

# 1 **Comparing miRNA structure of mirtrons and non-mirtrons**

2 Igor I. Titov<sup>1,2\*</sup>, Pavel S. Vorozheykin<sup>2</sup>

3 <sup>1</sup>Federal State Budget Scientific Institution “The Federal Research Center Institute of Cytology  
4 and Genetics of Siberian Branch of the Russian Academy of Sciences“, Novosibirsk, Russia

5 <sup>2</sup>Novosibirsk State University, Novosibirsk, Russia

6 \*Corresponding author

7 E-mail: titov@bionet.nsc.ru

8

## 9 **Abstract**

### 10 **Background**

11 MicroRNAs proceeds through the different canonical and non-canonical pathways; the most  
12 frequent of the non-canonical ones is the splicing-dependent biogenesis of mirtrons. We compare  
13 the mirtrons and non-mirtrons of human and mouse to explore how their maturation appears in  
14 the precursor structure around the miRNA.

### 15 **Results**

16 We found the coherence of the overhang lengths what indicates the dependence between the  
17 cleavage sites. To explain this dependence we suggest the 2-lever model of the Dicer structure  
18 that couples the imprecisions in Drosha and Dicer. Considering the secondary structure of all  
19 animal pre-miRNAs we confirmed that single-stranded nucleotides tend to be located near the  
20 miRNA boundaries and in its center and are characterized by a higher mutation rate. The 5' end  
21 of the canonical 5' miRNA approaches the nearest single-stranded nucleotides what suggests the  
22 extension of the loop-counting rule from the Dicer to the Drosha cleavage site. A typical  
23 structure of the annotated mirtron pre-miRNAs differs from the canonical pre-miRNA structure  
24 and possesses the 1- and 2nt hanging ends at the hairpin base. Together with the excessive  
25 variability of the mirtron Dicer cleavage site (that could be partially explained by guanine at its

1 ends inherited from splicing) this is one more evidence for the 2-lever model. In contrast with the  
2 canonical miRNAs the mirtrons have higher snp densities and their pre-miRNAs are inversely  
3 associated with diseases. Therefore we supported the view that mirtrons are under positive  
4 selection while canonical miRNAs are under negative one and we suggested that mirtrons are an  
5 intrinsic source of silencing variability which produces the disease-promoting variants. Finally,  
6 we considered the interference of the pre-miRNA structure and the U2snRNA:pre-mRNA  
7 basepairing. We analyzed the location of the branchpoints and found that mirtron structure tends  
8 to expose the branchpoint site what suggests that the mirtrons can readily evolve from occasional  
9 hairpins in the immediate neighbourhood of the 3' splice site.

## 10 **Conclusion**

11 The miRNA biogenesis manifests itself in the footprints of the secondary structure. Close  
12 inspection of these structural properties can help to uncover new pathways of miRNA biogenesis  
13 and to refine the known miRNA data, in particular, new non-canonical miRNAs may be  
14 predicted or the known miRNAs can be re-classified.

15

## 16 **Keywords**

17 microRNA, miRNA, mirtron, biogenesis, secondary structure, Dicer structure, overhangs,  
18 branchpoint, splicing

19

## 20 **Background**

21 Canonical pathway of animal miRNA begins with transcription. RNA polymerase II (or  
22 polymerase III for some miRNAs) creates long primary transcript (pri-miRNA) which contains  
23 one or more hairpins (pre-miRNAs), poly(A) tail and 7-methylguanosine cap [1, 2]. MiRNA  
24 genes are dispersed in various genomic locations (intronic, exonic or intergenic regions) and can  
25 be transcribed independently or as a part of other host genes [3-5]. A cluster brings together

1 miRNAs with inter-miRNA distance up to 10kb and can form a polycistronic transcriptional unit  
2 (for example, mir-100/let-7/mir-125 and mir-71/mir-2 clusters) [6-8]. MiRNAs can be located in  
3 both DNA strands (for example, hsa-miR-3120 and hsa-miR-214, dme-miR-iab-4): although  
4 these miRNAs are close to each other, they can be regulated post-transcriptionally either united  
5 or independent [9, 10].

6 After the transcription, animal pri-miRNAs are cleaved by the Microprocessor complex of the  
7 RNase III enzyme Drosha and its co-factor DGCR8 [3]. The complex releases pre-miRNA  
8 hairpin by cropping the stem-loop [11-13]. This step can be regulated by a variety of ways: in  
9 some of them proteins are recruited to protein-protein interactions, in others the pri-miRNA  
10 primary and secondary structures are involved in RNA-protein or RNA-RNA bindings.

11 Further, the Drosha product is moved from the nucleus to the cytoplasm by the protein  
12 Exportin-5 (EXP5) and the cofactor Ran-GTP [14]. Some other proteins (for example, XPO1)  
13 can transport non-canonical pre-miRNAs [15]. EXP5 does not only transfer the precursors, but  
14 also prevents them from degradation [16].

15 In the cytoplasm, the pre-miRNA must be cleaved by RNase III enzyme Dicer near the  
16 terminal loop. The cleavage releases a double stranded miRNA duplex with typical 2nt 3'  
17 overhangs [17]. Usually, animal Dicer contains the following domains: helicase, PAZ, dsRNA  
18 binding and two RNase III (A and B) domains [18]. Each of these domains are involved in the  
19 miRNA maturation process. The helicase domain promotes the pre-miRNA recognition by  
20 interacting with the terminal loop and facilitates the processing [19]. The PAZ domain identifies  
21 the precursor's termini and binds to them. Each of the two RNase III domains cuts one of the two  
22 pre-miRNA strands and releases the miRNA duplex from the terminal loop [20, 21].

23 After the Dicer has produced the miRNA duplex, a miRNAs-induced silencing complex  
24 (miRISC) is formed and targets mRNAs [22-24] or non-coding RNAs [25-27].

25 In addition to the canonical miRNA biogenesis described above, another pathways can  
26 generate miRNAs in a Drosha- and/or Dicer-independent manner [28, 29]. Most of the non-

1 canonical miRNAs are mirtrons which bypass the Drosha cleavage step and are derived through  
2 the mRNA splicing, the lariat debranching and refolding into a canonical-like stem-loop  
3 structure [30-34]. If this stem-loop contains extra-nucleotides at 5' or 3' ends (the so-called  
4 “tailed” mirtron), they are trimmed by exonucleases to gain an appropriate structure for the  
5 exportin complex. From this time on, the processing goes on the canonical pathway. At present,  
6 the hundreds of mirtrons have been found [32], however the features of non-canonical  
7 maturation are still poorly studied in contrast with the canonical one.

8 Mammalian mir-1225 and mir-1228, initially predicted as mirtrons, are actually splicing-  
9 independent [35]. Moreover, biogenesis of these miRNAs does not require the most of the  
10 canonical components (DGCR8, Dicer, Exportin-5 or Ago2) but still involves Drosha [36, 37].  
11 This class of miRNAs, termed “simtrons” (splicing-independent mirtron-like miRNAs), reveals a  
12 new pathway of small regulatory RNA production [36, 37]. Another Dicer-independent pathway  
13 is observed for the mir-451 family, this pathway involves the catalytic activity of the Ago2  
14 protein [38- 40]. Drosha generates pre-mir-451 with ~18-nt stem which is too short to be  
15 processed by Dicer. Therefore the pre-miRNAs are processed by Ago2 which cleaves the hairpin  
16 in the middle of its 3' strand and yields a ~30nt long RNA product [38- 40]. Then poly(A)-  
17 specific ribonuclease PARN trims the 3' RNA end to release the mature 5' miRNA [41].

18 On each step of miRNA maturation, the biogenesis leaves its footprints as the specific pre-  
19 miRNA and miRNA features. The canonical miRNAs are usually located near the terminal loop  
20 of the pre-miRNA hairpin. Simultaneously some non-canonical miRNAs (e.g. originated from  
21 simtrons mir-1225 and mir-1228) are distant from the terminal loop. The mir-451 family can be  
22 processed by the canonical pathway as well as by the non-canonical one in which Drosha  
23 produces a short hairpin with the miRNA that overlaps within the terminal loop. Based on the  
24 pri-/pre-miRNA structural properties, new non-canonical miRNAs may be predicted or the  
25 known miRNAs can be re-classified. Also, close inspection of these characteristics can help to

1 uncover the new pathways of the miRNA biogenesis and to discover the errors in annotated  
2 miRNA data.

3 It is commonly believed [42] that miRNA genes have been either evolved from random  
4 hairpins in intergenic regions or in intronic regions of protein-coding genes or duplicated the  
5 miRNA genes and the transposable elements (e.g. the fraction of the human TE-derived miRNAs  
6 in miRBase had been constantly growing [43]). Most of new miRNAs disappeared over time  
7 while the survived ones adapted and then came under purifying selection like the old miRNAs  
8 until they could start an another cycle of adaptive-conservative evolution in other tissues [44].  
9 Many of these new miRNAs are mirtrons, they have been often evolved in clade- and species-  
10 specific ways and more quickly than the canonical miRNAs [45, 46].

11 In this paper we consider the structural properties of the animal miRNAs and compare  
12 mirtrons with non-mirtrons, most of the latter are the canonical miRNAs. First, we study the  
13 miRNA pair layout which shows itself in overhang lengths. Second, we investigate the distances  
14 from the miRNA ends to the nearest single-stranded nucleotide. Then we inspect the loop  
15 positional frequencies in the miRNA and its flanks and correlate these frequencies with the  
16 mutation rate. Next, we study SNP density in miRNA and its flanks. Finally, we consider how  
17 the branchpoints are located within the mirtron pre-miRNAs.

