary cells clustered by their robust miRNA expression profile using Miru²⁸. An interactive heatmap of the expression data is available at http://fantom.gsc.riken.jp/5/suppl/De_Rie_et_al_2017/vis_viewer/#/heatmap.

The expression levels of miRNAs varied greatly and were highly skewed, with on average five miRNAs contributing half of the total miRNA expression in a given library (Figure 2b, human; Figure S17, mouse), whereas most known and candidate miRNAs were expressed at low levels (Figure 2c, human; Figure S18, mouse). The extremely wide distribution of miRNA expression across miRNAs and cell types was confirmed by qPCR (Figure S19).

Cell ontology analysis

A cell type specificity index, analogous to the previously defined tissue specificity index²⁹, was calculated to quantify the cell type specificity of miRNA expression across the FANTOM5 collection of primary cell types (Table S13). Previously described highly cell-type-specific miRNAs included miR-122-5p, miR-142-5p, and miR-302a-5p, which were enriched in hepatocytes, leukocytes, and pluripotent stem cells, respectively (Figure 2d). In contrast, miRNAs such as miR-100-5p and miR-29a-3p were broadly expressed but specifically depleted in particular cell types (leukocytes and pluripotent stem cells, respectively; Figure 2d). Candidate miRNAs tended to be restricted to specific cell types, with 80% of the robust candidate set and 96% of the permissive candidate set having a higher cell type specificity index than the median value for robust known miRNAs (Table S13).

We then calculated the statistical significance of expression enrichment or depletion of each miRNA (Table S13) with respect to cell ontology clusters (Table S14) defined by the FANTOM5 cell ontology annotation^{30,31}, which organizes FANTOM5 samples by cell type in a hierarchical framework. Of miRNAs in the robust set, 54% had enriched expression in their statistically most significant cell ontology cluster, whereas 27% were broadly expressed, with depleted expression in their statistically most significant cell ontology cluster. The remaining 19% were expressed at low levels without statistically significant enrichment or depletion in any cell ontology cluster; understanding their functionality may need profiling in further cell types or states.

Pluripotent stem cells were characterized by cell-type-specific miRNAs, whereas cell-type-specific depletion of broadly expressed miRNAs was predominantly found in leukocytes. Examples of enriched expression not reported previously included miR-488-5p in neural cells, miR-506-3p in light melanocytes, and miR-205-5p in epithelial cells. miRNAs previously not reported as broadly expressed included miR-887-3p, which was present in most samples but was depleted in leukocytes.

Identification of miRNA promoters

We developed an automatic pipeline to identify miRNA promoters using Gencode v19 and RefSeq transcripts as candidate pri-miRNAs and the FANTOM5 CAGE data as putative

transcription start sites. This pipeline predicted promoters for 539 robust, 623 permissive, and 3,951 candidate pre-miRNAs in human (Table S15), and for 358 robust, 446 permissive, and 994 candidate pre-miRNAs in mouse (Table S16). Manual curation by two independent annotators confirmed the selected promoter for 512 (95%) robust pre-miRNAs; the computationally selected promoter was corrected for 26 pre-miRNAs and dropped for 1 pre-miRNA. Manual curation furthermore identified the promoter for an additional 196, mostly intergenic, pre-miRNAs, thereby generating the—to our knowledge—largest miRNA promoter collection to date (Table S17, Figure S20a). Across the human robust set, an associated ENCODE RAMPAGE (RNA Annotation and Mapping of Promoters for the Analysis of Gene Expression³²) 5' end was found within 300 base pairs of more than 75% of the FANTOM5 curated promoters both for intergenic and intronic miRNAs, outperforming the miRGen³³, Chang *et al.*¹⁴, miRStart³⁴, and TSmiR³⁵ collections of miRNA promoter annotations (Figure S20b). The median distance between the FANTOM5 annotated miRNA promoter and the associated RAMPAGE 5' end was 1 nucleotide, and was thereby closer than any of the existing miRNA promoter annotations (Figure S20c). RACE experiments confirmed that the transcripts generated at the identified promoter extended to the mature miRNA for 6 out of 7 miRNAs (Figure S21, Table S18). RNA-seq data¹⁴ of cells expressing a dominant-negative Drosha protein provided additional evidence for the FANTOM5 annotated pri-miRNAs, with 483 out of 607 pri-miRNAs (80%) having a 5' end within 300 base pairs of an RNA-seq transcript assembly extending to the mature miRNA locus (Figure S22).

Both in human and in mouse, promoter sequences of intronic and intergenic miRNAs, like those of transcription factor genes, were highly conserved across species compared to the promoter sequences of protein-coding genes and of long non-coding RNAs (Figures 3a and S23, human; Figure S24, mouse). The distance between the transcription start site of the pri-miRNA and the mature miRNA locus was strongly conserved between human and mouse both for intronic miRNAs (Spearman correlation = 0.90; $n = 78$; Student $t = 18.27$; $P = 2 \times 10^{-29}$ two-sided) and for intergenic miRNAs (Spearman correlation = 0.86; $n = 27$; Student $t = 8.33$; $P = 1 \times 10^{-8}$ two-sided) (Figures 3b and S25). While this suggests that pri-miRNA transcripts may have some functional role beyond providing the substrate for pre-miRNA excision, there was no evidence of substantially elevated sequence conservation across species in pri-miRNAs (Figure S26).

Correlation of mature miRNA and pri-miRNA expression levels

The expression levels of mature miRNAs correlated with the CAGE expression levels of the associated promoter, with comparable correlation values for intergenic (average Spearman $r = 0.27$; $n = 180$; $P = 1 \times 10^{-30}$, Mann-Whitney U test, one-sided) and intronic (average Spearman $r = 0.25$; $n = 362$; $P = 2 \times 10^{-53}$, Mann-Whitney U test, one-sided) miRNAs (Figure 3c and S27; Table S19). The correlation was substantially higher for highly differentially expressed miRNAs (average Spearman $r = 0.59$; $n = 159$; $P = 1 \times 10^{-35}$, Mann-Whitney U test, one-sided), and exceeded correlations found for previously published^{14,33–35} miRNA promoter annotations (Figure S20d). About 11% of pri-miRNAs in human were polycistronic, containing multiple mature miRNAs with highly correlated expression levels (average Spearman $r = 0.74$; $n = 1,372$; $P < 10^{-100}$, Mann-Whitney U test,

one-sided) (Figure 3d and S28). Together this suggests that miRNA expression is primarily regulated at the transcriptional level.

Using the CAGE expression level of the pri-miRNA as a proxy for the mature miRNA expression level, we extended the FANTOM5 miRNA expression atlas to the full breadth of the 1,829 (human) and 1,029 (mouse) libraries in the FANTOM5 CAGE expression compendium^{17–19}. This allowed us to assess miRNA expression also in samples for which only a CAGE library was available, covering an additional 49 primary cell types, 245 cell lines, 138 tissue types, and 13 time courses in human, and an additional 48 primary cell types, 1 cell line, 234 tissue types, and 12 time courses in mouse. A cell ontology analysis was performed using the CAGE expression pattern of each human pri-miRNA (Tables S15 and S17) across 338 cell ontology clusters (Table S20), encompassing 636 CAGE libraries. This showed enriched expression of mir-202 in gonad, of mir-208a, known to be a key regulator of cardiac function³⁶, in heart, as well as of multiple miRNAs in brain, including mir-488, mir-556, and mir-885. Lastly, the CAGE data allowed us to measure the individual contribution of each paralog to the expression of miRNAs encoded multiple times in the human genome, providing evidence for differential regulation of paralogs in different cell types and tissues. For example, we found that mir-128-1 was expressed in most samples, while its paralog mir-128-2 was highly enriched in brain (Figure 3e).

Transcriptional regulation of miRNA expression

The accuracy of the predicted miRNA promoter regions was assessed using the Motif Activity Response Analysis (MARA) framework³⁷ (Figure S29). Using this framework, we predicted the expression levels of mature miRNAs based on the presence of putative transcription factor binding sites in the identified miRNA promoter region, and compared to the expression levels of the mature miRNAs observed in the sRNA libraries. The prediction accuracy of the FANTOM5 miRNA promoter atlas outperformed those of previously published miRNA promoter annotations^{14,33–35} (Figure S20e).

