

SOFTWARE

Open Access



# CircMarker: a fast and accurate algorithm for circular RNA detection

Xin Li<sup>1</sup>, Chong Chu<sup>2</sup>, Jingwen Pei<sup>1</sup>, Ion Măndoiu<sup>1</sup> and Yufeng Wu<sup>1\*</sup>

From 13th International Symposium on Bioinformatics Research and Applications (ISBRA 2017) Honolulu, Hawaii, USA. 30 May - 2 June 2017

## Abstract

**Background:** While RNA is often created from linear splicing during transcription, recent studies have found that non-canonical splicing sometimes occurs. Non-canonical splicing joins 3' and 5' and forms the so-called circular RNA. It is now believed that circular RNA plays important biological roles such as affecting susceptibility of some diseases. During the past several years, multiple experimental methods have been developed to enrich circular RNA while degrade linear RNA. Although several useful software tools for circular RNA detection have been developed as well, these tools are based on reads mapping may miss many circular RNA. Also, existing tools are slow for large data due to their dependence on reads mapping.

**Method:** In this paper, we present a new computational approach, named CircMarker, based on k-mers rather than reads mapping for circular RNA detection. CircMarker takes advantage of transcriptome annotation files to create the k-mer table for circular RNA detection.

**Results:** Empirical results show that CircMarker outperforms existing tools in circular RNA detection on accuracy and efficiency in many simulated and real datasets.

**Conclusions:** We develop a new circular RNA detection method called CircMarker based on k-mer analysis. Our results on both simulation data and real data demonstrate that CircMarker runs much faster and can find more circular RNA with higher consensus-based sensitivity and high accuracy ratio compared with existing tools.

**Keywords:** Circular RNA, High-throughput sequencing, Genomics, RNA-Seq

## Background

In most eukaryotic genes, coding regions (exons) are separated from noncoding regions (introns) [1]. During RNA splicing, introns are removed and exons are joined to form a contiguous coding sequence called messenger ribonucleic acid (mRNA). This “mature” mRNA is ready for translation, and those contiguous coding sequences are called transcripts [2]. Splicing often occurs in a linear way, which generates the so-called linear RNA. Recent studies show that sometimes circular RNA may be generated during transcription [3]. Circular RNA (or circRNA) is a type of RNA which forms a covalently closed continuous loop.

That is, the 3' and 5' ends normally present in an RNA molecule are joined together [4, 5]. This feature leads to numerous properties of circular RNAs [6]. However, since the amount of circular RNA is often much lower than linear RNA, circular RNA has not been thoroughly studied until recently. During the past several years, several papers report that circular RNA may be associated with diseases and traits [7]. More and more circular RNAs have been identified recently [8, 9].

Since circular RNAs do not have 5' or 3' ends, they are resistant to exonuclease-mediated degradation and are presumably more stable than most linear RNAs in cells. Based on this feature, some benchmark experimental methods have been developed to degrade the linear RNA while enriching the circular RNA. For example, one method is treating samples with RNase R, an enzyme

\*Correspondence: [yufeng.wu@uconn.edu](mailto:yufeng.wu@uconn.edu)

<sup>1</sup>Department of Computer Science and Engineering, University of Connecticut, Storrs 06269, CT, USA

Full list of author information is available at the end of the article



which degrades linear RNAs but not circular RNAs. This treatment can enrich circular RNAs [10, 11].

Computational tools for circular RNA detection have been developed. Currently, there are several existing tools for circular RNA detection, such as Find\_circ [12], CIRCexplorer [13] and CIRC [14]. Find\_circ is one of the first tools for circular RNA detection. Since it is difficult to map the joint position of circular splicing back to the reference genome, Find\_circ tries to collect all un-mapped reads based on reads mapping results from Bowtie [15]. Then, all unmapped reads are converted to new short reads by combining the head and the tail parts of current reads together. Then Find\_circ maps the new short reads back to the reference. CIRCexplorer performs reads mapping using Bowtie and TopHat. The main idea is using the concept of fusion gene to detect circular RNA. First, CIRCexplorer tries to find the un-mapped reads. Then, those un-mapped reads are mapped back to the reference using TopHat-Fusion [16] to detect potential circular RNA candidates with the back-spliced junction reads. CIRC uses BWA [17] for reads mapping, trying to find circular RNA by analyzing CIGAR signatures in the SAM file. Some of these tools such as CIRCexplorer depend on transcriptome annotation, while others support de novo circular RNA detection, such as Find\_circ. Note that often circular RNA comes with the splicing signals of “AG” or “AC” as starting while “GT” or “CT” as the ending for direction “+” and “-” respectively [18]. This can be useful for circular RNA detection. Prior literature also tries to evaluate these tools in terms of their performance, such as precision and sensitivity [19].

