

SOFTWARE

Open Access



Software evaluation for de novo detection of transposons

Matias Rodriguez and Wojciech Makalowski*

Abstract

Transposable elements (TEs) are major genomic components in most eukaryotic genomes and play an important role in genome evolution. However, despite their relevance the identification of TEs is not an easy task and a number of tools were developed to tackle this problem. To better understand how they perform, we tested several widely used tools for de novo TE detection and compared their performance on both simulated data and well curated genomic sequences. As expected, tools that build TE-models performed better than k-mer counting ones, with RepeatModeler beating competitors in most datasets. However, there is a tendency for most tools to identify TE-regions in a fragmented manner and it is also frequent that small TEs or fragmented TEs are not detected. Consequently, the identification of TEs is still a challenging endeavor and it requires a significant manual curation by an experienced expert. The results will be helpful for identifying common issues associated with TE-annotation and for evaluating how comparable are the results obtained with different tools.

Introduction

The vast majority of eukaryotic genomes contain a high number of repetitive DNA sequences. These sequences can be broadly classified as tandem repeats or interspersed repeats. Tandem repeats are short sequences with a length up to a few dozen bases that lie adjacent one to another and are approximate copies of the same pattern of nucleotides. Similarly, interspersed repeats are homologous DNA sequences that can be found in multiple copies scattered throughout a genome and their lengths can vary immensely from a hundred nucleotides up to more than twenty-thousand nucleotides. Most of these interspersed repetitive sequences found in genomes originated in the proliferation of transposable elements.

Transposable elements (TEs) are mobile genetic sequences possibly related to viral components that have evolved the ability to increase their abundance in a genome by making copies of themselves. The fraction of TEs in a genome can vary widely and can represent more

than 80% of plant genomes [28]. To put into perspective how common they are, if we consider a well-studied case, such as the human genome, the annotated protein-coding genes represent only a very small fraction of approximately 5% of all the sequences, meanwhile TEs can make up to about 68% of the sequences [5].

Genomes and TEs have coevolved similarly to a host-parasite relationship and this led the genomes to develop multiple mechanisms to suppress TE activity as they can compromise the integrity of the genome and can cause deleterious mutations. Consequently, there is a constant evolutionary arms race between transposon activity and the host genome trying to suppress their proliferation [15]. Despite the parasitic nature of TEs, they play a fundamental role in genome evolution, contributing to plasticity, shaping, and altering the architecture of the genome. TEs contribute to gene regulatory networks as their activity can disrupt regulatory sequences modifying gene expression by altering chromatin structure, behaving as enhancers or promoters, or, when transcribed as part of a larger transcript, creating new transcript isoforms altering splicing and mRNA stability [18]. There are multiple examples of TEs that have been domesticated

*Correspondence: wojmak@uni-muenster.de
Institute of Bioinformatics, Faculty of Medicine, University of Münster,
48149 Münster, Germany



and proteins derived from them which were co-opted, such as the RAG1 gene from the somatic V(D)J recombination in humans and the retrotransposons that maintain the telomeres in *Drosophila*. RNA-mediated retrotransposition of transcribed genes is also a source for gene duplications that can lead to novel traits [17].

Historically, TEs were considered useless selfish sequences and their influence on genes and genomes was often dismissed [21]. It was not until the last two decades that they started to be considered as major components of genomes and important players of genome evolution, but due to the difficulties posed by their repetitive nature their annotation and role in genetic studies still continues to be neglected [3].

The correct identification of TEs is an important step in any genome project since their repetitive nature can create difficulties during de novo genome assemblies, breaking the continuity of contigs as a result of the same reads mapping to multiple loci [25]. They can also hinder annotation by creating conflicts with gene prediction programs if they can be found inside a host gene, carry part of a host gene when replicating, become pseudogenes, or contain spurious ORFs.

There are multiple tools for TE-detection but there are no clearly defined pipelines or software tools that could be considered as standard, as there are no clear metrics to compare the results obtained from each software [12]. Most tools also rely on a high copy number of elements for correct identification and are usually tested in organisms that have large genomes and a high abundance of TEs.