Discussion

miRNAs are key factors that contribute to cellular regulation by targeting specific transcripts for translational repression or for degradation. Advances in sequencing technology led to an increase in sequencing depth from nearly 1,300 reads per sRNA library in the first miRNA atlas³⁸ to nearly 4.4 million reads per library in FANTOM5, allowing an accurate measurement of the expression even of miRNAs not highly expressed. Such miRNAs may be abundant in a few cells in the population sampled, or in cell types, cell lines, or cellular conditions that are not included in our sample collection. Alternatively, they may be a signature of the ongoing evolution of the human miRNA repertoire. In particular, pervasive transcription of mammalian genomes^{22,39,40} generates a large number of hairpin secondary structures, which are widely encoded in the genome, that can act as substrates for processing by Drosha in the nucleus and Dicer in the cytoplasm. Whereas the majority of the sRNAs thus generated may be evolutionarily neutral and remain expressed at low levels, some of them may provide a selective advantage, develop higher expression levels during evolution⁴¹, and become fixed in the genome as core miRNAs. Finally, we note that in spite

of the breadth and depth of the FANTOM5 sRNA sequencing data, most sRNAs currently annotated as miRNAs failed to meet multiple high-confidence criteria, and may belong to a different class of short non-coding RNAs, such as transcription initiation RNAs²³ or DNA damage response RNAs⁴², or may be degradation products⁴³.

Compared to existing miRNA expression atlases^{29,38}, the FANTOM5 atlas covers the widest range of normal primary cells, enabling detailed analyses of miRNA expression and their contribution to establishing and maintaining cell type identity. The candidate miRNAs not reported previously were in particular highly specific to cell type, and may therefore be missed in miRNA profiling studies in tissues rather than in specific cell types.

We found extensive evidence that CAGE peaks observed at the Drosha cleavage site are due to the downstream RNA fragment generated by Drosha processing of the pri-miRNA. Analysis of these CAGE tags suggested that these RNA fragments do not have a 7-methylguanosine cap, but may instead be uncapped or, alternatively, have a non-canonical cap. For polycistronic pri-miRNAs, such a cap may play a role in preventing rapid degradation of the downstream fragment, which itself may contain miRNAs.

The MARA analysis allowed us to predict miRNA expression levels based on the regulatory motifs found in the miRNA promoter region, indicating that transcriptional regulation plays a central role in governing miRNA expression levels. Comparing the promoters of miRNAs, protein-coding genes, and long non-coding RNAs showed a similar prevalence of transcription factor binding sites in proximal promoter regions (data not shown), suggesting that the basic mechanisms of transcriptional regulation are largely the same for these three classes of gene products. The identification of miRNA promoter regions as described in this work therefore paves the way for a detailed analysis of the transcriptional regulation of miRNA expression using the same computational and experimental methods that have previously proven their efficacy in the analysis of gene expression.

Methods

Samples and library preparation

Short RNA libraries were prepared following the Illumina TruSeq Small RNA Sample Preparation protocol (catalog number RS-200-0012, RS-200-0024, RS-200-0036, RS-200-0048) using the same RNA samples from which CAGE libraries were produced previously^{17,18}, as well as one additional RNA sample without a matching CAGE library. RNA samples not previously described are listed in Table S2. TruSeq Small RNA Sample Prep Index Sequences were used as bar codes to allow pooling of multiple samples in one library. The short RNA libraries were sequenced using the Illumina HiSeq2000 sequencer in single-read, 50-base mode. The metadata of all FANTOM5 RNA samples, including those used for sRNA sequencing, are available in the FANTOM5 Semantic catalog of Samples, Transcription initiation And Regulators⁴⁵ (SSTAR; <http://fantom.gsc.riken.jp/5/sstar>). SSTAR sample pages also provide links to the FANTOM5 miRNA expression atlas web interface.

Data processing

We extracted the short RNA sequences from the raw sequences using in-house scripts. We removed linker artifact sequences using TagDust⁴⁶ version 1.13, ribosomal sequences using rRNA dust¹⁷ version 1.00, and filtered against mature tRNAs, ribosomal RNA, and 7SL RNA using global alignment. We mapped the remaining sequences using the Burrows-Wheeler Alignment (bwa) tool⁴⁷ version 0.5.9-r16 to genome assembly hg19 (human) or mm9 (mouse), including chromosome Y if the donor was known to be male. Table S3 shows the number of short RNA sequences mapped to the genome for each sample. Two samples had fewer than 100,000 mapped tags and were discarded from further analysis.

Short RNA annotation and filtering

We used release 21 of the miRBase database¹², lifted over to genome assembly hg19 (human) or mm9 (mouse), as our reference set of known miRNAs. Four pre-miRNAs in human that could not be lifted over to genome assembly hg19 and an additional six human pre-miRNAs that were lifted over to unplaced chromosomes were excluded from the analysis. We annotated all mapped short RNA reads mapping to genomic loci for ribosomal RNA, tRNAs, the RNA component 7SL of the signal recognition particle, small nuclear RNAs, small nucleolar RNAs, small Cajal body-specific RNA, small cytoplasmic RNAs, and piRNAs (piwi-interacting RNAs). We corrected for cross-mapping as described previously⁴⁸, discarding all mappings to unannotated loci if the short RNA sequence could be mapped to an annotated locus instead.

Drosha CAGE peak analysis

We calculated the total number of CAGE tags starting at each genomic position across all 1,885 (human) and 1,202 (mouse) FANTOM5 CAGE libraries^{17,18}, as available at <http://fantom.gsc.riken.jp/5/datafiles/latest/basic/>, as well as all 145 human ENCODE CAGE data²², which we downloaded from <http://hgdownload.cse.ucsc.edu/goldenpath/hg19/encodeDCC/wgEncodeRikenCage/>. We defined the 3' end of the pre-miRNA as the 3' nucleotide of the mature miRNA on the 3' arm of the pre-miRNA; the expected Drosha cleavage site is immediately downstream of this nucleotide. For each pre-miRNA in the robust set, we calculated, for each position with respect to the expected Drosha cleavage site, the total number of CAGE tags in the CAGE libraries. We normalized by dividing by the sum over the positions to obtain the CAGE profile with respect to the expected Drosha cleavage site for each pre-miRNA. We then summed the CAGE profiles across the pre-miRNAs to obtain the average CAGE profile with respect to the expected Drosha cleavage site. Based on this profile, we selected a 9-basepair window between -2 and +7 base pairs with respect to the expected Drosha cleavage site for the FANTOM5 CAGE data, and an 8-basepair window between -2 and +6 base pairs for the ENCODE CAGE data, as the Drosha CAGE peak window for a given pre-miRNA.

For each pre-miRNA, we counted the number of CAGE tags with a 5' end within this window, as well as the number of CAGE tags with a 5' end anywhere between the pre-miRNA boundaries. Since CAGE tags tend to occur in clusters on the genome, we expect the distribution of the CAGE tag counts to be heavily overdispersed compared to the Poisson distribution. We therefore used the negative binomial distribution instead, with the

dispersion parameter r estimated by fitting the distribution to the number of CAGE tags in any 8- or 9-basepair window on the human or mouse genome. This resulted in a dispersion of 1.856943 and 1.616542 for the FANTOM5 human and mouse CAGE data, respectively (using a 9-basepair window), and 0.325001 for the ENCODE CAGE data (using a 8-basepair window). Using these dispersion values, we calculated the statistical significance of the FANTOM5 and ENCODE CAGE peaks given the number k of CAGE tags within the window, the number K of CAGE tags within the pre-miRNA, the window size w , as well as the genomic extent L of the pre-miRNA as $I_p(k, r)$, where I is the regularized incomplete beta function and $p = \mu / (r + \mu)$, with $\mu = wK/L$ the expected number of tags at the Drosha CAGE peak under the null hypothesis.

For human, we calculated an overall statistical significance value by combining the FANTOM5 and ENCODE statistical significance into a single P -value using Fisher's method.

The relative occurrence of CAGE tags in different subcellular fractionations (Figure S9) and the bias in the first nucleotide of CAGE tags (Figure S12) were evaluated for pre-miRNAs in the robust set with a statistically significant Drosha CAGE peak.

Identification of candidate novel miRNA

Candidate novel miRNAs were identified using miRDeep2 (ref. 26), resulting in 7,461 (human) and 2,034 (mouse) predicted pre-miRNAs, including 918 (human) and 590 (mouse) known pre-miRNAs. To avoid predicted miRNAs from failing the miRBase high-confidence criteria owing to flaws in the predicted secondary structure, we repeated the secondary structure calculation for each predicted miRNA by applying RNAfold⁴⁹ version 2.1.2 on the sequence of the precursor miRNA while constraining the structure by allowing nucleotides in each arm of the hairpin to only base-pair to nucleotides in the other arm of the hairpin. In the comparison of the candidate novel miRNAs with the 3,524 tissue- and primate-specific miRNAs published recently²⁷, we required the pre-miRNAs to overlap by at least 80%.