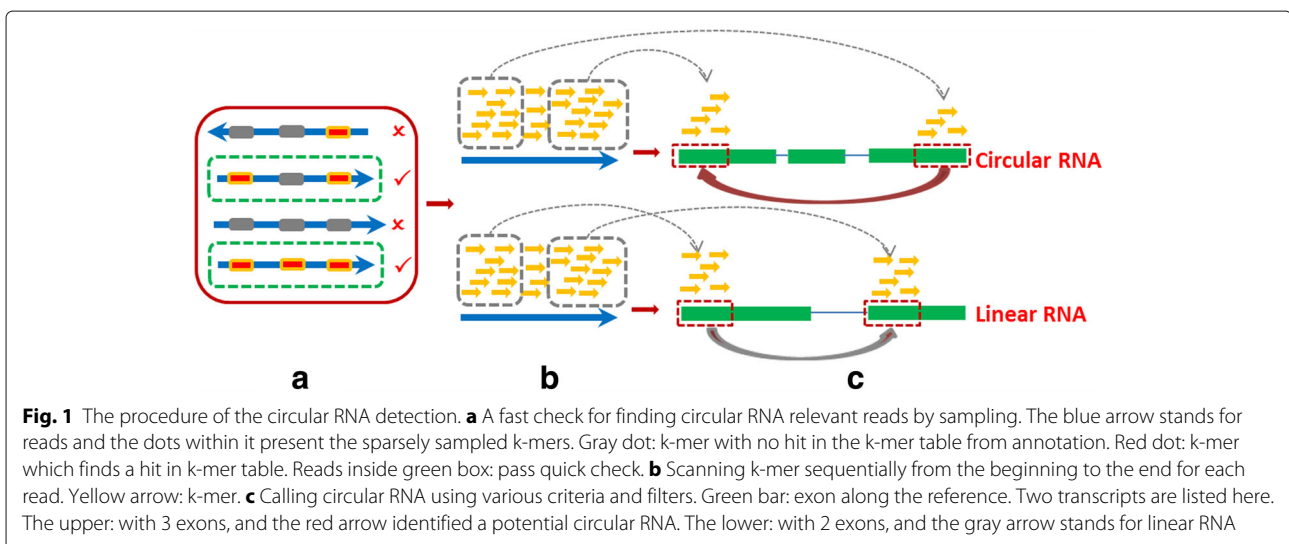
All of these methods mentioned above depend on reads mapping. These mapping based methods have some inherent issues. The first issue is computational efficiency: the existing tools use BWA, Bowtie or TopHat for reads

mapping. Although BWA and Bowtie are widely used in sequence analysis, reads mapping is still time-consuming for circular RNA detection. This is because reads mapping tries to map every read, even when the read is not relevant for circular RNA detection. In addition, since some new short sequences may be created in the middle by circRNA detection tools for the second round mapping, reads mapping can become very slow when TopHat-fusion is used, due to the large length of sequences. Moreover, these tools may miss circular RNA in some cases due to errors in reads mapping. For example, some reads related to circular RNA may be un-mapped due to reads error.

In this paper, we develop a new computational method, called CircMarker, for circular RNA detection. The objective of CircMarker is finding the presence of circular RNA (in particular the join of two known exons). CircMarker doesn't reconstruct the complete sequence of circular RNA. The key idea of CircMarker is that it doesn't rely on reads mapping. Instead, CircMarker analyzes short sequence segments, called k-mers, for circular RNA detection. The main advantage of using k-mers is efficiency: finding k-mers from reads is much faster than reads mapping. Another advantage is that k-mer tolerates more errors in reads and carries useful information about the presence of circular RNA, which may be missed by reads mapping. Empirical results show that CircMarker is more accurate than (or as accurate as) existing methods on simulated and real datasets in calling circular RNA. CircMarker runs much faster than existing methods.

### High-level approach

The overall approach of CircMarker is shown in Fig. 1. CircMarker is based on analyzing k-mers in the sequence reads. That is, CircMarker doesn't perform reads mapping. CircMarker only considers the circular RNA which



**Fig. 1** The procedure of the circular RNA detection. **a** A fast check for finding circular RNA relevant reads by sampling. The blue arrow stands for reads and the dots within it present the sparsely sampled k-mers. Gray dot: k-mer with no hit in the k-mer table from annotation. Red dot: k-mer which finds a hit in k-mer table. Reads inside green box: pass quick check. **b** Scanning k-mer sequentially from the beginning to the end for each read. Yellow arrow: k-mer. **c** Calling circular RNA using various criteria and filters. Green bar: exon along the reference. Two transcripts are listed here. The upper: with 3 exons, and the red arrow identified a potential circular RNA. The lower: with 2 exons, and the gray arrow stands for linear RNA

comes from the exons identified by annotation file. We do not consider de novo circular RNA cases in this paper. CircMarker uses three types of inputs, including the reference genome, the transcription annotation file and sequence reads. Note that all circular splicing that we consider here occurs at the boundary of exons identified by the given annotation file. CircMarker first processes the annotation file and the reference genome. It extracts and stores all k-mers that are located near the exon boundaries. To speed up, CircMarker first performs a fast check to find the reads that are likely to be relevant for circular RNA detection. Then it processes each read and compares k-mers in the read with the stored k-mers to identify circular RNA based on the signatures from circular RNA. When two k-mers from a single read are *out of order* relative to the reference, CircMarker considers this as an evidence for the existence of circular RNA.

### Implementation

#### Processing the reference genome and annotations

CircMarker creates a table for storing the k-mers within the reference genome that are near the exon boundaries as specified by the annotations. The k-mer table is designed to be space-efficient. We only record the following five types of information for each k-mer, including chromosome index, gene index, transcript index, exon index and the part tag as shown in Fig. 2. The “part tag” specifies whether the k-mer comes from the head (i.e. beginning) part or the tail (i.e. ending) part of the exon. Due to the relative small ranges of index, a record on a k-mer only needs eight bytes. We call it the annotation position. One k-mer may contain a group of annotation positions. 32 bits integer is used to store the information of a k-mer, which means the maximum length of a k-mer should be shorter than 16 bp, and all k-mers which contain invalid letters such as “N” are discarded.