18 Our observations support the current view that RNA secondary structure plays a crucial role in  
19 miRNA maturation and exemplify how the biogenesis peculiarities become apparent in this  
20 structure. A lot of miRNA/pre-miRNA prediction methods use the RNA secondary structure [47-  
21 54], so our results can be useful for further improvement of the existing methods. The excessive  
22 SNP density and the branchpoint locations within mirtron precursors demonstrate that mirtrons  
23 represent new miRNAs which could be easily recruited from introns. The further inspection of  
24 the mirtron branchpoints can help in better understanding the role of the secondary structure in  
25 splicing.

26

## 1 **Methods**

2 The sequences and structures of the pre-miRNAs and miRNAs were downloaded from the  
3 miRBase database (release 21.0) [55]. There are 15 731 unique animal pre-miRNAs which  
4 contain 22 603 experimentally validated mature miRNAs approximately equally in both arms of  
5 the precursors. We excluded few pre-miRNAs with non-canonical nucleotides and with more  
6 than two annotated miRNAs.

7 We selected those mirtrons which are simultaneously presented both in miRBase-21.0 and in  
8 the paper [56]. These data contain 464 human and mouse mirtron pre-miRNAs, while a number  
9 of non-mirtron human and mouse pre-miRNAs is 2438.

10 The revisited miRNA sequences were taken from the miRBase-21.0 whose identifiers are  
11 simultaneously presented in [57]. This set contains more than one thousand animal pre-miRNAs.

12 We used the SNP database miRNASNP-2.0 (based on miRBase-19.0 and dbSNP137) [58] to  
13 calculate the SNP densities in human miRNA genes and their flanks. The dataset contains both  
14 common (minor allele frequency  $> 0.01$ ) and rare SNP variants. The SNP density was defined  
15 as:  $N_{snp} \times 1000/L$ , where  $N_{snp}$  was the number of the SNPs in the RNA region,  $L$  was the length  
16 of the region (seed, miRNA excluding seed, pre-miRNA excluding miRNA). SNP occurrence  
17 per sequence for disease and non-disease human pre-miRNAs was calculated as in [59] and was  
18 based on miRNA associated diseases from [60] and on miRNASNP-2.0 [58].

19 Branchpoint data of human and mouse introns were taken from the supplemental materials of  
20 the paper [61].

21 Data of the animal nucleotide substitutions were taken from [62].

22 The unpaired nucleotide frequency (UNF) for each RNA position was calculated as the  
23 portion of miRNAs that had a single-stranded nucleotide at the position.

24 The distance between miRNA end and its nearest single-stranded nucleotide was calculated as  
25 the minimal number of the nucleotides between the miRNA boundary and the single-stranded  
26 region in the same miRNA strand.

1

## 2 **Results and Discussion**

### 3 **Overhang lengths**

4 The diversity of the miRNA cleavage site leads to overhang variety, therefore they can shed  
5 new light on the nature and mechanism of the cleavage process. The overhangs are the miRNA  
6 ends hanging from its miRNA duplex thus reflecting the miRNAs disposition. Both Dicer and  
7 Drosha cut the miRNA precursor, especially the 3' variable miRNA ends, in a number of  
8 neighbouring positions and can form other than canonical 2nt overhangs. Each cleavage variant  
9 produces its own version of the miRNA duplex and in some miRBase records the variant with  
10 2nt overhangs is not the most observable.

11 The Dicer and Drosha interactions with pre-miRNA, especially sensitivity of their RNase  
12 domains (RIIIA and RIIB) to RNA sequence, defines the miRNA ends. These domains prefer to  
13 generate the U-ended miRNAs as the main fraction [63]. The G-ended miRNAs rarely occur and  
14 the G-avoiding generates the atypical 1nt and 3nt overhangs [64]. For the “homogeneous”  
15 cleavage (as it was defined by [63]) the overhang shortening appears as a 1nt context shift at the  
16 3' miRNA ends (compare figure 2A and figure 2C of the second-most frequent miRNA fraction  
17 in [63]). Sometimes these shortened miRNAs are found in miRBase in another species: compare,  
18 for example, tgu-let-7b, aca-let-7b and mml-let-7b miRNAs (Additional file 1). Occurrence of  
19 the A/G at the neighbourhood of miRNA ends could lead to the heterogenous cleavage, i.e. to  
20 levelling of the miRNA fractions (figure 2 from [63]). Based on the works of Starega-Roslan  
21 et.al. [63- 66], one can conclude that the overhang lengths are sequence-dependent; they are also  
22 structure-dependent as it was observed in [67] where the sliding (bulge) loops induced cleavage  
23 heterogeneity.

24 Since not only the sequence but also the structure of the pre-miRNA can trigger this cleavage  
25 heterogeneity, we measure the overhang lengths by a number of excessive nucleotides beyond

1 the closing pair of the miRNA duplex regardless of its structural state, single-stranded or double-  
2 stranded, rather than a length of hanging end.

3 To estimate the occurrence of the atypical overhangs we study the overhang lengths,  
4 considering the mirtrons as a separate miRNA class. Unlike the canonical miRNAs, the mirtrons  
5 use splicing to bypass Drosha cleavage. The mirtron database consists mainly of human and  
6 mouse mirtrons [56], therefore we consider four miRNA sets: animal miRNAs (1), animal  
7 miRNAs without human and mouse ones (2), human and mouse non-mirtrons (3) and mirtrons  
8 (4). The third set contains only few numbers of known non-canonical miRNAs which could not  
9 significantly influence on.

10

11 **Figure 1. The overhang lengths of miRNA duplexes.** The frequency of the overhang lengths of  
12 miRNA duplexes: animal miRNAs without human/mouse ones (A), human and mouse non-  
13 mirtrons (B) and mirtrons(C). The overhang lengths occurrence of both cleavage sites for animal  
14 miRNA duplexes; in each quarter-square the miRBase pre-miRNA structure which leads to the  
15 corresponding overhang types is schematically shown (D). Negative values correspond to an  
16 atypical 5' overhangs. The long overhangs on the panel D correspond to structure prediction  
17 errors and are described further in the text. These overhangs are not shown on the panels A and  
18 C. In mirtron case the splicing overhangs are considered instead of the Drosha ones.

19

20 Figures 1A-1C show the overhang length distributions of the miRNA duplexes from three of  
21 four sets. The data on figure 1A includes the canonical miRNAs, Drosha-independent mirtrons  
22 and a small number of other non-canonical miRNAs, e.g. Dicer-independent mir-451 family  
23 [29], simtrons [37], etc. As we see, figure 1A does not significantly differ from the figure 1B  
24 which represents the distributions of human and mouse non-mirtrons. The data for all animal  
25 miRNAs are also similar to figure 1A and figure 1B and therefore are not shown here. Figure 1C  
26 displays the corresponding distributions of the most abundant non-canonical class, mirtrons.



1 The canonical overhang lengths for both cleavage sites are well-known to be equal to 2nt,  
2 what is actually observed on figures 1A and 1B where all the length distributions peak at 2nt.  
3 The overhang distributions of the Dicer and Drosha cleavage sites are similar and asymmetrical  
4 (figures 1A and 1B). The overhangs are more readily shortening what could be explained either  
5 by more frequent exosome cutting of the 3' miRNA end than of the 5' one [68] and/or by the 3'  
6 cleavage site shifting inside the duplex.

7 The Drosha overhang distribution closely matches the Dicer's one (figure 1B) what supports  
8 the observations that Drosha and Dicer process the canonical miRNAs in a similar manner. In  
9 the mirtron case the splicing replaces the Drosha step, but the overhang statistics are surprisingly  
10 changed for both cleavage sites (figures 1B and 1C). First, splicing overhang distribution is  
11 shifted relatively to the Dicer distribution (figure 1C) what can reflect the exonuclease trimming  
12 of 5'-tailed mirtrons which are most frequently observed [56]. Second, although the same  
13 complex cleaves both mirtrons and non-mirtrons at the Dicer site, the overhang distributions  
14 differ (figures 1B and 1C): mirtrons distribution blurs what suggests lower Dicer cleavage  
15 precision, presumably due to the dependence of Dicer result on the output of splicing and further  
16 exonuclease editing.

17 The introns usually have sequence conservations at the 5' end (GU) and at the 3' end (AG).  
18 While the majority of mirtrons are tailed and lose one of the intron ends during the exosome  
19 cutting, the remaining end can contribute to the Dicer cleavage heterogeneity as it follows from  
20 the figure 2 in [63].

21 If the Dicer heterogeneity for mirtrons is caused not only by guanine at their ends, but also by  
22 a heterogeneity of the overhang lengths of splice site, such a dependence should appear in a  
23 coordinated variation of Dicer and Drosha overhangs of canonical miRNAs. This contradicts to  
24 the fact that Drosha and Dicer cleave independently and to clear this discrepancy up we plotted  
25 the dependence of the overhang lengths for both cleavage sites of animal miRNAs (figure 1D).  
26 Indeed, we observe significant linear dependence ( $\rho=0.338$ ,  $P = 2.42 \times 10^{-183}$ , Spearman's rank

1 correlation test) of the Dicer and Drosha overhang lengths (figure 1D). This dependence arises  
2 due to the several reasons. First, the way the pre-miRNA structure is formed, it excludes big  
3 bulge loops within miRNA duplex and, consequently, the pairs of long overhangs of the opposite  
4 sign. Second, the pairs of long overhangs of the same sign are observed due to the incorrect  
5 prediction of terminal loop (Additional file 2) or to the presence of false miRNAs in the  
6 miRBase. And the last possibility is the guanine avoiding at the first position of 5' end of 3'  
7 miRNA [63, 64]. To exclude these three reasons, we further considered only the canonical  
8 overhangs and the overhangs with minimal 1nt deviations from them. The first two reasons  
9 disappeared due to the near-canonical overhang lengths, and the last reason (at least as a main  
10 factor) was excluded after comparing guanine frequencies in three neighbouring positions at the  
11 miRNA boundary (see table 1 in Additional file 3). The remained duplexes with these weakly  
12 varying overhang lengths compose the most part (67.7%) of all duplexes and their overhangs still  
13 significantly correlate ( $\rho=0.139$ ,  $P = 2.2 \times 10^{-21}$ , Spearman's rank correlation test).