Validation of candidate novel miRNA expression by qPCR

Fresh buffy coat was obtained from the Red Cross following approval from the human research ethics committee of The University of Melbourne (ethics ID 1646608.1) and material supply agreement with the Red Cross (16-05VIC-21).

Peripheral blood mononuclear cells (PBMC) were isolated from buffy coat using Ficoll hypaque (GE Healthcare, Uppsala, Sweden) as described previously⁵⁰. CD14+ human monocytes were isolated from PBMC using human CD14+ magnetic beads (Milteny Biotec, Sydney, New South Wales, Australia). CD14+ monocytes were differentiated to macrophages in complete RPMI1640 media supplemented with 10% fetal calf serum and 100 ng ml⁻¹ human macrophage colony-stimulating factor (M-CSF) (PeproTech, Rehovot, Israel) for 5 days. Suspended cells were removed and adherent cells were washed with PBS before macrophages were collected.

miRNAs were isolated from monocytes and macrophages using mirVana miRNA Isolation Kit (Life Technologies, Melbourne, Victoria, Australia) following the manufacturer's

protocol. Briefly, cells were lysed in lysis buffer followed by phenol extraction, and miRNAs were isolated from the phenol aqueous phase using a spin column followed by elution in RNase-free water. Following manufacturer's protocol, cDNA synthesis was performed using miScript PCR Starter Kit (Qiagen, Hilden, Germany) by ligating a poly(A) tail to the miRNA followed by reverse transcription in the presence of universal tag. Samples without reverse transcriptase but with all other components were included and used as negative controls.

Forward primers specific to the candidate novel miRNAs were designed using miRprimer⁵¹ (Table S10). Real-time PCR was performed using miScript PCR Starter Kit (Qiagen, Hilden, Germany) and following the manufacturer's protocol. The PCR reaction was set up with the custom-made forward primers and the universal reverse primer supplied with the kit. No-template controls and cDNA samples without reverse transcriptase were included as negative controls. Thermal cycling was performed as suggested by the manufacturer's protocol.

The expression levels of a wide range of miRNAs have been analyzed using our miRNA PCR assay in order to evaluate the sensitivity of the assay and determine the confidence of our results. Short RNAs commonly used as a reference, including RNU6 and let-7a-5p (ref. 52), showed relatively high expression levels. Other miRNAs that are highly conserved in metazoans or known to be expressed in myeloid cells, including miR-191-5p (ref. 53), miR-15a-5p (ref. 54), miR-206 (ref. 55), miR-335-5p (ref. 56), and miR-339-3p (ref. 56), were included and used as positive controls, and showed moderate expression levels. Expression levels of miRNAs reported to be cell markers for other cell types and assumed to be lowly expressed in myeloid cells, including miR-153-3p (ref. 57) and miR-345-5p (ref. 58), were also analyzed in order to determine the detection limit of the assay. Our results demonstrate that the miRNA PCR assay could specifically detect the presence of the target miRNAs, and measure a wide spectrum of expression levels. The expression levels of the selected candidate novel miRNAs fell within the detection spectrum of our miRNA PCR assay, proving the reliability of our results.

Evaluation of miRNA expression normalization strategies

We counted the number of short RNA sequences with a length of 18–25 nucleotides overlapping the mature miRNA loci in each of the primary cell samples. We then applied the following normalization strategies:

- CPM (counts per million): Divide the count by the sum of counts for mature miRNAs in the robust set, and multiply by 1,000,000;
- TMM (trimmed mean of M values): Apply the “calcNormFactors” function in edgeR⁵⁹ with method “TMM” to the table of counts;
- RLE (relative log expression): Apply the “calcNormFactors” function in edgeR⁵⁹ with method “RLE” to the table of counts;
- DESeq (effective library size): Apply the “estimateSizeFactorsForMatrix” function in DESeq⁶⁰ to the table of counts;

- UQ (upper quantile normalization): Divide the count by the sum of the counts of the top-25% most abundant miRNAs in each sample;
- UD (upper decile normalization): Divide the count by the sum of the counts of the top-10% most abundant miRNAs in each sample.

To evaluate each normalization strategy, we divided the primary cell samples in FANTOM5 into groups ($n = 96$) of independent donors of the same cell type. For each cell type group, we calculated the variance for each miRNA across the donors. To find the error between different cell types, we first calculated the average expression for each miRNA across donors in each cell type group, and then calculated the difference in the average expression between each pair of cell type groups ($n = \frac{1}{2} \times 96 \times 95 = 4,560$) for each miRNA. To evaluate the total error, we calculated the mean squared error across miRNAs for each cell type group, as well as the mean squared error across miRNAs for each pair of cell type groups, and took the square root of each to find the root mean square (RMS) with cell type groups and between cell type groups (Figure S16a). We averaged the RMS error over the $n = 96$ cell type groups, and over the $n = 4,560$ pairs of cell type groups, and calculated the ratio of the average RMS error within cell types to the average RMS error between cell types (Figure S16b). To evaluate the standard error (Figure S16c), we calculated the mean squared error across cell type groups for each miRNA, as well as the mean squared error across pairs of cell type groups for each miRNA, took the square root, and plotted the resulting RMS value for each miRNA against its mean expression level. We then used linear regression to calculate the slope of the RMS error within each cell type and between different cell types as a function of the miRNA expression level. Dividing these two slopes yielded the ratio in RMS error within cell types and between different cell types, normalized by miRNA expression level (Figure S16d).

Clustering and visualization of miRNA expression patterns

miRNAs were clustered based on their expression patterns using the network visualization and analysis tool Miru²⁸ (<http://kajeka.com/miru/miru-about/>). The Pearson correlation was calculated for each pair of miRNAs. A modified Fruchterman-Rheingold algorithm was used to lay out the network graph in 3-dimensional space, in which 502 nodes representing miRNAs were connected by 3,369 weighted, undirected edges representing correlations of at least 0.6 between expression patterns. Areas of high connectivity and correlation, representing groups of miRNAs with similar expression profiles, were identified using the Markov clustering algorithm (MCL) with an MCL inflation value of 2.2. Clusters were manually annotated based on the cell type or tissue of greatest expression. All nodes in one cluster and the label describing the cluster are shown in the same color. The smallest labeled cluster contains six nodes; for clarity, smaller clusters have not been labeled but can be identified by groups of nodes of the same color.

Validation of miRNA expression quantitation by qPCR

Expression of selected miRNAs was measured using the TaqMan[®] MicroRNA Assay (Applied Biosystems) according to its protocol. RNA samples 11544, 11624, 11705 (CD19+ B cells, donor1, 2, and 3), 11269, 11346, 11418 (dermal fibroblast donor1, 2, and 3), 12626

(H9 embryonic stem cells), and 11523, 11603, 11684 (hepatocyte donor1, 2, and 3) (Table S3) were used after confirmation of the RNA quality by measuring the RNA Integrity Number (RIN) value using a TapeStation and the 260/280 and 260/230 ratios using NanoDrop. The Ct values obtained were normalized against the Ct value of small nucleolar RNA SNORD48.

Cell type specificity index

Following the definition of the tissue specificity index (TSI)²⁹, we defined the cell type specificity index of miRNA j as

$$\text{index}_j = \frac{1}{N-1} \sum_{i=1}^N \left(1 - \frac{x_{j,i}}{\max_{i'} x_{j,i'}} \right)$$

where N is the number of primary cell types in FANTOM5, and $x_{j,i}$ is the expression in counts-per-million of miRNA j in cell type i , averaged over independent donors.

Guide strand selection

For each pre-miRNA, we designated the hairpin arm with the highest expression level (in counts-per-million) in any of the FANTOM5 samples as the guide strand, and refer to the opposite arm as the passenger strand.

Cell ontology analysis

We used the FANTOM5 cell ontology^{30,31} to create cell ontology clusters (Tables S14 and S20). We performed a likelihood-ratio test comparing the expression data between the samples in each cell ontology cluster and the background, consisting of all other samples listed in Tables S14 and S20, modeling the tag counts by a negative binomial distribution. For each miRNA, we selected the three cell ontology terms for which the expression in the cell ontology cluster compared to the background was statistically most significantly higher, and the three cell ontology terms for which the expression in the cell ontology cluster compared to the background was statistically most significantly lower. The P -values listed in Tables S13, S15, and S17 for each miRNA for specific cell ontology clusters were not corrected for multiple testing.