When extracting k-mers from annotated exons in the reference genome, we only consider the exons with circular splicing signal in either head or tail part. And we only consider the k-mers which come from the left and right boundaries of the exon. The length of the boundary region  $L_B$  is defined as below:

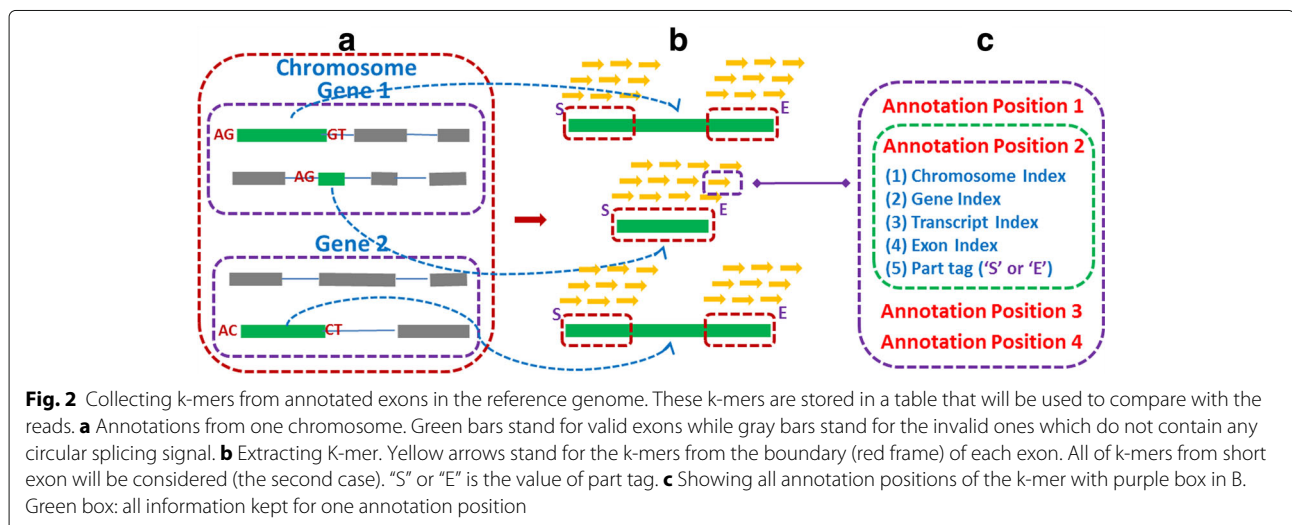
$$L_B = L_R \times R_c$$

$$L_{em} = 2 \times L_B + 2 \times K - 1$$

Here  $K$  is the length of the k-mer.  $L_R$  is the length of reads, and  $R_c$  represents the percentage of reads that should be covered in each boundary. Since we expect more than half reads to be considered, we set the default value of  $R_c$  as 30% ( $2 * 30\% = 60\% > 50\%$ ).  $L_e$  is the length of current exon. If  $L_e \leq L_{em}$ , we use the whole exon to create k-mer and set the part tag as “S” if it located in the first half part and “E” for the second half. Otherwise, we use the head boundary and tail boundary of current exon to extract k-mers and set the part tag to “S” or “E” respectively.

#### Processing sequence reads

Once the k-mer table of the annotated exons is created, we now process each sequence read. Here we examine k-mers contained in a read and search for a match in the k-mer table. This way, we obtain the “hitting status”, which means which transcript can be hit by current reads. “hit exon” means the exon that is hit by the k-mer in the reads in the k-mer table. Each read may be related to more than one hitting status. Each hitting status contains at least one hit exon, and each hit exon should be supported by at least one k-mer. We scan all reads to check their hitting status. In order to skip irrelevant reads, we first perform basic check by sampling eight k-mers from 10% to 80% position of the current read. The read passes the sampling check only if at least two k-mers can find a hit in the k-mer table. If so, we examine all of k-mers from start to end, collecting



all hitting status in this order. Since each hitting status is contributed by multiple k-mers, the “best hitting case” is defined as:

$$\begin{cases} N_h - PreBestHitNum > 5 \\ |N_h - PreBestHitNum| \leq 5 \text{ and } \sum L_{h_e} < PreBestHitLen \end{cases} \quad \text{or}$$

$N_h$  specifies the number of k-mers which supports all of hit exons in one hitting status. The  $PreBestHitNum$  means the  $N_h$  of previous best hitting status.  $L_{h_e}$  means the length of one hit exon in current hitting status, and  $PreBestHitLen$  means the summary length of the hit exons in previous best hitting status. If  $N_h$  is larger than  $PreBestHitNum + 5$ , the current hitting status will be set as the previous best hitting case, which means we prefer the hitting status with conditional larger number of k-mers supporters. Otherwise, the hitting status with the shorter total length of the hit exon will be chosen. We set the previous best hitting status as the final best hitting case when all of hitting status been processed. Finally, the  $N_h$  of best hitting case should be at least 5. Otherwise it is discarded.

**Filtering**

The previous step identifies best hitting cases. Due to the inherent noise in the data (e.g. read errors and duplications), we perform the following filtering step to improve the accuracy. There are two main filtering procedures.

**Filtering procedure one.** The first filtering procedure is checking the hitting number. The minimum hitting number  $N_{h_m}$  is defined as below:

$$N_{h_m} = \begin{cases} \sum L_e - K + 5, & \text{if } \sum L_{h_e} \leq L_{e_m} \\ L_B \times 1.2, & \text{if } \sum L_{h_e} > L_{e_m} \end{cases}$$

The key is that short exons should be fully covered by the reads more than one time. Otherwise, we need to ensure the reads to be within both boundaries of the hit exons. Any best hitting case is discarded if the  $N_h$  is smaller than  $N_{h_m}$ .