14 The overhang lengths interdependence can be also verified in biochemical studies (for  
15 example [69]). Unfortunately, their paper does not provide the data on joint occurrence of  
16 miRNA/miRNA\* and therefore can not be used to test our hypothesis.

17 To understand the nature of this correlation, as the null hypothesis we considered the 2-  
18 parametric model of independent overhang lengths and fitted the length frequencies (table 2,  
19 Additional file 3). For the both (Dicer and Drosha) cleavage sites the long (short) overhangs are  
20 observed about 7 (4) times less often than the canonical ones (table 3, Additional file 3) as it was  
21 already seen in figure 1A and figure 1B This model describes well all length frequencies except  
22 for the 1nt/1nt and 3nt/3nt length pairs which are observed twice as often as expected. So, these  
23 pairs are what induce the length correlation. Moreover, this correlation is robust to definition of  
24 the overhang length (Additional file 4).

25 We suggest that this correlation reflects the organization of the pre-miRNA/Dicer complex.  
26 The pre-miRNA/Dicer complex consists of the sub-units that move as a whole, what appears as a

1 collective movements of large scale around a hinges. The PAZ domain is the main moving  
2 domain for the Dicer [70] this domain adapts to the pre-miRNA ends. We speculate that among  
3 all possible collective movements pre-miRNA/Dicer complex undergoes smaller scale  
4 movements of 2-lever type which are responsible for the coherence trend of the overhang lengths  
5 (Additional file 5). The tips of two levers bind to the pre-miRNA ends in the PAZ domain. The  
6 another two tips of the levers are located in the RNase IIIA and RIIIB determining the distance  
7 between the cleavage sites. As a result, the tips of the levers can close in and move away in  
8 concert what leads to such a number of states of the cleavage complex that the Dicer overhang  
9 tends to vary cooperatively with the Drosha one.

10 Unfolding this 2-lever model we note that these levers should differ in their rigidity. One  
11 lever (associated with Dicer RNase IIIA) is connected with the 5' ends of the miRNAs and  
12 tightly bound to RNA which, as well as the connector helix, ensures the lever rigidity and  
13 manifests itself in less variability of the 5' end of the 3' miRNA and in G-avoiding on both 5'  
14 ends. In contrast, the other lever (associated with Dicer RNase RIIIB) is soft and its free  
15 movement forms the cleavage variability and the variety of the overhang lengths.

16 Another evidence in favor of the lever mechanism is the more pronounced heterogeneity of  
17 the Dicer cleavage site for mirtrons in which the overhang lengths of the splicing site are more  
18 variable and the ends are often freely hanging (see next two sections of this paper) thus forming  
19 different spatial distances between each other.

20

## 21 **miRNA end distance to the nearest single-stranded region**

22 The pre-miRNA secondary structure, as well as the nucleotide sequence, can influence the  
23 miRNA boundaries. The Dicer cleavage depends either on the stem size or the terminal loop [71]  
24 and its precision (i.e., the fraction of the most probable miRNA) also fulfills the so-called “loop-  
25 counting rule”: Dicer cleaves precisely at 2nt distance to any upstream loop, in other cases Dicer  
26 produces variable 5' end of the 3' miRNA [69].

1 Gu and co-authors were focused on the Dicer cleavage site. We inspect how the loop-counting  
2 rule makes itself evident in the distance between miRNA end and the single-stranded regions and  
3 answer the question “Is there something like the loop-counting rule for the Drosha cleavage  
4 site?”. Expecting structural difference between canonical and non-canonical pre-miRNAs we  
5 explore separately mirtrons and non-mirtrons.

6 As one can see on figure 2A the loop-counting rule appears as a pronounced peak at 2nt  
7 which contains approximately 40% of animal miRNAs (excluding human and mouse ones). The  
8 remaining cases are partially referred to the incorrect prediction of the pre-miRNA secondary  
9 structures (Additional file 2) and to the fact that the structure is locally unstable in the loop  
10 neighbourhood. The blue peak (5' end of the 5' miRNA) is even more pronounced, this miRNA  
11 end is located immediately before the single-stranded region (figure 2A). This suggests that the  
12 loop-counting rule for the Drosha cleavage site exists as well.

13

14 **Figure 2. Distance between miRNA end and its nearest single-stranded nucleotide in the**  
15 **same miRNA strand.** Considered are only those miRNA ends in which their terminal nucleotide  
16 is double-stranded. The frequencies of the 5' miRNA ends are shown on panels A and C. The  
17 frequencies of the 3' miRNA ends are shown on panels B and D. The data are presented for 5'  
18 and 3' miRNA sequences separately: for animal miRNAs excluding human and mouse ones (A  
19 and B) and for the human and mouse mirtrons and non-mirtrons (C and D). The positive values  
20 correspond to the distances to the nearest single-stranded nucleotide outside the miRNA. The  
21 negative values are the numbers of nucleotides that must be cut off from the miRNA to reach the  
22 nearest loop in the miRNA. The distance 0 is observed for those miRNA ends that are exactly at  
23 the boundary of the single-stranded region. The distance frequencies for all animal miRNAs (not  
24 shown) are the sum of results for all datasets and are only slightly different from the observations  
25 A and B.

26

1        Figure 2B shows the distance frequencies to the 3' miRNA end. The broader distribution on  
2        figure 2B comparing with the figure 2A suggests that the cleavage complex for the 3' miRNA  
3        end is less accurate than for the 5' end as it was already proposed by [63]. Both cleavage sites are  
4        processed by Drosha and Dicer RNase domains (RIIIA and RIIB). The RIIIA domain processes  
5        the 3' miRNA while the RIIB domain handles the 5' one. Therefore, the RNA site affects the  
6        cleavage precision to a greater extent than the particular RNase domain. Specifically, the 5' end  
7        is controlled by well-defined neighbouring structure (figure 2A) and nucleotide sequence [63]. In  
8        contrast, the 3' end may be determined mostly by the overhang length (or distance between  
9        RNase IIIA/IIIB cleavage sites according to the lever model) thus providing the required size of  
10       the miRNA.

11       Some of these animal miRNAs are mirtrons which are though rare but the most abundant non-  
12       canonical miRNAs. To reveal the structural difference between canonical and non-canonical  
13       miRNAs we consider miRNAs of human and mouse where mirtrons are better identified.

14       The frequency distributions near the hairpin terminal loop are similar for mirtrons and non-  
15       mirtrons (red and pink bars on figure 2C, blue and light blue bars on figure 2D) what reflects the  
16       fact that the Dicer processes the both classes in the same way. In contrast, the mirtron and non-  
17       mirtron frequency distributions strongly differ for another cleavage site where the different  
18       processing complexes cleave the RNA molecules (red and pink bars on figure 2D, blue and light  
19       blue bars on figure 2C). More pronounced peaks of mirtrons stems from the observation that  
20       their ends at the hairpin base are located at 0-1nt distance from the single-stranded region  
21       (figures 2C and 2D). Some of these mirtron pre-miRNAs are not produced immediately by  
22       splicing, but are rather derived by further cropping the single-stranded regions by exonucleases.

23       Thus the pre-miRNA secondary structure is an important factor for precise recognition of both  
24       miRNA ends. In particular, the loop-counting rule [69] could be extended to the Drosha cleavage  
25       site. The structural signature near the hairpin base of the mirtron pre-miRNAs is so clearly  
26       defined that may be useful for their validation.

1

## 2 **The unpaired nucleotide frequency across miRNA**

3 As we have seen above, loop position is an important factor of miRNA processing. Therefore,  
4 we consider the unpaired nucleotide frequencies (UNFs) within miRNA and how these  
5 frequencies relate to nucleotide substitutions. Going along the secondary structure of miRNA  
6 classes we compare mirtron and non-mirtron unpaired nucleotide frequencies (UNFs) within  
7 miRNA and its nearest neighbourhood (figures 2C and 2D).

8 Figure 3A shows that the UNF varies across the miRNA sequence. The miRNA positions can  
9 be roughly divided into two groups (figures 3A and 3B). The first group, loop-rare positions  
10 (grey region), contains the seed region (positions 2-8) and the additional binding site (positions  
11 13-16) with few nucleotides downstream (positions 17-19). The UNFs of these regions are close  
12 because they are typically paired to each other in the opposite strands of the miRNA duplex. The  
13 loop-frequent positions (white region) fall into miRNA center and its ends.

14

15 **Figure 3. The unpaired nucleotide frequency (UNF) across the miRNA sequence.** 5' end of  
16 the miRNA starts from position one. Negative positions correspond to the miRNA flank. The  
17 UNF is not shown at the very ends of several long miRNAs. (A) Animal miRNAs. (B) The UNF  
18 dependence on the relative rate of nucleotide substitutions in animal miRNAs [62]. The seed  
19 points concentrate near the very UNF-axis. Spearman's rank correlation test was used to estimate  
20 the significance of the correlation between the UNF and the rate of nucleotide substitutions  
21 ( $\rho=0.81$ ,  $P=2.76 \times 10^{-6}$ ). (C-D) The UNF profile of 5' miRNAs and of 3' miRNAs of human and  
22 mouse (mirtrons and non-mirtrons) and of animal excluding human and mouse.

23

24 Wheeler reported that the substitution rate reflects the importance of the positions 2-8 and 13-  
25 16 [62]. To reveal the relation between secondary structure and mutations within miRNA  
26 sequence we plotted the dependence of the UNF on the nucleotide substitution rate taken from

1 [62]. Although the correlation between the UNF and the substitution rate is significant  
2 ( $P=2.76 \times 10^{-6}$ ), this dependence is stepwise rather than linear (figure 2B). The step is formed by  
3 two groups of positions (grey and white regions) in each of them the dependence does not exist.  
4 The seed region (positions 2-8) contains the most conserved positions in the miRNA [62]. Four  
5 neighboring positions of the seed and of the additional binding site (positions 9 and 17-19) form  
6 the transition from the conserved double-stranded to the more variable single-stranded regions of  
7 the miRNA. This agrees with the fact that a RNA base-pairs near a loop use to be partially  
8 unwinded. Taking together figures 3A and 3B we conclude that the secondary structure is one of  
9 the miRNA evolutionary constraints in the same way as for the majority of structural RNAs  
10 where the loops are more variable.