The identification of TEs can be really daunting and a time-consuming endeavor for the amount of data that needs to be processed and compared and the challenges inherent to their complex nature. TEs are extremely diverse, they comprise multiple classes of elements that can vary immensely in sequence, length, structure, and distribution [31]. Some TE families found in eukaryotic genomes can be very old with a majority of inactive copies due to accumulation of mutations or fragmentation during the insertion process. This means that remains of antique copies from a family can be very divergent from active TEs, making the detection of the remnants of decayed copies or the definition of consensus sequences a real challenge that is hindered by the great variability of TEs within the same family. The proliferation of TEs can also result in the generation of nested TEs and some families show a clear preference for jumping into other TEs that act as hotspots for insertion [8], making the detection and correct annotation of them even more difficult.

There are well curated TE databases, such as RepBase [1] or Dfam [14], with libraries of consensus sequences. A homology-based approach relies on the TE sequences

from these libraries which are then mapped against the studied genome. To identify new TEs a de novo approach is used and there are abundant software alternatives which rely on different strategies ranging from structural information, periodicity, k-mer counting, or repetitiveness, among others [19]. When a new species is sequenced, a strategy which uses only information from curated databases is not enough and it is necessary to use a de novo strategy to identify novel families and species-specific TEs.

In this work we compare TE detection software which are widely used by researchers and we assess their performances on genomes with well curated TE annotations. We ran a number of de novo TE detection software packages on simulated sequences and genome sequences and then compared and evaluated their performance in detecting a wide variety of divergent TE families. A particular scenario that we tried to consider is the detection of transposons in smaller genomes of around a hundred million bases. In all cases the software for identification of TE rely on the presence of a large number of elements of the same family and that is usually not common in smaller genomes where a lower number of copies of TEs is expected.

Methods

Datasets

Genomic data with annotated TEs were downloaded from the UCSC Genome Browser database [11]. The TE annotation provided by UCSC Genomes was obtained from mapping TEs from the RepBase database [1] against each genome using RepeatMasker [29]. The sequence data sets we used varied from 46.7 Mb for the human chromosome 21 and 137.5 Mb for the fruit fly genome (see Table 1).

Simulated data

We used a Python script to simulate an ideal scenario where the composition, coordinates, and divergence of all the TEs are already known. This script takes input from a configuration file for GC content, TE sequences, number of copies, expected divergence (mutations and indels), percent of fragments, and nesting. The script starts by simulating a random sequence of a predefined length and a GC content that constitutes the base sequence where TEs are going to be inserted. Then it obtains the name of each TE and the number of copies from the configuration file and random positions are chosen and assigned to each TE. In the next step, TE sequences are loaded from a library and the information about divergence and fragmentation is taken into account to generate random mutations and fragments that are inserted into the base sequences. The last step takes all of this

Table 1 Datasets used for testing the TE de novo detection tools

Dataset	Genome version	Sequences	Length	TE-fraction
Human genome	GRCh38/hg38 p12	Chromosome 21	46.7Mb	36.74%
Zebrafish	GRCz11/danRe 11	Chromosome 1	59.6Mb	46.58%
Fruit fly	BDGP Release 6	Whole genome	137.5 Mb	17.55%
Simulated data	n/a	Simulation	100.1 Mb	60.00%

information to generate a fasta file with the whole new sequence with TEs and a GFF file with all the coordinates and relevant information. For the simulated dataset we used a base sequence of 40Mb with the GC content similar to the human genome (42%) and inserted 60 Mb of TE-sequences from 20 different families downloaded from the Dfam database [14]. Although the divergence threshold for individual copies inserted into simulated genome was set to 30%, the majority of the sequences were between 90 and 95% identical to the cognate TE consensus sequences (see Fig. 1). Detailed information on inserted sequences is provided in Table S1.