**Filtering procedure two.** Based on the number of hit exons in best hitting case, we divide all cases into two types: the case of self-circular if the number of hit exons is equal to 1, and the regular-circular case otherwise.

For self-circular case, only the exon containing the circular splicing signal in both sides will be considered. Then, the best hitting case will be considered as the self-circular RNA candidate if  $L_e \leq L_{e_m}$ . Otherwise, we collect the part tags from begin to end, and condense the tags which belong to the same exon based on the number of hitting. For example, we define  $S(n)$  as n continuous tag “S” in one exon (similar for  $E(n)$ ). If we have  $S(1)$  and  $E(10)$  in one exon, then we condense them to  $E(10)$ . This may help us to filter some random hits. We consider a candidate a valid self-circular RNA only if there are two tags that are

arranged from E to S sequentially (i.e. going backward at the circular RNA join junction).

For the regular-circular case, the best hitting case will be considered only if it contains two exons. First, an exon will be skipped if its hitting time is at most 3 in order to remove some random hits. Then, we try to condense the tags. Here the method described in the self-circular case will be applied at first. After that, for the first exon we condense SE to E while condensing SE to S for the second exon. This condensation logic may help for the case where some of exons are fully covered by current reads. The best hitting case will be kept only if the number of condensed tag in both exons is equal with 1 and the tags arranged from E to S sequentially.

**Calling circular RNA**

There are two cases for calling circular RNA: the self-circular case and the regular-circular case.

**Self-circular RNA.** First, a self-circular RNA candidate will be discarded if the length of current exon is shorter than the read length while the  $N_h$  is smaller than  $L_e - K + 1$ . Otherwise, the best hitting case will be considered to be a valid self-circular RNA candidate if it contains circular splicing signals in both sides.

**Regular-circular RNA.** For the direction “+”, the candidate will be dropped if the exon index increases monotonically. Otherwise, we try to identify the breakpoint at the position of the first decreasing and set it to be the joint junction of circular RNA. We call the exon with large index as the head exon while another one as the tail exon. Based on this definition, the head exon is located in the later part of the reference, while the tail exon is located in the earlier part, and the circle should connect the head exon back to the tail exon. The candidate will be viewed as a valid regular-circular RNA candidate only if the head exon and tail exon have the tail and head circular splicing signal respectively. We set the end position of the head exon and the start position of the tail exon as the position of this called regular-circular RNA.

For the direction “-”, the procedures is almost the same as the direction “+”. The only difference is how to choose the joint junction. In this case, the candidate will be dropped if the exon index decreases monotonically. Otherwise we try to identify the breakpoint at the the first increasing and set it to be the joint junction of circular RNA. The exon with small index is viewed as the head exon while the big index exon is set as the tail exon.

**Refining circular RNA candidates.** We count how many reads support each circular RNA candidate. Only the candidate with support number smaller than the pre-defined threshold will be viewed as the correct one. Since

the maximum coverage of circular RNA is unknown in most cases, we set the default value to be a large number to allow all of valid circular RNA candidates.

## Results

Since the study of circular RNA is still at an early stage, there is no widely accepted benchmark data for evaluating the circular RNA calling at present. Recently, there are some public circular RNA databases which collect different types of circular RNA from published papers. Some databases come with the recommended circular RNA detection tool, such as CircBase [20]. Others focus on collecting the relationship between circular RNA and diseases or traits, such as Circ2Traits [21].

In this paper, we use both simulated and real data to compare CircMarker with three existing tools, including CIRI, Find\_circ, and CIRCexplorer in terms of the number of called circular RNA, accuracy, consensus-based sensitivity, bias and running time. We note that all three tools we compare have customizable parameters. In addition, since all these methods depend on the mapping results coming from different mapping tools, including BWA, Bowtie, TopHat and Tophat fusion, mapping results may impact the accuracy as well. Some tools, such as CIRI, discussed how to optimize the parameters, while other didn't. None of them provided explicit guidelines on how to set parameters for different types of genomes. Therefore, we use the default parameters for these tools including CircMarker in comparison, and we notice that this may lead to some biases to comparison. When comparing the genomic positions of circular joint junction, we allow up to five bp tolerance. Since CircMarker is based on k-mers and each chromosome has its own k-mer table, the running time can be reduced significantly by parallelization (i.e. running analysis on each chromosome in parallel). We compare the performance of these tools on the first three chromosomes individually. Because some existing tools do not support parallelization, we use a single core to run each program for circular RNA detection, and use 10 to 12 cores to run the reads mapping programs such as BWA, Bowtie and TopHat.

### Simulated data

We first use simulated data for evaluation. To generate simulated data, we use the simulation script (called "CIRIsimulator.pl") released by CIRI. The reference genome is the chromosome 1 in human genome (GRCh37). The annotation file is the version 18 (Ensembl 73). Two different cases are simulated as follows: (1) pair-end reads with 13,856,032 sequences, which roughly lead to 10X coverage for circular RNA and 100X coverage for linear RNA, and (2) pair-end reads contains with 9,400,036 sequences, which lead to us 50X coverage for both circular and linear RNA. The goal of the case 1

simulation is simulating the regular RNA-seq, while the case 2 focuses on the situation when the coverage of circular RNA is higher. The reads length is 101 bp and the insert size is 252 bp in both cases. The total number of simulated circular RNA in benchmark is 8033 and 8071 for those two cases respectively. Note that the true circular RNA is known in simulated data, which can be used in comparison. Since the coverage of circular RNA is known in simulated data, we set the "maximum support reads" to be 10 and 50 in CircMarker respectively. We use the following three statistics for comparison: (1) hit number  $N_h$ ; the number of called circular RNA that are true, (2) accuracy:  $\frac{N_h}{N}$  where  $N$  is the total number of called circular RNAs by a method, (3) running time.