11 On figures 3A and 3B we observe three variable and, at the same time, single-stranded  
12 miRNA regions: the center and the ends. The 5' miRNA end together with the mismatches at  
13 miRNA center are responsible for Ago-protein sorting, as it was shown for some animal species  
14 [71]. Figures 3A and 3B (as well as figure 2) represent the miRNA ends effort to be bounded by  
15 single-stranded nucleotides. As a part of this tendency, the left (right) peak on figure 3C (3D)  
16 supports again the existence of the loop-counting rule for the Drosha cleavage site. As for  
17 mirtrons they have a similar UNF profile over the whole sequence except the positions by the  
18 hairpin base (figures 3C and 3D), where the mirtrons are much more single-stranded what  
19 characterizes their unique biogenesis. In particular, the first nucleotide of the most 5' mirtrons  
20 (79%) is single stranded and quite often survives after the intron cropping by exonucleases [31].  
21 Its opposite end of the 3' pre-miRNA is also often single-stranded and uses to be immediately  
22 formed by splicing [31].

23

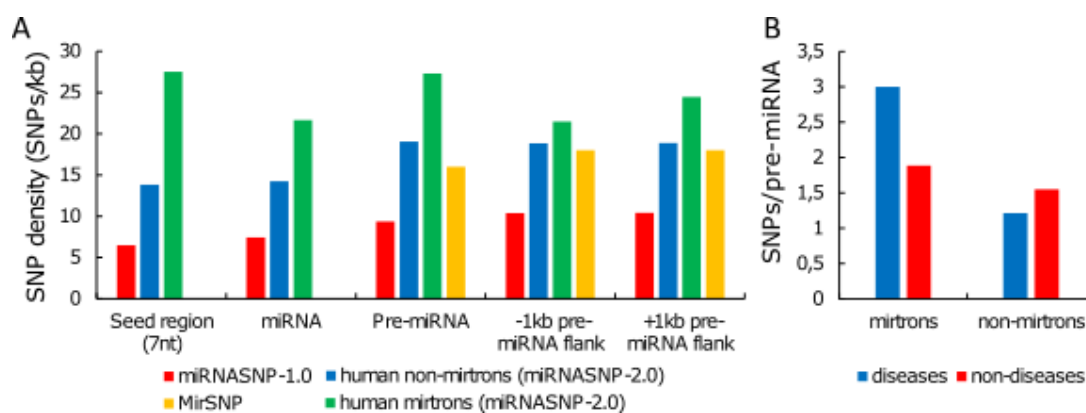
## 24 **SNP density of human miRNA and its neighbourhood**

25 SNPs are the most frequent genetic changes in human genome. The miRNA-related SNPs  
26 may affect the miRNA functions and subsequently result in the phenotype changes and diseases

1 so that some miRNAs could appear as the diseases-prediction biomarkers. These SNPs could  
 2 potentially alter miRNA maturation, silencing machinery, pri-/pre-/miRNA structure, miRNA  
 3 expression and target binding [72].

4 The previous papers have produced the controversial values of the SNP densities (red and  
 5 yellow bars, figure 4A) [73, 74]. Therefore, we recalculated the SNP density using the latest  
 6 human SNPs in the pre-miRNAs and their flanks. Returning to the keynote of our paper we  
 7 consider separately the SNP densities of non-mirtrons and mirtrons (blue and green bars,  
 8 correspondingly, figure 4A).

9



10

11 **Figure 4. SNP in pre-miRNA and its link with diseases.** (A) SNP density in human pre-  
 12 miRNAs and their flanking regions. Blue (green) bars correspond to human non-mirtrons  
 13 (mirtrons). The results are based on miRNASNP-2.0 (miRBase 19.0 and dbSNP137) [58]. Red  
 14 bars were calculated by [73] using miRBase-16.0 and dbSNP132 (miRNASNP-1.0 database).  
 15 Yellow bars were calculated by [74] using miRBase 18.0 and dbSNP135 (MirSNP database).  
 16 SNP density is shown separately for seed region, microRNA (without seed region), pre-miRNA  
 17 (without microRNA) and upstream and downstream pre-miRNA flanks. Note, that the MirSNP  
 18 data (yellow bars) provide the densities of slightly different regions, namely, the entire pre-  
 19 miRNA sequence and both 200bp pre-miRNA flanks. (B) SNP occurrence per pre-miRNA for  
 20 disease or non-disease mirtrons and non-mirtrons. The disease pre-miRNAs are associated with



1 at least one disease and were taken from [59, 60], SNPs were extracted from miRNASNP-2.0  
2 [58].

3

4 The region regarding the canonical miRNAs can be divided into two parts of equal  
5 conservation: miRNA sequence and pre-miRNA (excluding the miRNA) with its flanks (figure  
6 4A). This confirms only one relation in the previously found hierarchy of the SNP densities  
7 (seed < miRNA < pre-miRNA < flanks) [60, 72, 73]. In particular, the small difference between  
8 the seed and the rest of the miRNA (figure 4A, red bars) became even weaker (blue bars).

9 The miRNA sequence is saturated with functional sites (seed, additional binding site, etc.).  
10 Besides, miRNA duplex holds a most part of its pre-miRNA and thus carries the main burden of  
11 responsibility for forming the pre-miRNA hairpin. The rest of the pre-miRNA and its flanks also  
12 include a number of important regulatory elements (e.g. miRNA binding sites, UG, CNNC and  
13 other motifs [68, 75]), which are sparsely distributed over RNA sequence. This agrees with our  
14 observation that the pre-miRNAs and their flanks have greater SNP density than the miRNA  
15 sequences (figure 4A).

16 The other reason of hierarchy disappearing could be that the updates of miRBase and dbSNP  
17 cause the SNP density to increase throughout the regions and the difference between them to  
18 smooth off. For example, the dbSNP can rapidly accumulate rare SNPs and the miRBase – less  
19 conserved sequences. We consider this in the last section where we test the robustness of our  
20 results. As it turned out, the hierarchy of the SNP densities is being restored for common SNPs  
21 in the robust miRNAs.

22 In the mirtrons the SNPs take place more frequently than in other miRNAs and pre-miRNAs  
23 (figure 4A) in accordance with the fact that SNPs most often occur in introns [76]. Intronic SNPs  
24 can affect mRNA expression and splicing [77] and thus may influence the mirtron processing. In  
25 contrast with the non-mirtrons, SNPs proceed in the mirtron pre-miRNAs more frequently than

1 in their flanks, in line with that the mirtron pre-miRNA flanks can often contain exons where  
2 SNPs rarely occur [76].

3 Lu et al. and later Han et al. explored the relation between miRNA conservation and diseases,  
4 and found out that SNPs occur less frequently in miRNAs associated with diseases than in non-  
5 associated ones [59, 60]. Using recent SNP data we re-calculated the SNP occurrence of non-  
6 mirtrons (figure 4B) and confirmed their observation that SNP-rare pre-miRNAs are often  
7 associated with a number of diseases while SNP-frequent ones are not [59, 60]. In contrast, the  
8 non-disease mirtron pre-miRNAs have lower number of the SNPs per sequence than the disease  
9 ones. Taking together the species specificity of mirtrons [45, 46], their increased SNP density  
10 (figure 4A) and their disease association, this suggests that mirtrons undergo positive selection  
11 while the most canonical miRNAs are under the negative one [78]. These also agrees with our  
12 observation that the mirtrons have wider distributions than the non-mirtrons (figures 2 and 4).

13 However, the results of such a simple SNP analysis as above should be taken with great  
14 caution: the carefully prepared samples, allele frequencies and other population genetic  
15 parameters are needed for a more profound analysis.

16

## 17 **Branchpoints in mirtrons and introns**

18 Mirtrons are confirmed by splicing dependence of miRNA expression. And vice versa, the  
19 hairpin of future pre-miRNA may affect splicing, in particular, the branchpoint site recognition.

20 In a recent paper [61] authors extensively studied the human, mouse and yeast branchpoints.  
21 Comparing their U2 basepairing modes we observe that the U2:mirtron model has been  
22 identified less frequently than the U2:intron model (table 1, Additional file 6): thus, more close  
23 inspection of mirtron splicing, in particular, the role of intronic hairpins immediately upstream  
24 the 3' splice site, can shed a new light on splicing in general.

25 Using human and mouse data from Taggart et al. [61] we compare the branchpoint  
26 distribution of mirtrons and introns. Despite the distributions similarity, mirtron branchpoints are

1 more often located in the expected region (10-40 nucleotides upstream from the 3' splice site)  
2 than the introns what suggests the more frequent constitutive splicing of the mirtron precursors  
3 (figure 1, Additional file 6). Remarkably, the mirtron branchpoints are most probably located at  
4 18-24nt away from 3' splice site, i.e. near the miRNA end (figure 1, Additional file 6).

5

6 **Figure 5. Distance from terminal loop to branchpoint (A) and to 3' pre-miRNA end (B).**  
7 **Distance between branchpoint and Dicer cleavage site of 3' miRNA (C).** Note that the 3'  
8 miRNAs are mainly located near the 3' splice site. The two separate points on the left figure  
9 show the branchpoint frequencies of the 5' pre-miRNA strand (green) and of the terminal loop  
10 (red) as a whole. The blue curve displays the branchpoint distribution along the 3' pre-miRNA  
11 strand. On the center figure shown are the distances of two mirtron groups: with branchpoint  
12 within the 3' pre-miRNA strand (blue) and with branchpoint within the terminal loop (red). On  
13 the right figure considered are the distances between the 5' end of the 3' miRNA and the  
14 branchpoint within the 3' pre-miRNA strand. Negative values correspond to the branchpoints  
15 into miRNA sequence.