Software

In this work we compared strategies for de novo detection of TEs using k-mer based tools and programs that construct TE-models. We tested three k-mer counting tools: Red, P-Clouds version 0.9, and phRAIDER. Due to the nature of the algorithms employed by k-mer counting software, these tools are extremely fast and usually don't require much computational power. Nevertheless, they usually require a big amount of RAM to store data structures, so they may not scale up well with large genomes. Red identifies candidate repetitive regions giving them a score, then processes these results using signal processing and the second derivative. These filtered data are used to train a Hidden Markov Model that scans the genome for candidate TEs. As described by the author, it is a novel repeat discovery system that trains itself automatically on an input genome [9]. P-Clouds counts oligonucleotides, then arranges them into clusters of "probability clouds" that are related oligonucleotides that occur as a group more often than expected by chance. Then it annotates the genome by finding stretches with a high density of oligos present in these "probability clouds" [10]. The premise of phRAIDER is to use spaced seeds to specify match patterns, i.e., to permit the search of substrings allowing mismatches in certain positions. Then it scans the genome searching for highly frequent seeds and how they overlap [27]. The biggest limitation of these tools is the fact that they don't make any attempt to classify found repeats. Moreover, there is no information provided on relation between the detected individual elements and no consensus sequences are computed.

We also compared three model-builders: RepeatScout version 1.0.6, REPET (TEdenovo pipeline) version 2.5, and RepeatModeler version 2.0. RepeatScout uses high frequency seeds and extends each seed to a progressively longer consensus sequence, following the dynamically inferred alignments between the consensus sequence and its occurrences in the genome. The alignment score encourages boundaries shared by some but not necessarily all alignments; it uses a standard SW-algorithm to extend until n-iterations fail to improve the score [23]. REPET is a package consisting of two pipelines, one for detection of TEs: TEdenovo, and another for their annotation: TEannot [6, 24]. Both of these pipelines are fully configurable and each step can be parametrized. The TEdenovo pipeline by default starts self-comparison of the input genome with BLASTER, a modified version of BLAST. Then it clusters the high scoring pairs using three tools: RECON, GROUPEr, and PILER, grouping closely related TE sequences. Finally, it performs a multiple alignment using MAFFT or MAP with the aim of having a consensus sequence for each TE family. Here, we are interested in the ability of the software to detect TEs, so we used only the TEdenovo pipeline. Finally, RepeatModeler is a pipeline that uses as an input the outputs of three other software, namely RECON, RepeatScout, and Tandem Repeats Finder. Additionally, it uses LTRHarvest or LTRretriever for LTR-TE detection [7]. It should be noted that RepeatModeler contains an optional module (option -LTRStruct) that enables clustering redundant LTR models. However, since we used only default parameters, this step was skipped in our analyses. RepeatScout, RepeatModeler, and REPET all give as a final result a fasta file with a consensus sequence for each type of TEs they could identify. Afterwards, we need to map back these consensus sequences to the genome to identify individual copies of TEs and get their coordinates. For this step we used the popular tool RepeatMasker, version 4.1.0 [29], using the three tools' outputs as libraries to screen the genomes for TEs.

Another tool used for detecting simple repeats was Tandem Repeat Finder (TRF) version 4.09 [2], a software which models tandem repeats using a probabilistic model. We used it to filter out simple repeats obtained by the k-mer counting tools and also to assess the ability