The results of the four tools being compared are shown in Fig. 3. Our results show that CircMarker outperforms the existing tools in terms of hit number, accuracy and running time. This is especially evident in case 1 (Fig. 3a), where CircMarker has fewer false positives and also calls more correct circular RNA than other tools. For case 2, the accuracy of CircMarker decreased to 32.04% from 70.90% in case 1. This is likely due to the weak performance of the option "coverage filter", for the similar coverage in both linear and circular RNA. Still, CircMarker is slightly more accurate than existing tools in this case. Moreover, CircMarker runs much faster than existing tools.

### Real data

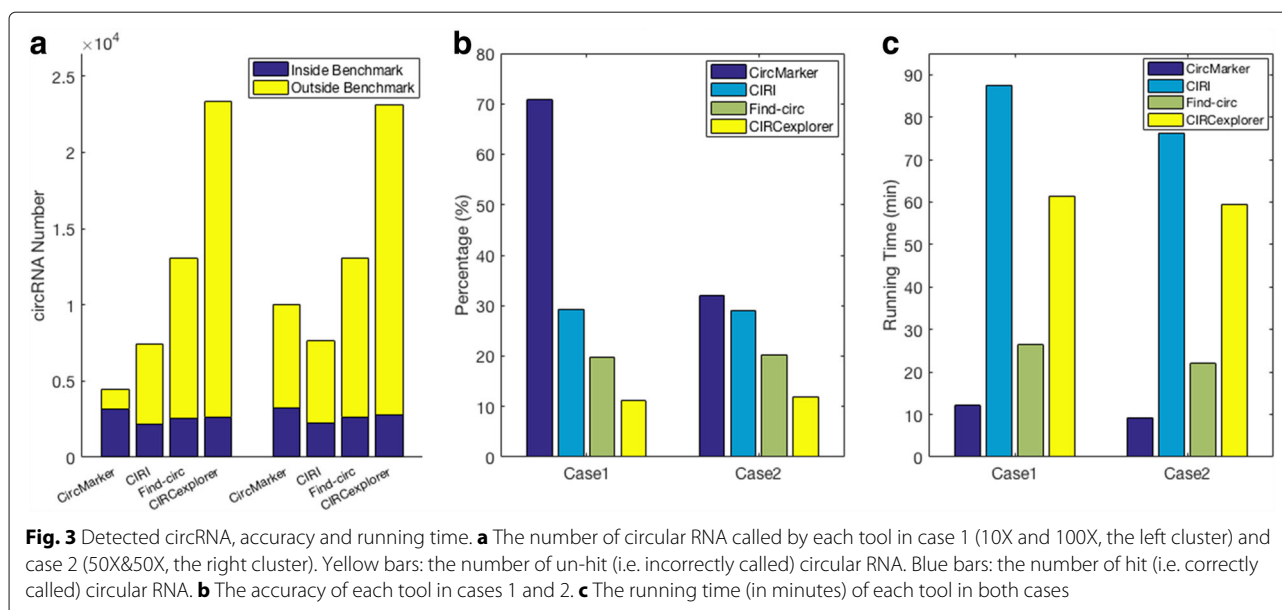
We use two types of real data to evaluate the performance of the four tools.

#### **Real RNase R treated sequence reads with public database information**

As described before, some public databases contain circular RNA called by published papers. In those papers, the authors usually only validate parts of the computationally detected circular RNA using biological experiments. The final result will be released only when the accuracy of those randomly chosen candidates meets certain standard. Therefore, we consider those released circular RNAs in these databases are reliable in this paper.

**Data collection** We choose CircBase [20] as the standard circular RNA database of *homo sapiens*. We use the circular RNAs recorded in this database as "benchmark". The reference genome and annotation file come from *homo sapiens* GRCh37 version 75. The RNA-Seq reads are from SRR901967. These RNA-Seq reads are used to examine circular RNAs from RNase R treated poly(A)-/ribo- RNAs in human embryonic stem cells. There are total 41,342,095 single-end reads in this data.

We use the first three human chromosomes for comparison and use four statistics for comparison. (1) Hit



number  $N_{h_{db}}$ : the number of circular RNA which has a matched circular RNA in the database. These matched circular RNA are called reliable circular RNA. (2) Intersection: the intersection of reliable circular RNA between CircMarker and other tools. This value could be used to evaluate the bias. (3) Reliability ratio:  $\frac{N_{h_{db}}}{N}$ . This measures the fraction of the number of matched circular RNA with regard to the total called ones  $N$ . (4) Running time. The best tool is expected to have large intersection with other tools (low bias), large number of reliable circular RNA with high reliability ratio and fastest running time.

The number of circular RNAs in CircBase from chromosome 1 to chromosome 3 is 9142, 7530 and 5320 respectively. The results show that CircMarker finds more “benchmarked” circular RNAs and runs much faster than others (Fig. 4a, c). For the reliability ratio, there is a trade off with hit number. CIRI obtains the highest reliable ratio, but has the smallest hit number. The reliability ratio of CircMarker is similar to those of Find\_circ and CIRCexplorer (Fig. 4a). CircMarker has the largest hit number. In addition, CircMarker has the large intersection with the results from other tools in all three chromosomes, which means it has low bias (Fig. 4b). As a result, CircMarker outperform the other tools in this data. Moreover, during this experiment, there is no preference for either database selection or comparison approach, therefore, “CircMarker” could be applied to other database as well.