16

17 To characterize the branchpoint distribution more closely we analyzed its locations along pre-  
18 miRNA sequences and found that the branchpoint often appears in terminal loop or in 3' pre-  
19 miRNA strand (Figure 5A). In those cases when it was found in terminal loop, the pre-miRNA  
20 hairpin is most frequently as short as possible (figure 5B). When branchpoint is located in 3'  
21 strand, its site is more likely to be 6-8 nt away from the terminal loop overlapping with the Dicer  
22 cleavage site as it follows from the mirtron length distributions (figure 5B) and from direct  
23 calculation of distances between branchpoint and Dicer cleavage site (figure 5C). The latter site  
24 attracts a nearest loop (figure 2C) so we conclude that the mirtron pre-miRNA secondary  
25 structure tends not to shield the branchpoint (therefore not to block the U2:mirtron basepairing)  
26 but rather to fix branchpoint and 3' splice site mutual disposition.

1

## 2 **Checking the stability of the results**

3 MiRBase is often criticized for including the transcriptional noise [57, 79-81] which  
4 sometimes is estimated to be as high as 2/3 of the annotated human miRNAs [57]. This criticism  
5 has motivated to filter the miRBase entries and to establish a new miRNA catalog (MirGeneDB  
6 database [57]).

7 To verify that our results are robust to the miRBase false positives, we repeated our  
8 computations for miRNAs whose identifiers are presented in MirGeneDB. The MirGeneDB  
9 authors rejected the most of the mirtrons, mainly by the improper mature/star offset and the  
10 undesirable heterogeneous processing. As a result, the MirGeneDB contains only 7 human and  
11 mouse mirtron entries (mmu-mir-1981, mmu-mir-3097, hsa-mir-3605, hsa-mir-3940, hsa-mir-  
12 4640, hsa-mir-5010, hsa-mir-6746). Therefore, we considered only the non-mirtrons.

13 We found that most of our results for non-mirtrons are stable against the miRBase reducing to  
14 robust miRNAs (Additional file 7). The main differences are the following. The miRNA regions  
15 became more contrasting, especially for the SNP densities. Generally, the SNP densities  
16 decreased up to 1.5-times and the difference in the densities of the pre-miRNA and its flanks re-  
17 appeared. After additionally removing the rare SNPs [82] the densities decreased much stronger  
18 and their hierarchy was completely restored (seed < miRNA < pre-miRNA < flanks, Additional  
19 file 7). The common SNPs are likely older, they have been subjected to selective forces over  
20 time [83] and produces the difference between seed and the rest of miRNA.

21 Most of variance, in particular within miRNAs, arises from rare variants most of which are  
22 either recently derived alleles or being selected against due to their deleterious nature [84]. This  
23 effect may be most pronounced in modern humans who live under relaxed selection and readily  
24 accumulate the deleterious rare alleles [85]. Mirtrons are more variable and likely carry more  
25 rare SNPs than the canonical miRNAs. Multiple rare SNPs often associate with complex

1 diseases. Therefore we speculate that mirtrons are a rich source of disease-promoting variants  
2 (figure 4B).

3

## 4 **Conclusions**

5 The pri-/pre-miRNA secondary structure plays an important role on each stage of biogenesis,  
6 in particular by positioning the miRNA excision sites. Drosha as well as Dicer can cleave  
7 imprecisely around the expected sites. For mirtrons splicing replaces the Drosha cleavage and  
8 manifests itself in the footprints of the secondary structure near the pre-miRNA base (dangling  
9 ends and less precise cleavage), while for the Dicer cleavage sites the characteristics of the  
10 canonical biogenesis matches the non-canonical ones. Both complexes recognize the  
11 inner/bulge-loop structure; therefore the loop-counting rule can help to predict not only Dicer  
12 cleavage site but also Drosha one. Dicer binds the pre-miRNA ends: as the result, the imprecise  
13 Drosha cleavage can induce Dicer error what appears in the dependence of the overhang lengths.  
14 To explain this interrelation of Dicer and Drosha precision we suggest the two-lever model of  
15 Dicer movements where the distance between RNase IIIA/IIIB cleavage sites fits the distance  
16 between pre-miRNA ends. In mirtron pre-miRNAs both ends are typically hanging and their  
17 distance varies widely thus increasing the Dicer cleavage imprecision. Also the mirtron hairpin  
18 brings together the splice sites for 60-80nt closer, exposes branchpoint site and adjusts it on the  
19 3' splice site. The mirtron structure appears to be well suited to the splicing and thus the mirtrons  
20 can evolve from the occasional hairpins (as readily as other miRNAs) in the immediate  
21 neighbourhood of the 3' splice site. Also through the splicing mirtrons can acquire guanine at  
22 their ends what induces the Dicer imprecision.

23 The secondary structure of pri-/pre-miRNAs appears also in their evolution, in particular, in  
24 clear difference of mutation rates between single- and double-stranded positions: thus the  
25 secondary structure shapes the functional subdivision of the precursor (seed, additional binding  
26 site, inner and terminal loops, etc.). For the SNP density extracted from the last miRNA SNP

1 database this division is reduced because of rare SNPs and non-robust miRNAs. In contrast to  
2 the canonical miRNAs the mirtrons exhibit higher SNP density and more SNPs per pre-miRNAs  
3 that are associated with diseases. This suggests that mirtrons unlike old canonical miRNAs are  
4 under positive selection and serve as an inherent source of silencing variability.

5

## 6 **List of abbreviations**

7 UNF – unpaired nucleotide frequency.

8

## 9 **Declarations**

### 10 **Ethics approval and consent to participate**

11 Not applicable.

### 12 **Consent for publication**

13 Not applicable.

### 14 **Availability of data and material**

15 All data analysed during this study are freely available from the sources in this published  
16 article.

### 17 **Competing interests**

18 The authors declare that they have no competing interests.

### 19 **Funding**

20 This work was supported by budget funding of governmental task (project № 0324-2016-  
21 0008).

### 22 **Authors' contributions**

1 PSV performed bioinformatics calculations. IIT and PSV contributed to the interpretation of  
2 results, and were involved in manuscript editing. IIT designed and coordinated the study. All  
3 authors have read and approved the final manuscript.

#### 4 **Acknowledgements**

5 Not applicable.

#### 6 **Additional files**

7 Additional file 1. The particular examples of the 3' end shifting inside miRNA. docx.

8 Additional file 2. Change of the pre-miRNA secondary structure reduces a pair of long  
9 overhangs to nearly canonical one. docx.

10 Additional file 3. The model of independent overhangs and its results. docx.

11 Additional file 4. Overhang length dependence are robust against the overhang definition.  
12 docx.

13 Additional file 5. The scheme of the proposed 2-lever model for Dicer. png.

14 Additional file 6. The branchpoint locations across mirtrons and introns. docx.

15 Additional file 7. The results of the stability test using MiRGeneDB database. docx.

16

#### 17 **References**

18 1. Lee Y, Kim M, Han J, Yeom KH, Lee S, Baek SH, Kim VN. MicroRNA genes are transcribed  
19 by RNA polymerase II. *The EMBO journal*. 2004;23:4051-4060.

20 2. Borchert GM, Lanier W, Davidson BL. RNA polymerase III transcribes human microRNAs.  
21 *Nature structural & molecular biology*. 2006;13:1097-1101.

22 3. Lee Y, Jeon K, Lee JT, Kim S, Kim VN. MicroRNA maturation: stepwise processing and  
23 subcellular localization. *The EMBO journal*. 2002;21:4663-4670.

24 4. Ozsolak F, Poling LL, Wang Z, Liu H, Liu XS, Roeder RG, Fisher DE. Chromatin structure  
25 analyses identify miRNA promoters. *Genes & development*. 2008;22:3172-3183.

- 1 5. Monteys AM, Spengler RM, Wan J, Tecedor L, Lennox KA, Xing Y, Davidson BL. Structure  
2 and activity of putative intronic miRNA promoters. *RNA*. 2010;16:495-505.
- 3 6. Altuvia Y, Landgraf P, Lithwick G, Elefant N, Pfeffer S, Aravin A, Margalit H. Clustering  
4 and conservation patterns of human microRNAs. *Nucleic Acids Research*. 2005;33:2697-2706.
- 5 7. Titov II, Vorozheykin PS. Analysis of miRNA duplication in the human genome and the role  
6 of transposon evolution in this process. *Russian Journal of Genetics: Applied Research*.  
7 2011;1:308-314.
- 8 8. Marco A, Ninova M, Griffiths-Jones S. Multiple products from microRNA transcripts.  
9 *Biochemical Society transactions*. 2013;41:850-854.
- 10 9. Roush S, Slack FJ. The let-7 family of microRNAs. *Trends in cell biology*. 2008;18:505-516.
- 11 10. Tyler DM, Okamura K, Chung WJ, Hagen JW, Berezikov E, Hannon GJ, Lai EC.  
12 Functionally distinct regulatory RNAs generated by bidirectional transcription and processing of  
13 microRNA loci. *Genes & development*. 2009;22:26-36.
- 14 11. Han J, Lee Y, Yeom KH, Nam JW, Heo I, Rhee JK, Kim VN. Molecular basis for the  
15 recognition of primary microRNAs by the Drosha-DGCR8 complex. *Cell*. 2006;125:887-901.
- 16 12. Nguyen TA, Jo MH, Choi YG, Park J, Kwon SC, Hohng S, Woo JS. Functional anatomy of  
17 the human microprocessor. *Cell*. 2005;161:1374-1387.
- 18 13. Kwon SC, Nguyen TA, Choi YG, Jo MH, Hohng S, Kim VN, Woo JS. Structure of human  
19 DROSHA. *Cell*. 2016;164:81-90.
- 20 14. Yi R, Qin Y, Macara IG, Cullen BR. Exportin-5 mediates the nuclear export of pre-  
21 microRNAs and short hairpin RNAs. *Genes & development*. 2003;17:3011-3016.
- 22 15. Xie M, Li M, Vilborg A, Lee N, Shu MD, Yartseva V, Steitz JA. Mammalian 5'-capped  
23 microRNA precursors that generate a single microRNA. *Cell*. 2013;155:1568-1580.
- 24 16. Zeng Y, Cullen BR. Structural requirements for pre-microRNA binding and nuclear export  
25 by Exportin 5. *Nucleic Acids Research*. 2004;32:4776-4785.