Identification of miRNA promoters

Candidate pri-miRNAs consisted of transcripts annotated in Gencode⁶¹ v19 (human) or vM5 (mouse) or in the NCBI Entrez Gene database⁶². For each pre-miRNA, we selected all candidate pri-miRNAs with a transcription start site upstream of the pre-miRNA and a 3' end downstream of the 5' end of the pre-miRNA, and defined all FANTOM5 permissive CAGE peaks¹⁷ within the genomic region from 500 bp upstream of the 5' end of the pri-miRNA to the 5' end of the pre-miRNA as the set of candidate promoters associated with the pre-miRNA. We averaged the expression level (in tags per million) of each candidate promoter across all FANTOM5 CAGE samples, and selected the candidate promoter with

the highest average expression level as the (computationally predicted) promoter of the miRNA. Each human miRNA in the robust set was manually curated by two annotators.

Validation of miRNA promoters by RAMPAGE

We downloaded all 212 BAM files containing ENCODE RAMPAGE sequencing data mapped to human genome assembly hg19 that were not marked as “low read depth” or “low replicate concordance”. We retained the 5′ end positions of RAMPAGE transcripts with a 3′ end within 1,000 basepairs of a pre-miRNA locus, discarding 5′ end positions supported by fewer than 5 RAMPAGE transcripts, and associated the remaining 5′ end positions with the pre-miRNA as putative transcription start sites.

Validation of miRNA promoters by RACE

We mixed 4.0 µl 5X First-Strand Buffer, 0.5 µl DTT (100 mM; Invitrogen, catalog number 70726), 1.0 µl dNTP Mix (20 mM), spun briefly in a microcentrifuge, and kept at room temperature. We combined 1.0–10.0 µl with 1 µg total RNA from monocytes, macrophages, and dendritic cells, 1.0 µl Random Primer Mix (N-15) (20 µM), and 0–9 µl sterile water to reach a total volume of 11.0 µl in separate microcentrifuge tubes, mixed the contents and spun the tubes briefly. We incubated these tubes at 72 °C for 3 minutes, and then cooled to 42 °C for 2 minutes. After cooling, we spun the tubes for 10 seconds at 14,000 g to collect the contents at the bottom. Next, we added 1.0 µl of Smarter oligo (20 µM) per reaction, and mixed well by vortexing and spun the tube briefly in a microcentrifuge. We then added 0.5 µl RNase Inhibitor (40 U/µl; Invitrogen RNaseOUT™, catalog number 10777019) and 2.0 µl SMARTScribe Reverse Transcriptase (100 U; Clontech, catalog number 639537) to the buffer mix, and mixed these reagents at room temperature. Next, we added 8.0 µl of the master mix to the RNA solution, mixed the contents of the tubes by gently pipetting, and spun the tubes briefly. We incubated the tubes at 42 °C for 90 minutes and heated the tubes at 70 °C for 10 minutes in a hot-lid thermal cycler. We then added 90 µl Tricine-EDTA buffer to each tube.

We prepared the master mix for the first PCR by combining 2.5 µl of the cDNA solution, 5.0 µl 10X Advantage 2 PCR buffer (Clontech, catalog number 639207), 1.0 µl dNTP Mix (10 mM each) 50X Advantage 2 Polymerase Mix (Clontech), 1.0 µl of the smarterRACE_forward primer at 10 pmol/µl, 1.0 µl of the miRNA-specific outer primer (Table S18) at 10 pmol/µl, and added PCR-grade water to reach a volume of 50 µl. We ran a 2-step PCR program consisting of 1 minute at 95 °C, 25 cycles of 30 seconds at 95 °C followed by 70 seconds at 68 °C, 7 minutes at 68 °C, and finishing at 8 °C. We diluted 5 µl of the primary PCR product into 245 µl of Tricine-EDTA buffer.

We prepared the master mix for the second PCR by combining 5.0 µl of the product of the first PCR after dilution with 5 µl of the 10X Advantage 2 PCR buffer, 1.0 µl dNTP Mix (10 mM), 1.0 µl of 50X Advantage 2 Polymerase Mix (Clontech), 2.0 µl of the Nextera_i7 primer, 2.0 µl of the miRNA-specific inner primer (Table S18), and 34 µl of PCR-grade water. We ran a 2-step PCR program consisting of 1 minute at 95 °C, 20 cycles of 30 seconds at 95 °C followed by 70 seconds at 68 °C, 7 minutes at 68 °C, and finishing at 8 °C. We diluted 5 µl of the PCR product into 245 µl of Tricine-EDTA buffer.

We prepared the master mix for the third PCR by combining 5.0 μ l of the PCR product of the second PCR with 5.0 μ l of 10X Advantage 2 PCR buffer, 1.0 μ l dNTP mix (10 mM), 1.0 μ l of 40X Advantage 2 Polymerase Mix (Clontech), 2.0 μ l of the Nextera_i7 primer, 2.0 μ l of the Nextera_i5 primer, and 34 μ l of PCR-grade water. We purified by AMPure at a 1.8 ratio, checked 2 μ l of the second PCR product on a TapeStation, kept the libraries at -20 °C until sequencing, and pooled the PCR products, each with a different barcode combination before paired-end sequencing on a MiSeq sequencer (Illumina).

We mapped the sequencing data to the human genome using Blat⁶³, merged each pair into a single mapped transcript, and retained transcripts that overlap an inner primer. The histograms in Figure S21 show the position of the 5' end of these transcripts.

Promoter sequence conservation analysis

We previously compiled a list of transcription factors in human and mouse¹⁷. Protein-coding genes and lncRNAs consisted of all other genes annotated in the NCBI Entrez Gene database⁶² as protein coding or miscRNA, respectively. For each gene in these three categories, we selected the associated p1 CAGE peak as defined previously¹⁷ as the gene promoter, and discarded all genes without an associated CAGE peak. We then found the phastCons conservation score⁴⁴, obtained from the UCSC Genome Browser database⁶⁴, for the alignment of 99 vertebrate organisms against the human genome hg19, as a function of position relative to the transcription start site for each gene and miRNA, and averaged these scores for each category at each position.

Construction of the FANTOM5 miRNA expression atlas of miRNAs

CAGE tag start site (CTSS) files³¹, excluding universal and whole body RNA samples, were downloaded from <http://fantom.gsc.riken.jp/5/datafiles/latest/basic/>. CAGE tag counts for technical replicates of the same RNA sample were summed for each genomic position. CAGE libraries published by Fort *et al.*¹⁹ were downloaded from DDBJ, accession DRA000914. The number of CAGE tags at each genome position were counted to generate CTSS files, and pri-miRNA expression tables were generated by summing the CAGE tags under each promoter, calculating the total number of tags mapped to the genome, and using this number to normalize to tags per million (tpm).

Mature miRNA expression tables were generated by counting the number of sRNA tags to each miRNA locus, calculating the total number of tags mapping to the robust miRNAs, and using this number to normalize to counts per million (cpm).

The CAGE and sRNA expression tables are available for download at the miRNA expression viewer at http://fantom.gsc.riken.jp/5/suppl/De_Rie_et_al_2017/.

To generate the heatmap, we averaged the cpm-normalized expression values of each miRNA across donors for each cell type, and converted the expression profile of each miRNA to Z-scores by subtracting the mean and dividing by the standard deviation across cell types. The heatmap was sorted both for cell types and for miRNAs by performing

centroid-linkage hierarchical clustering, using the Pearson correlation as the similarity measure.

Motif activity response analysis (MARA)

The genome-wide predictions of transcription factor binding sites were produced as described previously¹⁸. Briefly, we downloaded the whole-genome alignment of the human genome hg19 against 99 other vertebrate genomes, and of the mouse genome mm9 against 29 other vertebrate genomes, from the UCSC Genome Browser database⁶⁴, and extracted the multiple alignments of human, macaque, mouse, rat, cow, horse, dog, opossum, and chicken. We divided the genome into segments and realigned each segment using T-Coffee⁶⁵, and generated genome-wide transcription factor binding site (TFBS) predictions using MotEvo⁶⁶ for the SwissRegulon set of position-weight matrix motifs⁶⁷ (Figure S29a). We then counted the number of predicted TFBSs for each motif in the -300 to +100 base pair proximal promoter regions of genes in the NCBI Entrez Gene database⁶², excluding all miRNA promoters (Figure S29b). Next, we used MARA³⁷ to decompose the FANTOM5 CAGE expression profiles of these promoters in terms of their associated motifs, yielding the activity profile of each motif across the FANTOM primary samples (Figure S29c). We then counted the number of TFBSs for each motif in the -300 to +100 base pair proximal promoter region of each miRNA (Figure S29d), and predicted the miRNA expression level by calculating the weighted sum of the activities for motifs found (Figure S29e). We compared the predicted expression levels to the expression levels of the mature miRNA observed in the FANTOM5 sRNA sequencing data (Figure S29f) and calculated their correlation (Figure S29g) as a measure of the accuracy of the miRNA promoter identification. Following the MARA procedure³⁷, we normalized the cpm expression values of miRNAs by adding 0.5, taking the base-2 logarithm, subtracting the mean across samples, and finally subtracting the mean across miRNAs. We defined strongly differentially expressed miRNAs, included in Figure S20d and S20e, as those that had a standard deviation in expression, after normalization, across samples larger than 2.