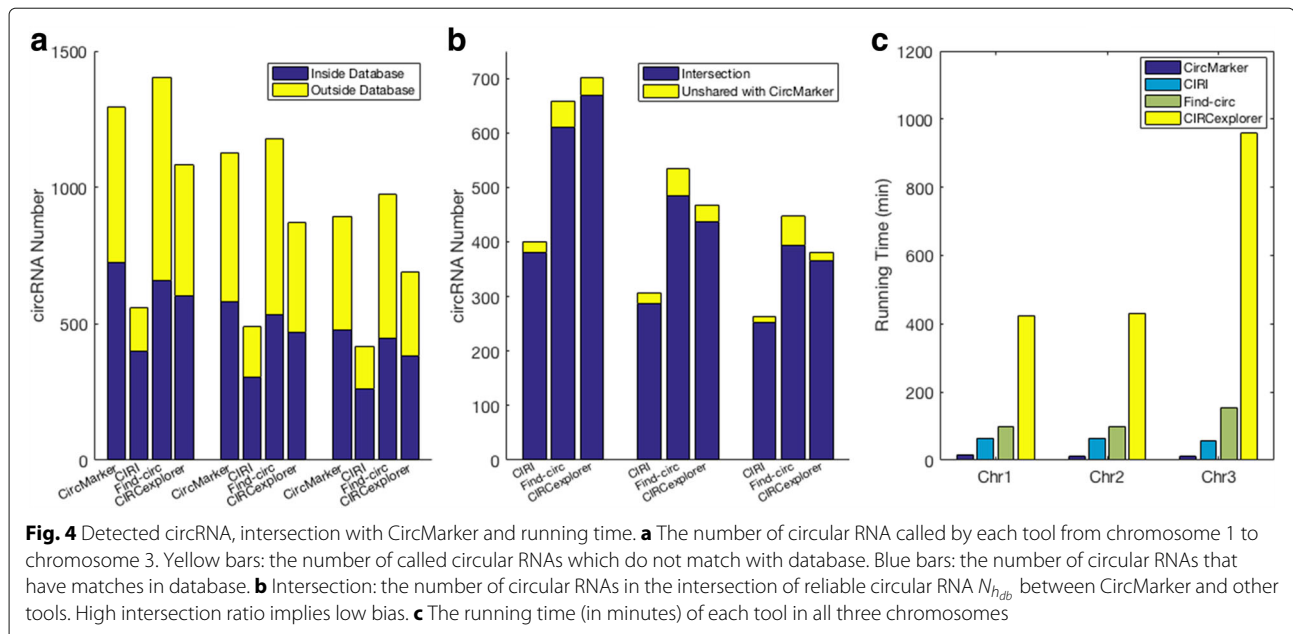
#### Real RNase R treated/untreated data

Recall that RNase R is an experimental technology that can break down linear RNA and enrich circular RNA. As a result, one popular way for validating a circular RNA detection tool is running the tool in two different types of reads: one from only rRNA eliminated sample (called

untreated), and the other from RNase R treated sample. The circular RNA which can be found in both types of reads is considered to be reliable.

**Data collection.** The reference genome and the annotation file are from *Mus Musculus GRCm38 Release79*. The RNase R treated reads are from SRR2219951 and the untreated reads are from SRR2185851. The library was prepared using the script Seq v2 Kit from Epicentre [22], and this data has been used to delineate the circular RNA complement of mouse brain at age 8 to 9 weeks. Both datasets contain pair-end reads, and SRR2219951 (treated) contains 44,661,952 sequences while SRR2185851 (untreated) contains 65,879,618 sequences.

We use the first three chromosomes of *Mus Musculus* with the two types of reads for comparison. We use the following three statistics. (1) Reliable circular RNA: the reliable circular RNAs are from the intersection of called circular RNAs between the treated and untreated reads. Each tool reports its own reliable circular RNA from chromosomes 1 to 3. (2) Consensus-based sensitivity: we say a called circular RNA to be trusted if this circRNA is called by at least two tools. These trusted RNAs are considered to be “benchmark”. We collect these trusted circular RNA for each chromosome. Then, we calculate the intersection between the reliable circular RNA and the benchmark for each tool respectively from chromosome 1 to 3. The consensus-based sensitivity is calculated by:  $\frac{|intersection|}{|benchmark|}$ . (3) Running time. Ideally, a circular RNA detection tool should obtain large number of reliable circRNA with high consensus-based sensitivity and fast running time in each chromosome.



The results are shown in Fig. 5. CircMarker finds larger number of reliable circular RNA than others in all three chromosomes (Fig. 5a). The number of circRNAs in benchmark (i.e. trusted circRNA supported by at least two tools) is 322, 353 and 186 for chromosomes 1 to 3 respectively. One can see that CircMarker gets the largest number of reliable circular RNA in all three chromosomes. In addition, it has the highest consensus-based sensitivity in chromosome 1 and 3, but has slightly lower consensus-based sensitivity than find\_circ in chromosome 2 (Fig. 5b). Moreover, CircMarker only needs around 15 minutes to finish the whole analysis of teated sample while other tools may take at least 1 hour (CIRCexplorer takes more than 9 hours). Overall, CircMarker outperforms the other tools on this data (Additional file 1: Table S1).

For both two verification experiments described above, the data we used here are randomly picked out without any preference. As a result, "CircMrker" could also be used to predict the circular RNAs in other dataset and species if the corresponding annotation file could be well obtained.