- 1 17. Bernstein E, Caudy AA, Hammond SM, Hannon GJ. Role for a bidentate ribonuclease in the  
2 initiation step of RNA interference. *Nature*. 2001;409:363-366.
- 3 18. Svobodova E, Kubikova J, Svoboda P. Production of small RNAs by mammalian Dicer.  
4 *Pflugers Archiv-European Journal of Physiology*. 2001;468:1089-1102.
- 5 19. Lau PW, Guiley KZ, De N, Potter CS, Carragher B, MacRae IJ. The molecular architecture  
6 of human Dicer. *Nature structural & molecular biology*. 2012;19:436-440.
- 7 20. MacRae IJ, Zhou K, Li F, Repic A, Brooks AN, Cande WZ, Doudna JA. Structural basis for  
8 double-stranded RNA processing by Dicer. *Science*. 2006;311:195-198.
- 9 21. MacRae IJ, Zhou K, Doudna JA. Structural determinants of RNA recognition and cleavage  
10 by Dicer. *Nature structural & molecular biology*. 2007;14:934-940.
- 11 22. Bartel DP. MicroRNAs: genomics, biogenesis, mechanism, and function. *Cell*.  
12 2004;116:281-297.
- 13 23. Shin C. Cleavage of the star strand facilitates assembly of some microRNAs into Ago2-  
14 containing silencing complexes in mammals. *Molecules & Cells*. 2008;26:3.
- 15 24. Okamura K, Phillips MD, Tyler DM, Duan H, Chou YT, Lai EC. The regulatory activity of  
16 microRNA\* species has substantial influence on microRNA and 3' UTR evolution. *Nature*  
17 *structural & molecular biology*. 2008;15:354-363.
- 18 25. Li JH, Liu S, Zhou H, Qu LH, Yang JH. StarBase v2. 0: decoding miRNA-ceRNA, miRNA-  
19 ncRNA and protein-RNA interaction networks from large-scale CLIP-Seq data. *Nucleic Acids*  
20 *Research*. 2014;42:D92-D97.
- 21 26. Hirata H, Hinoda Y, Shahryari V, Deng G, Nakajima K, Tabatabai ZL, Dahiya R. Long  
22 noncoding RNA MALAT1 promotes aggressive renal cell carcinoma through Ezh2 and interacts  
23 with miR-205. *Cancer research*. 2015;75:1322-1331.
- 24 27. Yao Y, Ma J, Xue Y, Wang P, Li Z, Liu J, Li Z. Knockdown of long non-coding RNA XIST  
25 exerts tumor-suppressive functions in human glioblastoma stem cells by up-regulating miR-152.  
26 *Cancer letters*. 2015;359:75-86.

- 1 28. Ruby JG, Jan CH, Bartel DP. Intronic microRNA precursors that bypass Drosha processing.
- 2 Nature. 2007;448:83-86.
- 3 29. Yang JS, Lai EC. Alternative miRNA biogenesis pathways and the interpretation of core
- 4 miRNA pathway mutants. Molecular cell. 2011;43:892-903.
- 5 30. Okamura K, Hagen JW, Duan H, Tyler DM, Lai EC. The mirtron pathway generates
- 6 microRNA-class regulatory RNAs in Drosophila. Cell. 2007;130:89-100.
- 7 31. Westholm JO, Lai EC. Mirtrons: microRNA biogenesis via splicing. Biochimie.
- 8 2011;93:1897-1904.
- 9 32. Ladewig E, Okamura K, Flynt AS, Westholm JO, Lai EC. Discovery of hundreds of mirtrons
- 10 in mouse and human small RNA data. Genome research. 2012;22:1634-1645.
- 11 33. Sibley CR, Seow Y, Saayman S, Dijkstra KK, El Andaloussi S, Weinberg MS, Wood MJ.
- 12 The biogenesis and characterization of mammalian microRNAs of mirtron origin. Nucleic acids
- 13 research. 2011;40:438-448.
- 14 34. Schamberger A, Sarkadi B, Orbán TI. Human mirtrons can express functional microRNAs
- 15 simultaneously from both arms in a flanking exon-independent manner. RNA biology.
- 16 2012;9:1177-1185.
- 17 35. Havens MA, Reich AA, Duelli DM, Hastings ML. Biogenesis of mammalian microRNAs by
- 18 a non-canonical processing pathway. Nucleic Acids Research. 2012;40:4626-4640.
- 19 36. Curtis HJ, Sibley CR, Wood MJ. Mirtrons, an emerging class of atypical miRNA. Wiley
- 20 Interdisciplinary Reviews: RNA. 2012;3:617-632.
- 21 37. Abdelfattah AM, Park C, Choi MY. Update on non-canonical microRNAs. Biomolecular
- 22 concepts. 2014;5:275-287.
- 23 38. Cheloufi S, Dos Santos CO, Chong MM, Hannon GJ. A dicer-independent miRNA
- 24 biogenesis pathway that requires Ago catalysis. Nature. 2012;465:584-589.

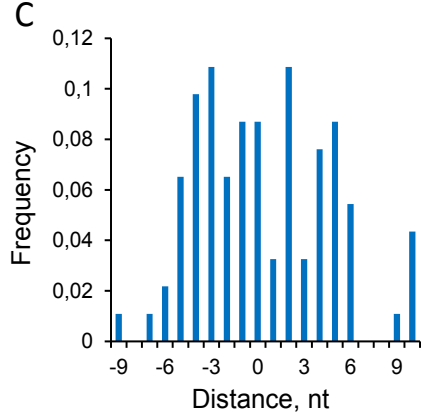
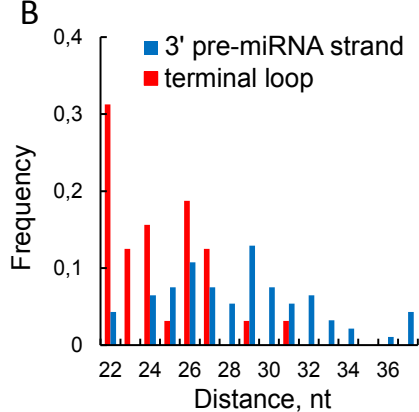
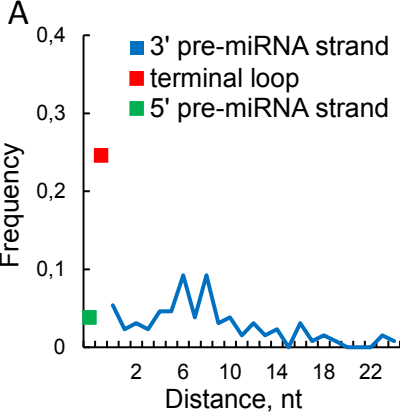
- 1 39. Cifuentes D, Xue H, Taylor DW, Patnode H, Mishima Y, Cheloufi S, Wolfe SA. A novel  
2 miRNA processing pathway independent of Dicer requires Argonaute2 catalytic activity.  
3 *Science*. 2012;328:1694-1698.
- 4 40. Liu YP, Schopman NC, Berkhout B. Dicer-independent processing of short hairpin RNAs.  
5 *Nucleic Acids Research*. 2013;41:3723-3733.
- 6 41. Yoda M, Cifuentes D, Izumi N, Sakaguchi Y, Suzuki T, Giraldez AJ, Tomari Y. Poly (A)-  
7 specific ribonuclease mediates 3'-end trimming of Argonaute2-cleaved precursor microRNAs.  
8 *Cell reports*. 2013;5:715-726.
- 9 42. Nozawa M, Miura S, Nei M. Origins and evolution of microRNA genes in *Drosophila*  
10 species. *Genome biology and evolution*. 2010;2:180-189.
- 11 43. Титов ИИ, Ворожейкин ПС. мРНК-содержащие транспозоны человека. Вавиловский  
12 журнал генетики и селекции. 2011;15: 323-326.
- 13 44. Lyu Y, Shen Y, Li H, Chen Y, Guo L, Zhao Y, Tang T. New microRNAs in *Drosophila* -  
14 birth, death and cycles of adaptive evolution. *PLoS genetics*. 2014;10:e1004096.
- 15 45. Berezikov E, Chung WJ, Willis J, Cuppen E, Lai EC. Mammalian mirtron genes. *Molecular*  
16 *cell*. 2007;28:328-336.
- 17
- 18 46. Berezikov E, Liu N, Flynt AS, Hodges E, Rooks M, Hannon GJ, Lai EC. Evolutionary flux  
19 of canonical microRNAs and mirtrons in *Drosophila*. *Nature genetics*. 2010;42:6-9.
- 20 47. Gkirtzou K, Tsamardinos I, Tsakalides P, Poirazi P. MatureBayes: a probabilistic algorithm  
21 for identifying the mature miRNA within novel precursors. *PloS one*. 2010;5:e11843.
- 22 48. Geis M, Middendorf M. Particle swarm optimization for finding RNA secondary structures.  
23 *International Journal of Intelligent Computing and Cybernetics*. 2011;4:160-186.
- 24 49. Xuan P, Guo M, Huang Y, Li W, Huang Y. MaturePred: efficient identification of  
25 microRNAs within novel plant pre-miRNAs. *PloS one*. 2011;6:e27422.