Data availability

Raw sequencing data of the sRNA libraries are available at the DNA Data Bank of Japan (DDBJ; <http://www.ddbj.nig.ac.jp/>) under accession numbers DRA001101, DRA002711, DRA003804, and DRA003807, and for the RACE experiments at the NCBI Gene Expression Omnibus (NCBI GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE98695.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Authors

Derek de Rie^{1,2}, Imad Abugessaisa¹, Tanvir Alam³, Erik Arner^{1,4}, Peter Arner⁵, Haitham Ashoor³, Gaby Åström⁵, Magda Babina⁶, Nicolas Bertin^{1,4,7}, A. Maxwell Burroughs^{1,4,8}, Ailsa J. Carlisle⁹, Carsten O. Daub^{1,4}, Michael Detmar¹⁰, Ruslan

Deviatiiarov^{1,11}, Alexandre Fort^{1,4}, Claudia Gebhard^{12,13}, Daniel Goldowitz¹⁴, Sven Guhl⁶, Thomas J. Ha¹⁴, Jayson Harshbarger^{1,4}, Akira Hasegawa^{1,4}, Kosuke Hashimoto^{1,4}, Meenhard Herlyn¹⁵, Peter Heutink¹⁶, Kelly J. Hitchens¹⁷, Chung Chau Hon¹, Edward Huang^{18,19}, Yuri Ishizu^{1,4}, Chieko Kai²⁰, Takeya Kasukawa¹, Peter Klinken²¹, Timo Lassmann^{1,4,22}, Charles-Henri Lecellier^{14,23}, Weonju Lee²⁴, Marina Lizio^{1,4}, Vsevolod Makeev^{25,26,27}, Anthony Mathelier¹⁴, Yulia A. Medvedeva^{25,28,29}, Niklas Mejhert⁵, Christopher J. Mungall³⁰, Shohei Noma^{1,4}, Mitsuhiro Ohshima³¹, Mariko Okada-Hatakeyama^{32,33}, Helena Persson³⁴, Patrizia Rizzu¹⁶, Filip Roudnický¹⁰, Pål Sætrom³⁵, Hiroki Sato²⁰, Jessica Severin^{1,4}, Jay W. Shin^{1,4}, Rolf K. Swoboda¹⁵, Hiroshi Tarui^{1,4}, Hiroo Toyoda³⁶, Kristoffer Vitting-Seerup³⁷, Louise Winteringham²¹, Yoko Yamaguchi³⁸, Kayoko Yasuzawa¹, Misako Yoneda²⁰, Noriko Yumoto³³, Susan Zabierowski³⁹, Peter G. Zhang¹⁴, Christine A. Wells^{18,19}, Kim M. Summers^{9,40}, Hideya Kawaji^{1,4,41}, Albin Sandelin³⁷, Michael Rehli^{12,13}, the FANTOM consortium, Yoshihide Hayashizaki^{4,41}, Piero Carninci^{1,4}, Alistair R. R. Forrest^{1,4,21,#}, and Michiel J. L. de Hoon^{1,4,#}

Affiliations

¹Division of Genomic Technologies, RIKEN Center for Life Science Technologies, Yokohama, Japan ²Centre for Integrative Bioinformatics (IBIVU), VU University Amsterdam, Amsterdam, The Netherlands ³Computational Bioscience Research Center, Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi Arabia ⁴RIKEN Omics Science Center (OSC), Yokohama, Japan ⁵Department of Medicine, Karolinska Institutet at Karolinska University Hospital, Huddinge, Sweden ⁶Department of Dermatology and Allergy, Charité Campus Mitte, Universitätsmedizin Berlin, Berlin, Germany ⁷Cancer Science Institute of Singapore, National University of Singapore, Singapore ⁸National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA ⁹The Roslin Institute and Royal (Dick) School of Veterinary Studies, University of Edinburgh, Edinburgh, UK ¹⁰Institute of Pharmaceutical sciences, Swiss Federal Institute of Technology (ETH) Zürich, Zürich, Switzerland ¹¹Institute of Fundamental Medicine and Biology, Kazan Federal University, Kazan, Russia ¹²Department of Internal Medicine III, University Hospital Regensburg, Regensburg, Germany ¹³Regensburg Centre for Interventional Immunology (RCI), Regensburg, Germany ¹⁴Department of Medical Genetics, Centre for Molecular Medicine and Therapeutics, Child and Family Research Institute, University of British Columbia, Vancouver, British Columbia, Canada ¹⁵Melanoma Research Center, The Wistar Institute, Philadelphia, Pennsylvania, USA ¹⁶German Center for Neurodegenerative Diseases (DZNE), Tübingen, Germany ¹⁷Australian Infectious Diseases Research Centre (AID), University of Queensland, Brisbane, Queensland, Australia ¹⁸The University of Melbourne Centre for Stem Cell Systems, School of Biomedical Sciences, The University of Melbourne, Victoria, Australia ¹⁹Walter and Eliza Hall Institute of Medical Research, Melbourne, Victoria, Australia ²⁰Laboratory Animal Research Center, Institute of Medical Science, The University of Tokyo, Tokyo, Japan ²¹Harry Perkins Institute of Medical Research, and the Centre for Medical

Research, University of Western Australia, QEII Medical Centre, Perth, Western Australia, Australia ²²Telethon Kids Institute, The University of Western Australia, Subiaco, Western Australia, Australia ²³Institute of Molecular Genetics of Montpellier, Montpellier, France ²⁴Department of Dermatology, Kyungpook National University School of Medicine, Daegu, South Korea ²⁵Vavilov Institute of General Genetics, Russian Academy of Sciences, Moscow, Russia ²⁶Engelhardt Institute of Molecular Biology, Russian Academy of Sciences, Moscow, Russia ²⁷Moscow Institute of Physics and Technology, Dolgoprudny, Russia ²⁸IMPPC, Institute of Predictive and Personalized Medicine of Cancer, Ctra. de Can Ruti, Badalona, Spain ²⁹Institute of Bioengineering, Research Center of Biotechnology, Moscow, Russia ³⁰Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, California, USA ³¹Department of Biochemistry, Ohu University School of Pharmaceutical Sciences, Koriyama, Japan ³²Laboratory of Cell Systems, Institute for Protein Research, Osaka University, Osaka, Japan ³³RIKEN Center for Integrative Medical Sciences, Yokohama, Japan ³⁴Division of Oncology and Pathology, Department of Clinical Sciences, Lund University, Lund, Sweden ³⁵Department of Computer and Information Science, Norwegian University of Science and Technology, Trondheim, Norway ³⁶Department of Clinical Molecular Genetics, School of Pharmacy, Tokyo University of Pharmacy and Life Sciences, Tokyo, Japan ³⁷The Bioinformatics Centre, Department of Biology, and Biotech Research and Innovation Centre (BRIC), University of Copenhagen, Copenhagen, Denmark ³⁸Department of Biochemistry, Nihon University School of Dentistry, Tokyo, Japan ³⁹The SKI Stem Cell Research Facility, The Center for Stem Cell Biology and Developmental Biology Program, Sloan Kettering Institute, New York, New York, United States ⁴⁰Mater Research Institute - University of Queensland, Translational Research Institute, Brisbane, Australia ⁴¹RIKEN Preventive Medicine and Diagnosis Innovation Program, Wako, Japan

Acknowledgments

FANTOM5 was made possible by the following grants: Research Grant for RIKEN Omics Science Center from MEXT to Y.H.; Grant of the Innovative Cell Biology by Innovative Technology (Cell Innovation Program) from the MEXT to Y.H.; Research Grant from MEXT to the RIKEN Center for Life Science Technologies; Research Grant to RIKEN Preventive Medicine and Diagnosis Innovation Program from MEXT to Y.H. K.V.S. and A.S. were supported by the Lundbeck and Novo Nordisk Foundations. A.R.R.F. is supported by a Senior Cancer Research Fellowship from the Cancer Research Trust, funds raised by the MACA Ride to Conquer Cancer, and the Australian Research Council's Discovery Projects funding scheme (DP160101960). Y.A.M. was supported by the Russian Science Foundation, grant 15-14-30002. R.D. was supported by the Russian Science Foundation, grant 14-44-00022. We would like to thank Lucia Schwarzfischer for technical assistance and Norbert Eichner and Gunter Meister for sequencing RACE products. We would also like to thank GeNAS for data production.