## Discussion

CircMarker takes advantage of annotation file to determine the position of the junction point caused by back splicing. Some existing circular RNA calling tools don't use annotation files. There are several advantages of using annotation file. First of all, annotation file contains the boundary positions of each exon, which may help to identify the junction point more accurate than only use the position where back splicing occur, especially for the case when the reads error is near the junction point. Secondly, it can help to filter some false positive cases if the exons involved in the back splicing do not contain

the reasonable splicing signal as expected. Finally, we can choose some parts of sequence from the boundary of each exon identified by the annotation file to build the k-mer table, which may improve the speed significantly. Moreover, the circular RNA which is supported by annotated exons should be considered as more reliable than the de novo one.

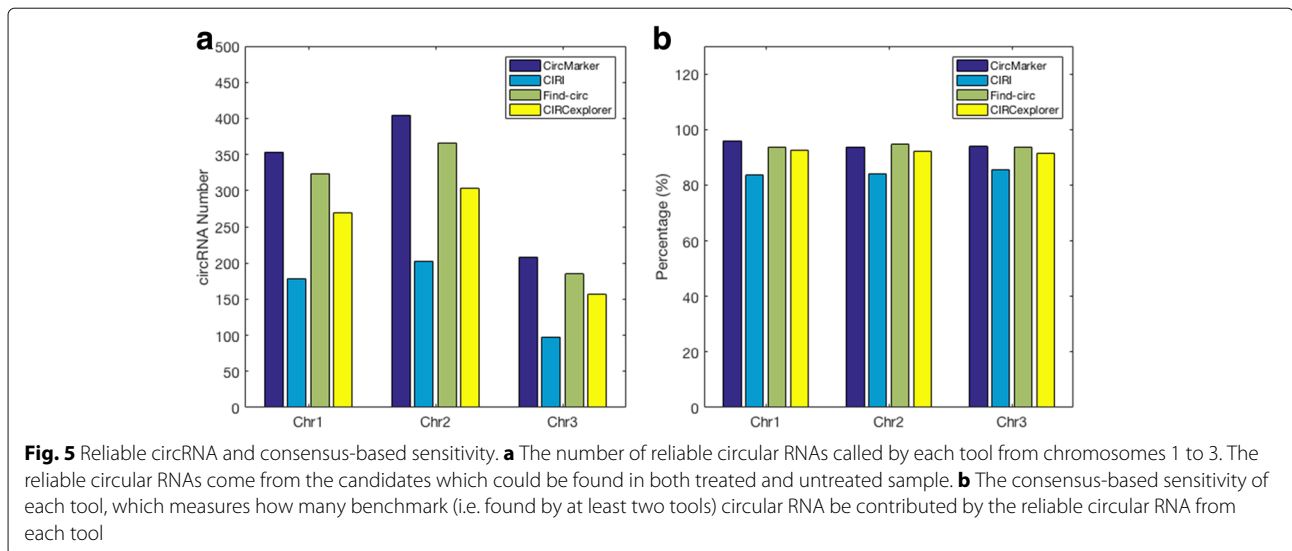
On the other hand, a major disadvantage is that, since CircMarker depends on annotation file, it may miss the de novo circular RNAs which occur in unannotated exons. As a result, CircMarker cannot handle the case or perform as well as expected if the annotation file is not given or the quality of annotation file is not good. In addition, some circular RNAs with back splicing junction points within intron may not be detected by CircMarker.

In order to evaluate the limitation of annotation dependence in regular datasets, we made some statistical calculations for the circular RNAs recorded in circBase. We find there are 91.2% circular RNA recorded in circBase that could hit the boundary of exons recorded in annotation file. In addition, when we perform comparison, the benchmark dataset includes both de novo circRNA and annotation based circRNA. The results show that CircMarker still outperforms other tools even the results predicted by some of other tools include both de novo and annotation based circRNAs.

Based on our experiments described above, a large number of circular RNAs exactly occur in annotated exons.

## Conclusion

In this paper, we develop a new circular RNA detection method called CircMarker based on k-mer analysis. CircMarker runs much faster than other tools because it



doesn't perform reads mapping. Moreover, k-mers contain useful information about circular RNA detection. Our results on both simulation data and real data demonstrate that CircMarker can find more circular RNA. It has higher consensus-based sensitivity and high accuracy/reliable ratio compared with others. In addition, the circular RNAs called by CircMarker often contain most circular RNAs called by other tools in the real data we tested. This implies that CircMarker has low bias.

CircMarker is easy for use. CircMarker is a stand-alone tool (implemented by C++) and does not depend on any third party tools. The source code is available under the GPLv3 licence at <https://github.com/lxwgcool/CircMarker>.

### Availability and requirements

**Project name:** CircMarker.

**Project home page:** <https://github.com/lxwgcool/CircMarker>

**Operating system(s):** Unix.

**Programming language:** C++.

**Other requirements:** GCC.

**License:** GPLv3.

### Additional file

**Additional file 1: Table S1.** Contains the detailed results presented in this paper. These include seven tables: one for the results on simulated data, three for public database validation on real RNase R treated reads on the first three chromosomes respectively, and three for the consensus-based validation on real RNase R treated and untreated reads on the first three chromosomes respectively. (XLSX 18 kb)

### Abbreviations

BWA: Burrows-wheeler aligner; CircRNA: Circular RNA; mRNA: Messenger RNA; RNA: Ribonucleic acid

### Funding

Publication of this article was funded by grants IIS-0953563, IIS-1447711 and IIS-1526415 from US National Science Foundation to YW.

### Availability of data and materials

The software is available under the GPLv3 licence at <https://github.com/lxwgcool/CircMarker>.