- 1 50. Leclercq M, Diallo AB, Blanchette M. Computational prediction of the localization of  
2 microRNAs within their pre-miRNA. *Nucleic Acids Research*. 2013;41:7200-7211.
- 3 51. Karathanasis N, Tsamardinos I, Poirazi P. MiRduplexSVM: A high-Performing miRNA-  
4 duplex prediction and evaluation methodology. *PloS one*. 2015;10:e0126151.
- 5 52. Li J, Xu C, Wang L, Liang H, Feng W, Cai Z, Liu Y. Prediction of pre-microRNA secondary  
6 structure based on reverse complementary folding. *IFAC-PapersOnLine*. 2015;48:239-244.
- 7 53. Peace RJ, Biggar KK, Storey KB, Green JR. A framework for improving microRNA  
8 prediction in non-human genomes. *Nucleic Acids Research*. 2015;43:e138.
- 9 54. Marques YB, de Paiva Oliveira A, Vasconcelos ATR, Cerqueira FR. Mirnacle: machine  
10 learning with SMOTE and random forest for improving selectivity in pre-miRNA ab initio  
11 prediction. *BMC Bioinformatics*. 2016;17:53.
- 12 55. Griffiths-Jones S. The microRNA registry. *Nucleic Acids Research*. 2004;32:D109–D111.
- 13 56. Wen J, Ladewig E, Shenker S, Mohammed J, Lai EC. Analysis of nearly one thousand  
14 mammalian mirtrons reveals novel features of dicer substrates. *PLoS Comput Biol*.  
15 2015;11:e1004441.
- 16 57. Fromm B, Billipp T, Peck L, Johansen M, Tarver JE, King BL, Peterson KJ. A uniform  
17 system for the annotation of vertebrate microRNA genes and the evolution of the human  
18 microRNAome. *Annual review of genetics*. 2015;49:213-242.
- 19 58. Gong J, Liu C, Liu W, Wu Y, Ma Z, Chen H, Guo Y. An update of miRNA SNP database for  
20 better SNP selection by GWAS data, miRNA expression and online tools. *Database (Oxford)*.  
21 2015; doi: 10.1093/database/bav029.
- 22 59. Lu M, Zhang Q, Deng M, Miao J, Guo Y, Gao W, Cui Q. An analysis of human microRNA  
23 and disease associations. *PloS one*. 2008;3:e3420.
- 24 60. Han M, Zheng Y. Comprehensive analysis of single nucleotide polymorphisms in human  
25 microRNAs. *PloS one*. 2013;8:e78028.

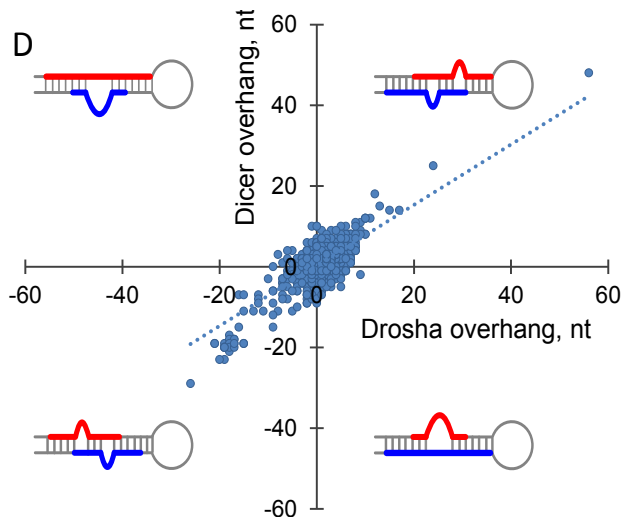
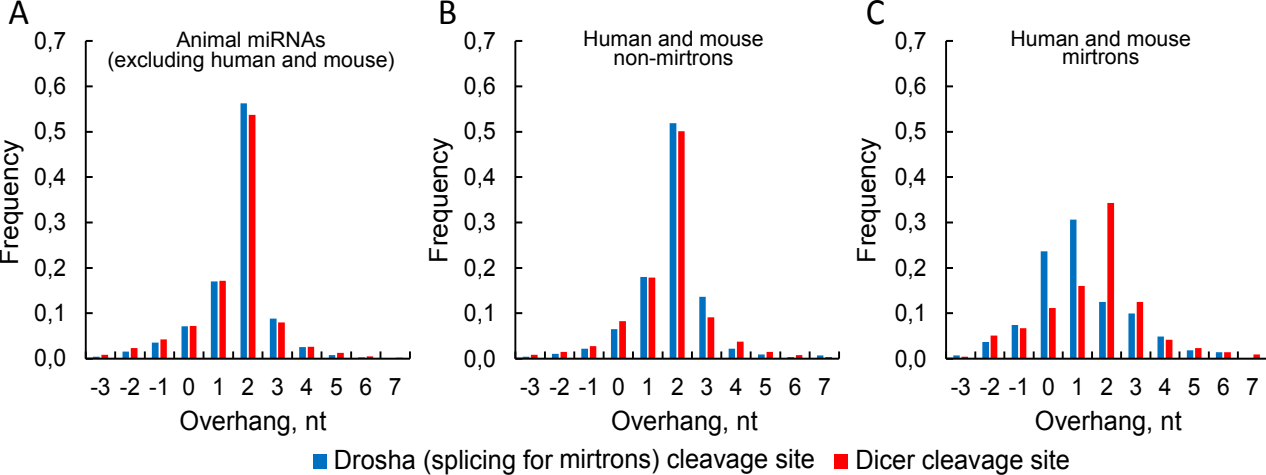
- 1 61. Taggart AJ, Lin CL, Shrestha B, Heintzelman C, Kim S, Fairbrother WG. Large-scale  
2 analysis of branchpoint usage across species and cell lines. *Genome research*. 2017;27:639-649.
- 3 62. Wheeler BM, Heimberg AM, Moy VN, Sperling EA, Holstein TW, Heber S, Peterson KJ.  
4 The deep evolution of metazoan microRNAs. *Evolution & development*. 2009;11:50-68.
- 5 63. Starega-Roslan J, Witkos TM, Galka-Marciniak P, Krzyzosiak WJ. Sequence features of  
6 Drosha and Dicer cleavage sites affect the complexity of isomiRs. *International journal of*  
7 *molecular sciences*. 2015;16:8110-8127.
- 8 64. Starega-Roslan J, Galka-Marciniak P, Krzyzosiak WJ. Nucleotide sequence of miRNA  
9 precursor contributes to cleavage site selection by Dicer. *Nucleic Acids Research*.  
10 2015;43:10939-10951.
- 11 65. Starega-Roslan, 2011. Starega-Roslan J, Krol J, Koscianska E, Kozlowski P, Szlachcic WJ,  
12 Sobczak K, Krzyzosiak WJ. Structural basis of microRNA length variety. *Nucleic Acids*  
13 *Research*. 2011;39:257-268.
- 14 66. Starega-Roslan J, Koscianska E, Kozlowski P, Krzyzosiak WJ. The role of the precursor  
15 structure in the biogenesis of microRNA. *Cellular and molecular life sciences*. 2011;68:2859.
- 16 67. Ma H, Wu Y, Niu Q, Zhang J, Jia G, Manjunath N, Wu H. A sliding-bulge structure at the  
17 Dicer processing site of pre-miRNAs regulates alternative Dicer processing to generate 5'-  
18 isomiRs. *Heliyon*. 2016;2:e00148.
- 19 68. Libri V, Miesen P, van Rij RP, Buck AH. Regulation of microRNA biogenesis and turnover  
20 by animals and their viruses. *Cellular and Molecular Life Sciences*. 2013;70:3525-3544.
- 21 69. Gu S, Jin L, Zhang Y, Huang Y, Zhang F, Valdmanis PN, Kay MA. The loop position of  
22 shRNAs and pre-miRNAs is critical for the accuracy of Dicer processing in vivo. *Cell*.  
23 2012;151:900–911.
- 24 70. Sarzyńska J, Mickiewicz A, Miłostan M, Łukasiak P, Błażewicz J, Figlerowicz M, Kuliński  
25 T. Flexibility of dicer studied by implicit solvent molecular dynamics simulations.  
26 *Computational Methods in Science and Technology*. 2010;16:97-104.