References

1. Bartel DP. MicroRNAs: genomics, biogenesis, mechanism, and function. *Cell*. 2004; 116:281–297. [PubMed: 14744438]
2. Shenoy A, Blelloch RH. Regulation of microRNA function in somatic stem cell proliferation and differentiation. *Nat Rev Mol Cell Biol*. 2014; 15:565–576. [PubMed: 25118717]
3. Li M, Izpisua Belmonte JC. Roles for noncoding RNAs in cell-fate determination and regeneration. *Nat Struct Mol Biol*. 2015; 22:2–4. [PubMed: 25565025]

4. Mehta A, Baltimore D. MicroRNAs as regulatory elements in immune system logic. *Nat Rev Immunol.* 2016; 16:279–294. [PubMed: 27121651]
5. Hasuwa H, Ueda J, Ikawa M, Okabe M. miR-200b and miR-429 function in mouse ovulation and are essential for female fertility. *Science.* 2013; 341:71–73. [PubMed: 23765281]
6. Sun K, Lai EC. Adult-specific functions of animal microRNAs. *Nat Rev Genet.* 2013; 14:535–548. [PubMed: 23817310]
7. Mendell JT, Olson EN. MicroRNAs in stress signaling and human disease. *Cell.* 2012; 148:1172–1187. [PubMed: 22424228]
8. Adams BD, Kasinski AL, Slack FJ. Aberrant regulation and function of microRNAs in cancer. *Curr Biol.* 2014; 24:R762–R776. [PubMed: 25137592]
9. Lin S, Gregory RI. MicroRNA biogenesis pathways in cancer. *Nat Rev Cancer.* 2015; 15:321–333. [PubMed: 25998712]
10. Jonas S, Izaurralde E. Towards a molecular understanding of microRNA-mediated gene silencing. *Nat Rev Genet.* 2015; 16:421–433. [PubMed: 26077373]
11. Ha M, Kim VN. Regulation of microRNA biogenesis. *Nat Rev Mol Cell Biol.* 2014; 15:509–524. [PubMed: 25027649]
12. Kozomara A, Griffiths-Jones S. miRBase: annotating high confidence microRNAs using deep sequencing data. *Nucleic Acids Res.* 2014; 42:D68–D73. [PubMed: 24275495]
13. Pritchard CC, Cheng HH, Tewari M. MicroRNA profiling: approaches and considerations. *Nat Rev Genet.* 2012; 13:358–369. [PubMed: 22510765]
14. Chang TC, Pertea M, Lee S, Salzberg SL, Mendell JT. Genome-wide annotation of microRNA primary transcript structures reveals novel regulatory mechanisms. *Genome Res.* 2015; 25:1401–1409. [PubMed: 26290535]
15. Kanamori-Katayama M, et al. Unamplified cap analysis of gene expression on a single-molecule sequencer. *Genome Res.* 2011; 21:1150–1159. [PubMed: 21596820]
16. Takahashi H, Lassmann T, Murata M, Carninci P. 5′ end-centered expression profiling using cap-analysis gene expression and next-generation sequencing. *Nat Protoc.* 2012; 7:542–561. [PubMed: 22362160]
17. Forrest ARR, et al. A promoter-level mammalian expression atlas. *Nature.* 2014; 507:462–470. [PubMed: 24670764]
18. Arner E, et al. Transcribed enhancers lead waves of coordinated transcription in transitioning mammalian cells. *Science.* 2015; 347:1010–1014. [PubMed: 25678556]
19. Fort A, et al. Deep transcriptome profiling of mammalian stem cells supports a regulatory role for retrotransposons in pluripotency maintenance. *Nat Genet.* 2014; 46:558–566. [PubMed: 24777452]
20. Fromm B, et al. A uniform system for the annotation of vertebrate microRNA genes and the evolution of the human microRNAome. *Annu Rev Genet.* 2015; 49:213–242. [PubMed: 26473382]
21. Nepal C, et al. Transcriptional, post-transcriptional and chromatin-associated regulation of pri-miRNAs, pre-miRNAs and moRNAs. *Nucleic Acids Res.* 2016; 44:3070–3081. [PubMed: 26673698]
22. Djebali S, et al. Landscape of transcription in human cells. *Nature.* 2012; 489:101–108. [PubMed: 22955620]
23. Taft RJ, et al. Tiny RNAs associated with transcription start sites in animals. *Nat Genet.* 2009; 41:572–578. [PubMed: 19377478]
24. Westholm JO, Lai EC. Mirtrons: microRNA biogenesis via splicing. *Biochimie.* 2011; 93:1897–1904. [PubMed: 21712066]
25. Matera AG, Terns RM, Terns MP. Non-coding RNAs: lessons from the small nuclear and small nucleolar RNAs. *Nat Rev Mol Cell Biol.* 2007; 8:209–220. [PubMed: 17318225]
26. Friedländer MR, Mackowiak SD, Li N, Chen W, Rajewsky N. miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids Res.* 2012; 40:37–52. [PubMed: 21911355]

27. Londin E, et al. Analysis of 13 cell types reveals evidence for the expression of numerous novel primate- and tissue-specific microRNAs. *Proc Natl Acad Sci USA*. 2015; 112:E1106–E1115. [PubMed: 25713380]
28. Freeman TC, et al. Construction, visualisation, and clustering of transcription networks from microarray expression data. *PLoS Comput Biol*. 2007; 3:2032–2042. [PubMed: 17967053]
29. Ludwig N, et al. Distribution of miRNA expression across human tissues. *Nucleic Acids Res*. 2016; 44:3865–3877. [PubMed: 26921406]
30. Meehan TF, et al. Logical development of the cell ontology. *BMC Bioinformatics*. 2011; 12:6. [PubMed: 21208450]
31. Lizio M, et al. Gateways to the FANTOM5 promoter level mammalian expression atlas. *Genome Biol*. 2015; 16:22. [PubMed: 25723102]
32. Batut P, Dobin A, Plessy C, Carninci P, Gingeras TR. High-fidelity promoter profiling reveals widespread alternative promoter usage and transposon-driven developmental gene expression. *Genome Res*. 2013; 23:169–180. [PubMed: 22936248]
33. Georgakilas G, et al. DIANA-miRGen v3.0: accurate characterization of microRNA promoters and their regulators. *Nucleic Acids Res*. 2016; 44:D190–D195. [PubMed: 26586797]
34. Chien CH, et al. Identifying transcriptional start sites of human microRNAs based on high-throughput sequencing data. *Nucleic Acids Res*. 2011; 39:9345–9356. [PubMed: 21821656]
35. Guo Z, et al. Genome-wide survey of tissue-specific microRNA and transcription factor regulatory networks in 12 tissues. *Sci Rep*. 2014; 4:5150. [PubMed: 24889152]
36. Ding J, et al. Trbp regulates heart function through microRNA-mediated Sox6 repression. *Nat Genet*. 2015; 47:776–783. [PubMed: 26029872]
37. Suzuki H, et al. The transcriptional network that controls growth arrest and differentiation in a human myeloid leukemia cell line. *Nat Genet*. 2009; 41:553–562. [PubMed: 19377474]
38. Landgraf P, et al. A mammalian microRNA expression atlas based on small RNA library sequencing. *Cell*. 2007; 129:1401–1414. [PubMed: 17604727]
39. Carninci P, et al. The transcriptional landscape of the mammalian genome. *Science*. 2005; 309:1559–1563. [PubMed: 16141072]
40. Kapranov P, et al. RNA maps reveal new RNA classes and a possible function for pervasive transcription. *Science*. 2007; 316:1484–1488. [PubMed: 17510325]
41. Meunier J, et al. Birth and expression evolution of mammalian microRNA genes. *Genome Res*. 2013; 23:34–45. [PubMed: 23034410]
42. Francia S, et al. Site-specific DICER and DROSHA RNA products control the DNA-damage response. *Nature*. 2012; 488:231–235. [PubMed: 22722852]
43. Valen E, et al. Biogenic mechanisms and utilization of small RNAs derived from human protein-coding genes. *Nat Struct Mol Biol*. 2011; 18:1075–1082. [PubMed: 21822281]
44. Siepel A, et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res*. 2005; 15:1034–1050. [PubMed: 16024819]
45. Abugessaisa I, et al. FANTOM5 transcriptome catalog of cellular states based on Semantic MediaWiki. *Database (Oxford)*. 2016; 2016:baw105. [PubMed: 27402679]
46. Lassmann T, Hayashizaki Y, Daub CO. TagDust—a program to eliminate artifacts from next generation sequencing data. *Bioinformatics*. 2009; 25:2839–2840. [PubMed: 19737799]
47. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*. 2009; 25:1754–1760. [PubMed: 19451168]
48. De Hoon MJL, et al. Cross-mapping and the identification of editing sites in mature microRNAs in high-throughput sequencing libraries. *Genome Res*. 2010; 20:257–264. [PubMed: 20051556]
49. Lorenz R, et al. ViennaRNA Package 2.0. *Algorithms Mol Biol*. 2011; 6:26. [PubMed: 22115189]
50. Vijayan D, Radford KJ, Beckhouse AG, Ashman RB, Wells CA. Mincle polarizes human monocyte and neutrophil responses to *Candida albicans*. *Immunol Cell Biol*. 2012; 90:889–895. [PubMed: 22641025]
51. Busk PK. A tool for design of primers for microRNA-specific quantitative RT-qPCR. *BMC Bioinformatics*. 2014; 15:29. [PubMed: 24472427]