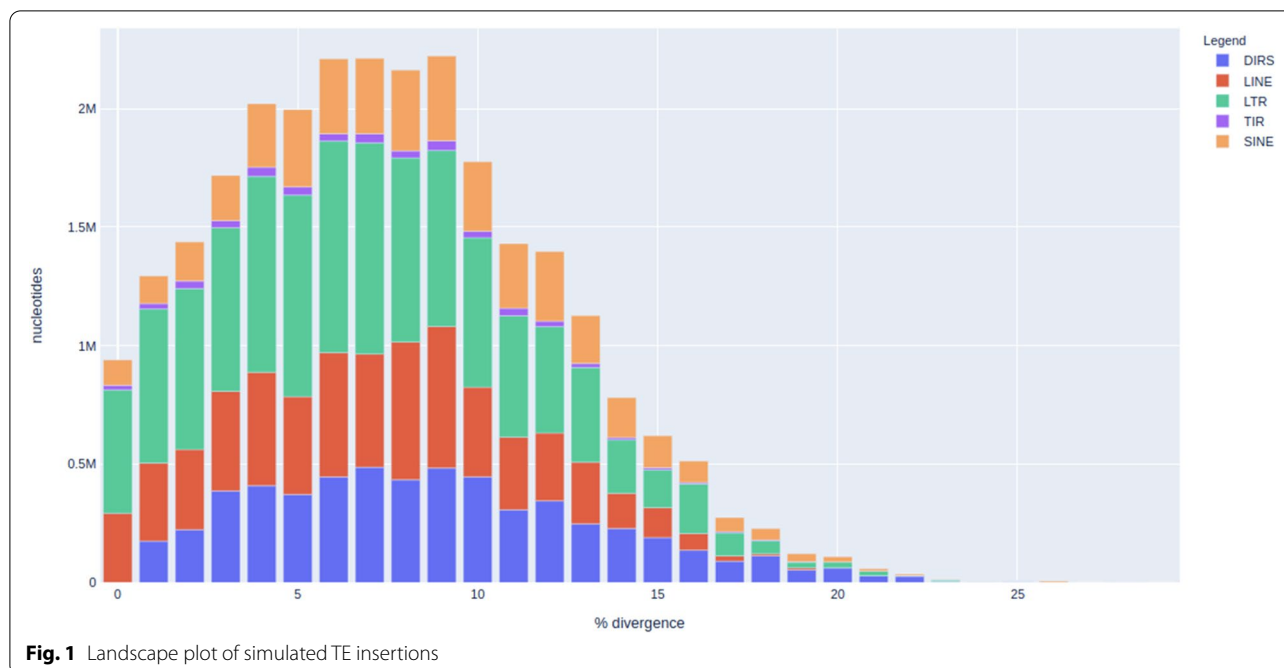


Fig. 1 Landscape plot of simulated TE insertions

of model-builders to cope with simple sequence repeats. TRF was run with default parameters on all the sequences analyzed and the results obtained were merged when an adjacent or overlapping annotation was reported. These results were converted into a GFF file for easing further analysis.

Pipeline

All the software were tested with default parameters as we intended to compare the typical performance of each tool without tuning their optional parameters (see Fig. 2). K-mer counting methods are expected to find all the high frequency k-mers, including simple repeats and interspersed repeats. For k-mer counting software after getting the results we ran Tandem Repeat Finder (TRF) with default parameters to filter out tandem repeats. The model-builders usually include some steps for filtering simple repeats and frequently make use of TRF, so this step was not replicated with these software.

The type of results produced by each set of tools is also quite different due to the different strategies used. K-mer counting tools return the coordinates in the genome with the regions where high frequency k-mers were found, meanwhile model-building software returns the consensus sequences of the TE candidates found as a fasta file. So for the latter set of tools we mapped back the consensus sequences against the original sequences using RepeatMasker, running it with default parameters. The results obtained from all the different tools were transformed to GFF format for further processing.