### About this supplement

This article has been published as part of *BMC Genomics* Volume 19 Supplement 6, 2018: Selected articles from the 13th International Symposium on Bioinformatics Research and Applications (ISBRA 2017): genomics. The full contents of the supplement are available online at <https://bmcbgenomics.biomedcentral.com/articles/supplements/volume-19-supplement-6>.

### Authors' contributions

XL designed algorithms, developed software, performed analysis and experiments, and wrote the paper. CC contributed to the original idea: using alignment-free method (not from alignment) to call circular events. JP contributed to experiments and wrote the paper. IM and YW designed the algorithms, wrote the paper and supervised the project. All authors have read and approved the final manuscript.

### Ethics approval and consent to participate

Not applicable

### Consent for publication

Not applicable

### Competing interests

The authors declare that they have no competing interests.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

### Author details

<sup>1</sup>Department of Computer Science and Engineering, University of Connecticut, Storrs 06269, CT, USA. <sup>2</sup>Department of Biomedical Informatics, Harvard Medical School, Boston 02115, MA, USA.

Published: 13 August 2018

### References

- Nigro JM, Cho KR, Fearon ER, Kern SE, Ruppert JM, Oliner JD, Kinzler KW, Vogelstein B. Scrambled exons. *Cell*. 1991;64(3):607–13.



2. Dignam JD, Lebovitz RM, Roeder RG. Accurate transcription initiation by rna polymerase ii in a soluble extract from isolated mammalian nuclei. *Nucleic Acids Res.* 1983;11(5):1475–89.
3. Ozsolak F, Milos PM. Rna sequencing: advances, challenges and opportunities. *Nat Rev Genet.* 2011;12(2):87–98.
4. Cocquerelle C, Mascrez B, Hetuin D, Bailleul B. Mis-splicing yields circular rna molecules. *FASEB J.* 1993;7(1):155–60.
5. Gao Y, Wang J, Zheng Y, Zhang J, Chen S, Zhao F. Comprehensive identification of internal structure and alternative splicing events in circular rnas. *Nat Commun.* 2016;7:.
6. Salzman J, Gawad C, Wang PL, Lacayo N, Brown PO. Circular rnas are the predominant transcript isoform from hundreds of human genes in diverse cell types. *PLoS ONE.* 2012;7(2):30733.
7. Houseley JM, Garcia-Casado Z, Pascual M, Paricio N, O'dell KM, Monckton DG, Artero RD. Noncanonical rnas from transcripts of the drosophila muscleblind gene. *J Hered.* 2006;97(3):253–60.
8. Hansen TB, Jensen TI, Clausen BH, Bramsen JB, Finsen B, Damgaard CK, Kjems J. Natural rna circles function as efficient microrna sponges. *Nature.* 2013;495(7441):384–8.
9. Westholm JO, Miura P, Olson S, Shenker S, Joseph B, Sanfilippo P, Celniker SE, Graveley BR, Lai EC. Genome-wide analysis of drosophila circular rnas reveals their structural and sequence properties and age-dependent neural accumulation. *Cell Rep.* 2014;9(5):1966–80.
10. Jeck WR, Sorrentino JA, Wang K, Slevin MK, Burd CE, Liu J, Marzluff WF, Sharpless NE. Circular rnas are abundant, conserved, and associated with alu repeats. *Rna.* 2013;19(2):141–57.
11. Szabo L, Salzman J. Detecting circular rnas: bioinformatic and experimental challenges. *Nat Rev Genet.* 2016;17(11):679–92.
12. Memczak S, Jens M, Elefsinioti A, Torti F, Krueger J, Rybak A, Maier L, Mackowiak SD, Gregersen LH, Munschauer M, et al. Circular rnas are a large class of animal rnas with regulatory potency. *Nature.* 2013;495(7441):333–8.
13. Zhang X-O, Wang H-B, Zhang Y, Lu X, Chen L-L, Yang L. Complementary sequence-mediated exon circularization. *Cell.* 2014;159(1):134–47.
14. Gao Y, Wang J, Zhao F. Ciri: an efficient and unbiased algorithm for de novo circular rna identification. *Genome Biol.* 2015;16(1):4.
15. Langmead B, Salzberg SL. Fast gapped-read alignment with bowtie 2. *Nat Methods.* 2012;9(4):357–9.
16. Kim D, Salzberg SL. Tophat-fusion: an algorithm for discovery of novel fusion transcripts. *Genome Biol.* 2011;12(8):72.
17. Li H, Durbin R. Fast and accurate short read alignment with burrows–wheeler transform. *Bioinformatics.* 2009;25(14):1754–60.
18. Jeck WR, Sharpless NE. Detecting and characterizing circular rnas. *Nat Biotechnol.* 2014;32(5):453–61.
19. Zeng X, Lin W, Guo M, Zou Q. A comprehensive overview and evaluation of circular rna detection tools. *PLoS Comput Biol.* 2017;13(6):1005420.
20. Glažar P, Papavasileiou P, Rajewsky N. circbase: a database for circular rnas. *Rna.* 2014;20(11):1666–70.
21. Ghosal S, Das S, Sen R, Basak P, Chakrabarti J. Circ2traits: a comprehensive database for circular rna potentially associated with disease and traits. *Front Genet.* 2013;4:283.
22. Li L, Pesavento PA, Leutenegger CM, Estrada M, Coffey LL, Naccache SN, Samayoa E, Chiu C, Qiu J, Wang C, et al. A novel bocavirus in canine liver. *Virology.* 2013;50(1):54.

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