52. Schwarzenbach H, Da Silva AM, Calin G, Pantel K. Data normalization strategies for microRNA quantification. *Clin Chem*. 2015; 61:1333–1342. [PubMed: 26408530]
53. Nagpal N, Kulshreshtha R. miR-191: an emerging player in disease biology. *Front Genet*. 2014; 5:99. [PubMed: 24795757]
54. Moon HG, Yang J, Zheng Y, Jin Y. miR-15a/16 regulates macrophage phagocytosis after bacterial infection. *J Immunol*. 2014; 193:4558–4567. [PubMed: 25261473]
55. Vinod M, et al. miR-206 controls LXRA expression and promotes LXR-mediated cholesterol efflux in macrophages. *Biochim Biophys Acta*. 2014; 1841:827–835. [PubMed: 24603323]
56. Cobos Jiménez V, et al. Next-generation sequencing of microRNAs in primary human polarized macrophages. *Genom Data*. 2014; 2:181–183. [PubMed: 26484091]
57. Zhang L, et al. miR-153 supports colorectal cancer progression via pleiotropic effects that enhance invasion and chemotherapeutic resistance. *Cancer Res*. 2013; 73:6435–6447. [PubMed: 23950211]
58. Srivastava SK, et al. MicroRNA-345 induces apoptosis in pancreatic cancer cells through potentiation of caspase-dependent and -independent pathways. *Br J Cancer*. 2015; 113:660–668. [PubMed: 26247574]
59. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010; 26:139–140. [PubMed: 19910308]
60. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol*. 2010; 11:R106. [PubMed: 20979621]
61. Harrow J, et al. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res*. 2012; 22:1760–1774. [PubMed: 22955987]
62. Brown GR, et al. Gene: a gene-centered information resource at NCBI. *Nucleic Acids Res*. 2015; 43:D36–D42. [PubMed: 2535515]
63. Kent WJ. BLAT—the BLAST-Like Alignment Tool. *Genome Res*. 2002; 12:656–664. [PubMed: 11932250]
64. Karolchik D, et al. The UCSC Genome Browser database: 2014 update. *Nucleic Acids Res*. 2014; 42:D764–D770. [PubMed: 24270787]
65. Notredame C, Higgins DG, Heringa J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. *J Mol Biol*. 2000; 302:205–217. [PubMed: 10964570]
66. Arnold P, Erb I, Pachkov M, Molina N, Van Nimwegen E. MotEvo: integrated Bayesian probabilistic methods for inferring regulatory sites and motifs on multiple alignments of DNA sequences. *Bioinformatics*. 2012; 28:487–494. [PubMed: 22334039]
67. Pachkov M, Balwierz PJ, Arnold P, Ozonov E, Van Nimwegen E. SwissRegulon, a database of genome-wide annotations of regulatory sites: recent updates. *Nucleic Acids Res*. 2013; 41:D214–D220. [PubMed: 23180783]

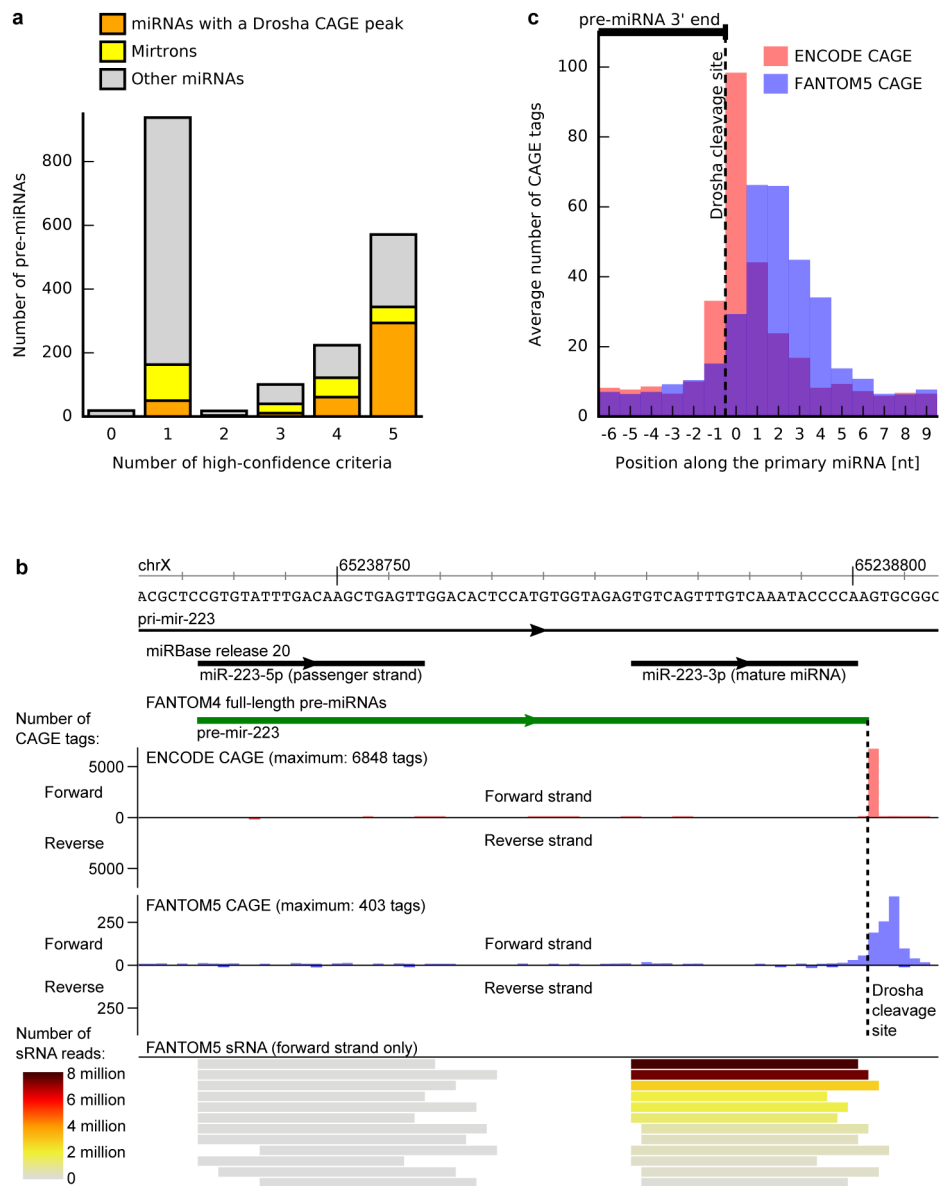


Figure 1. Selection of robust miRNAs and Drosha CAGE peak analysis

(a) Number of miRBase high-confidence criteria (Table 2) satisfied by human pre-miRNAs annotated in miRBase. Pre-miRNAs with a statistically significant ($P < 0.05$) Drosha CAGE peak are shown in orange; mirtrons are shown in yellow. (b) Genomic locus of mir-223 in human with the total number of FANTOM5 (blue) and ENCODE (red) CAGE tags as a function of the genomic position of their 5' end, showing a Drosha CAGE peak at the 3' end of the pre-miRNA. FANTOM5 sRNA reads are shown at the bottom, colored by their read count as defined by the color key. The exact 5' and 3' ends of the pre-miRNA were determined from FANTOM4 full-length sequencing data²³. (c) Number of CAGE tags as a function of their starting position relative to the 3' end of the pre-miRNA, averaged across human pre-miRNAs in the robust set ($n = 795$). The 3' end of the pre-miRNA was selected as the 3' end of the most prevalent sRNA on the 3' arm of the pre-miRNA in the