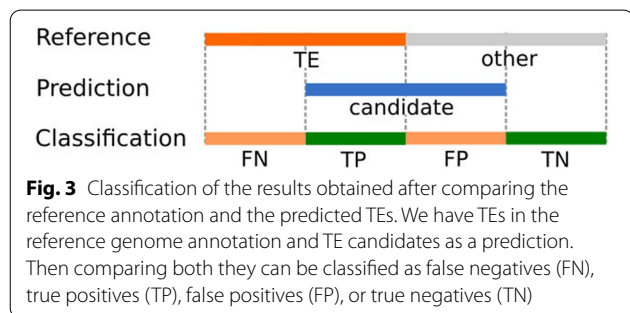
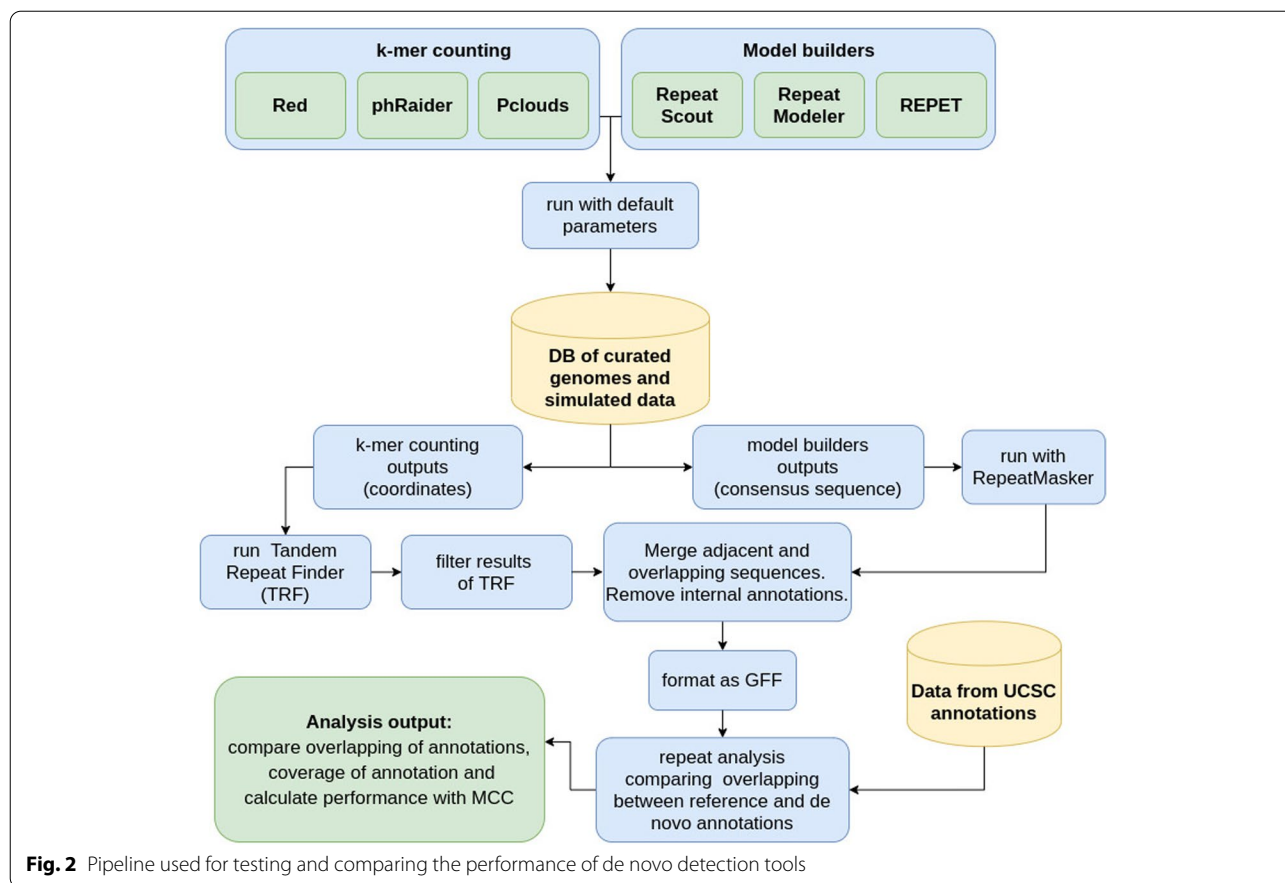
Then GFF files were sorted by coordinates and immediately adjacent, overlapped, or internal coordinates were merged as one, as the main idea is only the identification of transposon sequences. This step is necessary particularly in k-mer counting software which have the tendency to annotate many overlapping and fragmented repeats. For all the datasets tested we have as a reference the annotations downloaded from the UCSC Genomes database.

TE-models' comparison

TE-models generated by model-building software on simulated data were compared to original TE-families using blastn (version 2.10.1+) with default parameters.

Analysis

With the results of TE detection obtained from each software we created GFF files that were then compared against the original files from UCSC database with RepeatMasker mapping results. A custom Python script was used to obtain the overlapping regions of two GFF files and where there are differences in the annotation, it allowed us to compare the coordinates of the reference and the ones obtained by the TE de novo software. This way of evaluating the data is useful in order to create a confusion matrix that can be used as input for a binary classifier test that allows us to compare the performance of different software against a reference. When the reference annotation and the new annotation agree on the coordinates, these bases are counted as true positives,



or if nothing is annotated in both, these bases count as true negatives (Fig. 3). If the new annotation has bases not covered by the reference annotation, we consider them as false positives, and similarly if annotations in the reference are missing in the new annotation, these are counted as false negatives (see Fig. 3).