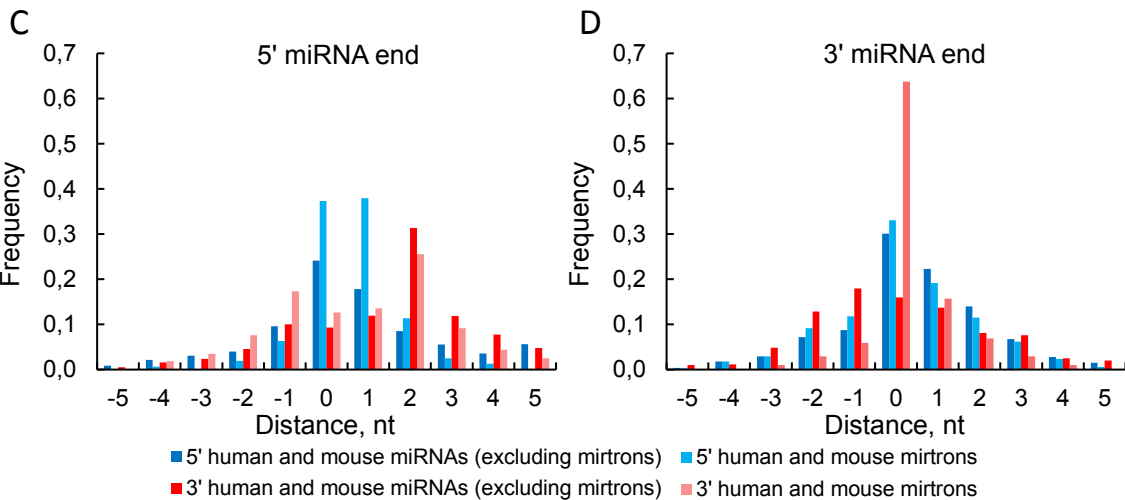
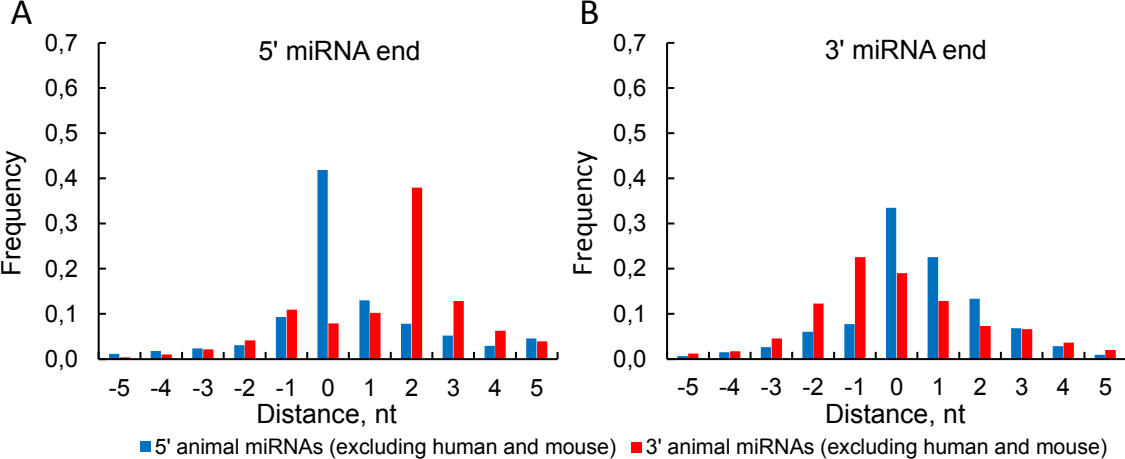
- 1 71. Ha M, Kim VN. Regulation of microRNA biogenesis. *Nature reviews Molecular cell*  
2 *biology*. 2014;15:509-524.
- 3 72. Jin Y, Lee CG. Single nucleotide polymorphisms associated with microRNA regulation.  
4 *Biomolecules*. 2013;3:287-302.
- 5 73. Gong J, Tong Y, Zhang HM, Wang K, Hu T, Shan G, Guo AY. Genome-wide identification  
6 of SNPs in microRNA genes and the SNP effects on microRNA target binding and biogenesis.  
7 *Human mutation*. 2012;33:254-263.
- 8 74. Liu C, Zhang F, Li T, Lu M, Wang L, Yue W, Zhang D. MirSNP, a database of  
9 polymorphisms altering miRNA target sites, identifies miRNA-related SNPs in GWAS SNPs  
10 and eQTLs. *BMC genomics*. 2012;13:1.
- 11 75. Krol J, Loedige I, Filipowicz W. The widespread regulation of microRNA biogenesis,  
12 function and decay. *Nature Reviews Genetics*. 2010;11:597-610.
- 13 76. Freedman ML, Monteiro AN, Gayther SA, Coetzee GA, Risch A, Plass C, James M.  
14 Principles for the post-GWAS functional characterisation of risk loci. *Nature Genetics*. 2011;  
15 doi:10.1038/ng.840.
- 16 77. Monlong J, Calvo M, Ferreira PG, Guigó R. Identification of genetic variants associated with  
17 alternative splicing using sQTLseeker. *Nature communications*. 2014;5:1.
- 18 78. Quach H, Barreiro LB, Laval G, Zidane N, Patin E, Kidd KK, Quintana-Murci L. Signatures  
19 of purifying and local positive selection in human miRNAs. *The American Journal of Human*  
20 *Genetics*. 2009;84:316-327.
- 21 79. Chiang HR, Schoenfeld LW, Ruby JG, Auyeung VC, Spies N, Baek D, Blelloch R.  
22 Mammalian microRNAs: experimental evaluation of novel and previously annotated  
23 genes. *Genes & development*. 2010;24:992-1009.
- 24 80. Wang X, Liu XS. Systematic curation of miRBase annotation using integrated small RNA  
25 high-throughput sequencing data for *C. elegans* and *Drosophila*. *Frontiers in genetics*. 2011;2:25.

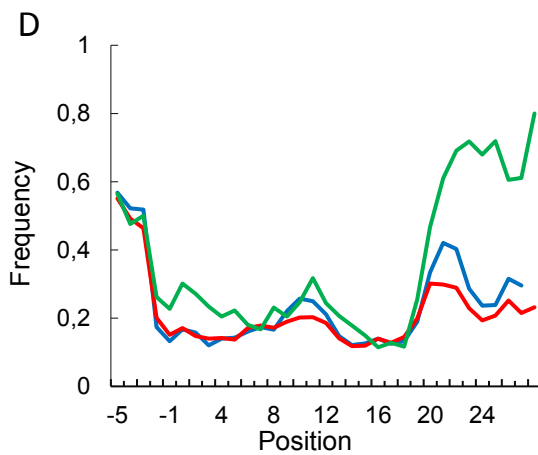
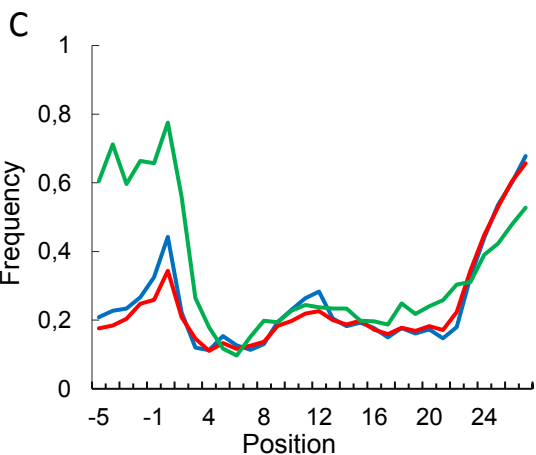
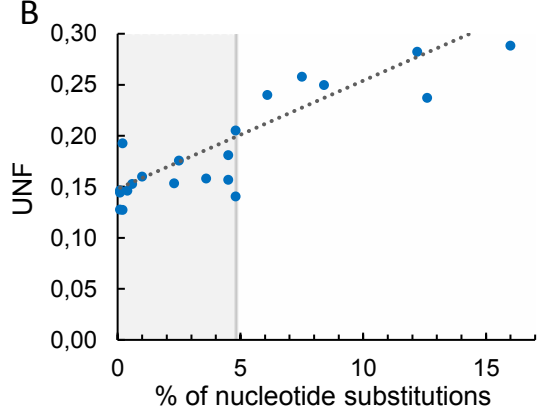
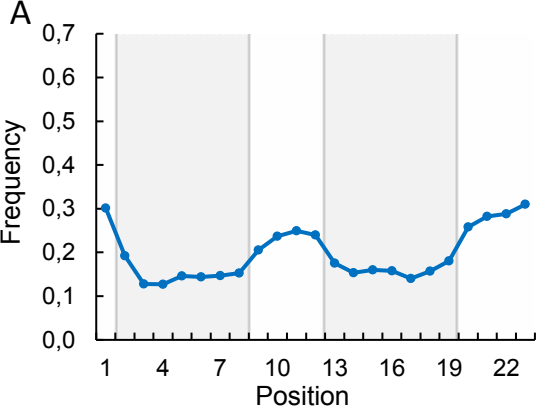
- 1 81. Meng Y, Shao C, Wang H, Chen M. Are all the miRBase-registered microRNAs true? A  
2 structure-and expression-based re-examination in plants. *RNA biology*. 2012;9:249-253.
- 3 82. Simovski B, Vodák D, Gundersen S, Domanska D, Azab A, Holden L, Johansen M. GSuite  
4 HyperBrowser: integrative analysis of dataset collections across the genome and  
5 epigenome. *GigaScience*, 2017;gix032.
- 6 83. Pritchard JK. Are rare variants responsible for susceptibility to complex diseases? *The*  
7 *American Journal of Human Genetics*. 2001;69:124-137.
- 8 84. Gibson G. Rare and common variants: twenty arguments. *Nature reviews. Genetics*.  
9 2011;13:135.
- 10 85. Lynch M. Rate, molecular spectrum, and consequences of human mutation. *Proceedings of*  
11 *the National Academy of Sciences*. 2010;107: 961-968.







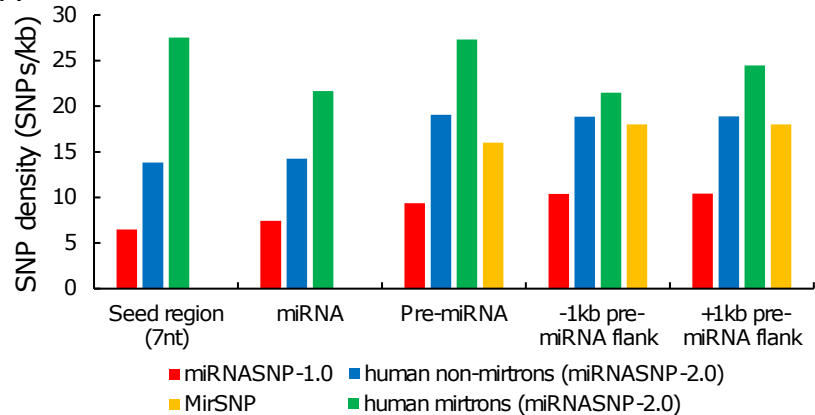




— 5' animal miRNAs (excluding human and mouse)  
 — 5' human and mouse miRNAs (excluding mirtrons)  
 — 5' human and mouse mirtrons

— 3' animal miRNAs (excluding human and mouse)  
 — 3' human and mouse miRNAs (excluding mirtrons)  
 — 3' human and mouse mirtrons

A



B