FANTOM5 sRNA data, with the position indicated as zero corresponding to the first nucleotide downstream of the 3' end of the pre-miRNA.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

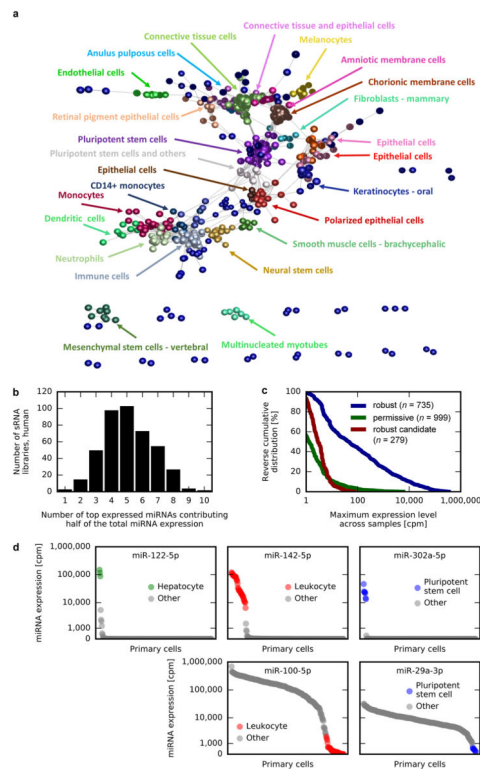


Figure 2. Expression profile and cell ontology analysis of mature miRNAs

(a) Miru²⁸ visualization of FANTOM5 primary cell samples based on their expression profile of robust mature miRNAs. (b) Number of most abundant miRNAs contributing at least 50% of the total miRNA expression in each human sRNA library in FANTOM5. (c) Reverse cumulative distribution of the maximum expression across the FANTOM5 samples of human miRNAs in the robust set, permissive set, and robust candidate set. (d) Examples of miRNAs enriched or depleted in specific primary cell samples. Expression of miR-122-5p, miR-142-5p, and miR-302a-5p was enriched in hepatocytes, leukocytes, and pluripotent stem cells, respectively; miR-100-5p and miR-29a-3p were broadly expressed, but depleted in leukocytes and pluripotent stem cells, respectively. cpm., counts per million.

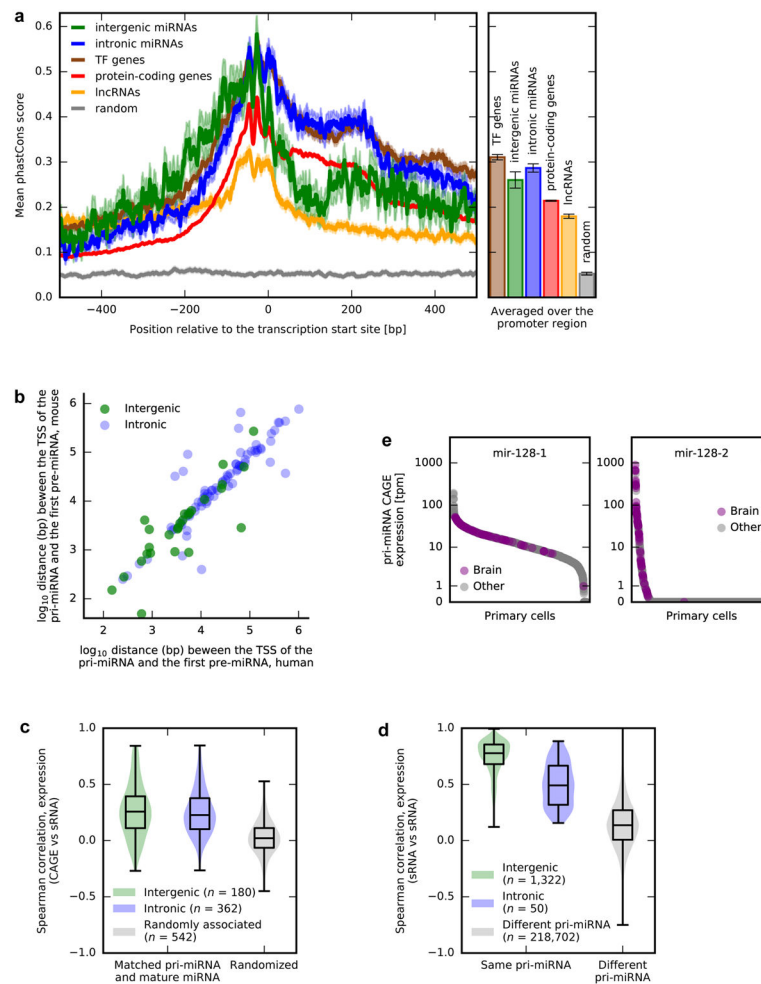


Figure 3. Analysis of the curated miRNA promoters of miRNAs in the robust set
(a) (left panel) Sequence conservation of the human genome, evaluated as the average phastCons⁴⁴ score, in the promoter region of non-coding pri-miRNAs (containing intergenic mature miRNAs; $n = 132$), coding pri-miRNAs (containing intronic mature miRNAs; $n = 415$), transcription factor (TF)-coding transcripts ($n = 1,651$), other protein-coding transcripts ($n = 15,350$), and long non-coding RNAs (lnc; $n = 1,461$). The sequence conservation of randomly selected genome regions is shown in gray. The shaded area corresponds to one standard deviation in the estimated mean phastCons score. (right panel) The average sequence conservation at promoter regions of miRNAs was higher than at the promoter regions of non-TF protein-coding genes (Mann-Whitney $P = 2 \times 10^{-16}$, two-sided) and of long non-coding RNAs (Mann-Whitney $P = 1 \times 10^{-35}$, two-sided). Error bars correspond to one standard deviation in the estimated mean phastCons score. **(b)** Distance between the transcription start site (TSS) of the pri-miRNA and the 5' end of the first pre-miRNA is highly correlated between human and mouse both for coding and non-coding pri-miRNAs, suggesting strong conservation of the genomic extent of pri-miRNAs. **(c)** Spearman correlation between the expression level of pri-miRNAs, as measured by CAGE, and mature miRNAs, as measured by sRNA sequencing, compared to a background distribution consisting of correlations between randomly paired pri-miRNAs and mature

miRNAs. Correlations for polycistronic pri-miRNAs were averaged across the mature miRNAs. **(d)** Spearman correlation in expression level between mature miRNAs originating from the same pri-miRNA, compared to a background distribution consisting of correlations between mature miRNAs originating from different pri-miRNAs. **(e)** Cell type-dependent expression of miRNA paralogs. While mir-128-1 was broadly expressed across most primary cell samples in FANTOM5, its paralog mir-128-2 was enriched in brain samples. tpm, tags per million. **(c-d)** The box extends from the lower to the upper quartile, with the center line at the median; the whiskers indicate the full range of the data.

Table 1

Human sRNA data sets analyzed in this study.

Origin	Data collection	Number of samples	Number of cell types
Primary cells	FANTOM5	293	119
	Fort <i>et al.</i> ¹⁹	6	
ES cells	Fort <i>et al.</i> ¹⁹	6	1
iPS cells	Fort <i>et al.</i> ¹⁹	6	1
Tissues	FANTOM5	6	4 tissues
Time courses	FANTOM5	87	6 time courses

Total number of sequenced reads: 1,519,621,910.

ES, embryonic stem; iPS, induced pluripotent stem.

Table 2

The miRBase high-confidence criteria¹². As a meaningful evaluation of the second, third, and fourth criteria relies on accurate knowledge of the position and extent of the mature miRNA on both strands of the pre-miRNA, we evaluated these three criteria only if the first criterion was satisfied.

1.	≥ 10 tags on each arm of the pre-miRNA, or ≥ 100 tags on one arm of the pre-miRNA, with ≥ 5 tags on the other arm
2.	$\geq 50\%$ of the tags on each arm of the pre-miRNA have the same 5' end
3.	0–4 nt overhang at the mature 3' end on each arm
4.	$\geq 60\%$ of nucleotides of the mature sequence on each arm are base-paired
5.	$G < -0.2$ kcal/mole/nucleotide