With this kind of data we have a binary classification problem, where each category can be classified using a confusion matrix. There are multiple tests to evaluate and compare the results obtained by a binary classifier which make use of a confusion matrix data and one of the most commonly used methods is the Matthews Correlation

Coefficient (MCC). The MCC has the advantage that it uses all four quadrants of a confusion matrix considering the proportions of each class and requires that in both classes negative and positive elements are correctly classified, performing well even when using imbalanced data and when one class is underrepresented [4]. The MCC evaluates the results obtained from a prediction, as in this case the TE de novo software TE candidates, against the known annotated data. The values of MCC range from -1 to $+1$, where a value of -1 is obtained when all the predictions are wrong, 0 when results are not better than random guessing, and 1 where all predictions are correct. In this work we used the MCC as a measure of the performance obtained from the different software tools. Additionally, we developed several R scripts for plotting GFF coordinates which visually compare the annotations obtained from each tool.

Results

As mentioned above, the two groups of programs provide different types of results. While k-mer counting software provided a list of regions that are occupied by repetitive sequences, model-building software analyzed here returned sets of repeats' models. These models can

be next used to scan a genome and annotate individual repeats, including TEs. For this step, we used a popular program, namely RepeatMasker (see Methods section).

Model building

Three different programs were used to create TE-models for both real genomic data and simulated sequences. The results of the latter are the most informative as we knew the exact number of expected TE-families. Interestingly, all three programs generated more TE-models than we used for the simulation (see Table 2).

RepeatScout generated the smallest number of TE-models (30) and only in few cases more than two models for a given TE-family: three for L1 and four for Polinton. However, it has a tendency to create homo-dimeric elements, for instance Copia, DIRS, HERVL (see Table S1). On the other extremum lies REPET, which created the highest number of models, although it failed to report an Alu model. This is a bit surprising since there were 730 Alu insertions of a single sub-family (Alu Y). REPET not only generated the highest number of models but some of them were dimeric and hybrid. The latter were caused by a few nested repeats, for instance a Jockey nested in

a Ngaro or a Tc1-Mariner nested in a HERVL. In general, longer elements tend to give rise to several models by both RepeatModeler and REPET. For instance, 5.5kb long L1 element is a source of six models in RepeatModeler analysis and nine models in the case of REPET. Polinton, which is 18.5kb, resulted in eight models in both RepeatModeler and REPET and four models in the case of RepeatScout (see Table S1). Interestingly, some of these models overlap each other, suggesting that they could be merged during manual curation.

In our simulation, we “mutated” individual TE-copies up to 30 % divergence from the reference sequence and many of the individual copies were truncated at the 5’ end. In general, a consensus sequence recovery at the nucleotide level was very good, with the average sequence identity of models to their respective reference sequences at 97.3% (stdev = 3.76). However, many of TEs were broken by a given software into several models. Probably the best example is Polinton, based on which RepeatModeler and REPET created nine models each and none of these covered the whole Polinton sequence that was used in the simulation (see Fig. 4). The shortest transposon inserted into the simulated data was the 311 nucleotide long AluY element. The individual sequences were “mutated” to average 13% divergence from the reference sequence and 67 of them were truncated at their 5’ end up to 30% of the sequence length. RepeatScout performed the best, returning a 308 nt long model with the sequence identical to the reference and just 3 nt missing from the 5’ end. Surprisingly, REPET didn’t report any models based on these sequences. Finally, RepeatModeler created three different models: one almost ideal with just a 5’ terminal guanosine missing and two others, a bit shorter but with extra eight and thirty-five nucleotides added to their 5’ end (see Table S1).

When analyzing the simulated data, one trend became clear, namely that on occasions multiple models are

Table 2 Number of TE models generated by each software from a simulated sequence containing TEs of 20 different families (see Table S1 for detailed information). Please note that REPET failed to generate an Alu model

Software	Number of models generated
RepeatScout	30
RepeatModeler	80
REPET	82
Number of TE families inserted in the simulated data	20

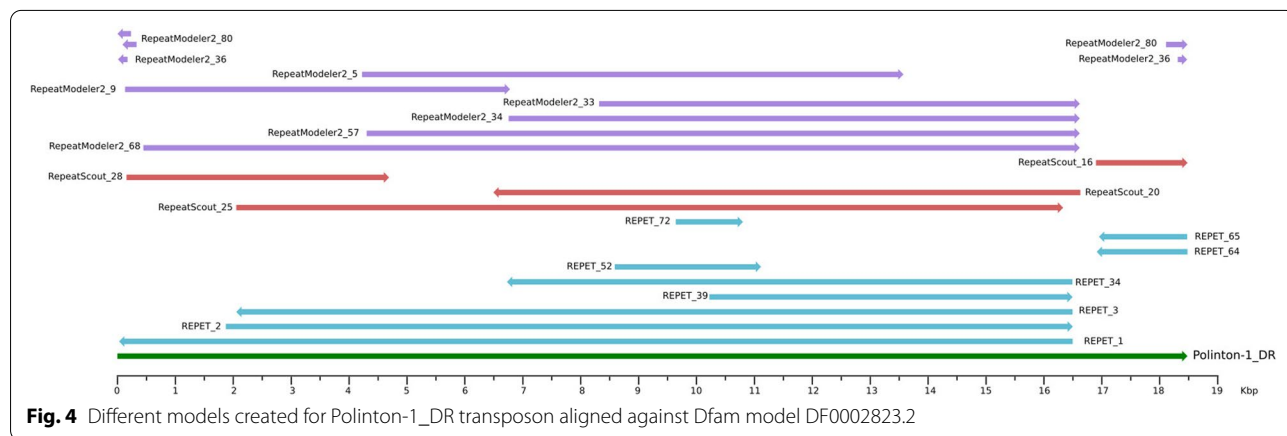


Fig. 4 Different models created for Polinton-1_DR transposon aligned against Dfam model DF0002823.2

generated for the same TE and there are common patterns observed for each software. A characteristic of RepeatModeler is that it tends to generate redundant models, with up to six to eight models for the longest TEs. This is also a common behavior observed with REPET, where many fragments were generated. Another interesting observation that only occurs with REPET is that in some models part of a nested TE was included into a model resulting in chimeric models. With RepeatScout there is much less redundancy with the number of models, but again something unique happens and some of the reported models are total or partial duplicates of the original TE.

An example of different models generated for a DIRS TE of 5.6kb which were particularly difficult to resolve is shown in Fig. 5. This particular TE RepeatModeler generated six different models of different lengths. These models are on average 15 % diverged at the sequence level. There were two models calculated by RepeatScout, one almost identical to the reference and another one almost twice the length of the original and consisting of a duplication leading to erroneous homodimer. Interestingly, the two copies of this homodimer are complementary to each other as they lie on opposite strands compared to the reference sequence. REPET reported four models in total. Two are relatively short, encompassing about one-third of the reference sequence and partially overlapping in a head-to-tail orientation. Another model is a hybrid TE consisting of two overlapping fragments of DIRS element in a head-to-tail orientation with a fragment of TRANSIB transposon. The fourth model closely resembles the original TE but is truncated by about 240 nucleotides at its 5' end (see Fig. 5).

In Fig. 4 we present models generated for another TE, namely Polinton-1_DR [16]. The full-length transposon is 18.5kb, including 350bp terminal inverted repeats. All three software systems compared in our study reported few models for this TE but none of the models recovered the full length TE (see Fig. 4 and Table S1). Interestingly, both REPET and RepeatModeler generated similar close-to-full-length models that are missing one of the inverted repeats, while RepeatScout's longest model misses both inverted repeats. However, inverted repeats were reported as separate models by each of the program.

In the real data from model organisms, RepeatScout created the highest number of models with almost three-thousand consensus for zebrafish chromosome 1 (see Table 3). This is in contrast to the simulated data where RepeatScout generated the least number of models. REPET lies on the other extreme of the spectrum with just 65 TE models for the human chromosome 21, including Alu model that was missing from the simulated data analysis. Interestingly, this chromosome is annotated with almost 1000 different TE families. The smaller

Table 3 Number of models generated by three software in real sequence analysis and number of TE-families annotated in publicly available data on the same sequences

Software	Zebrafish chromosome 1	Fruit fly genome	Human chromosome 21
RepeatScout	2919	2593	464
RepeatModeler	1779	686	428
REPET	342	557	65
Public annotation	875	219	979

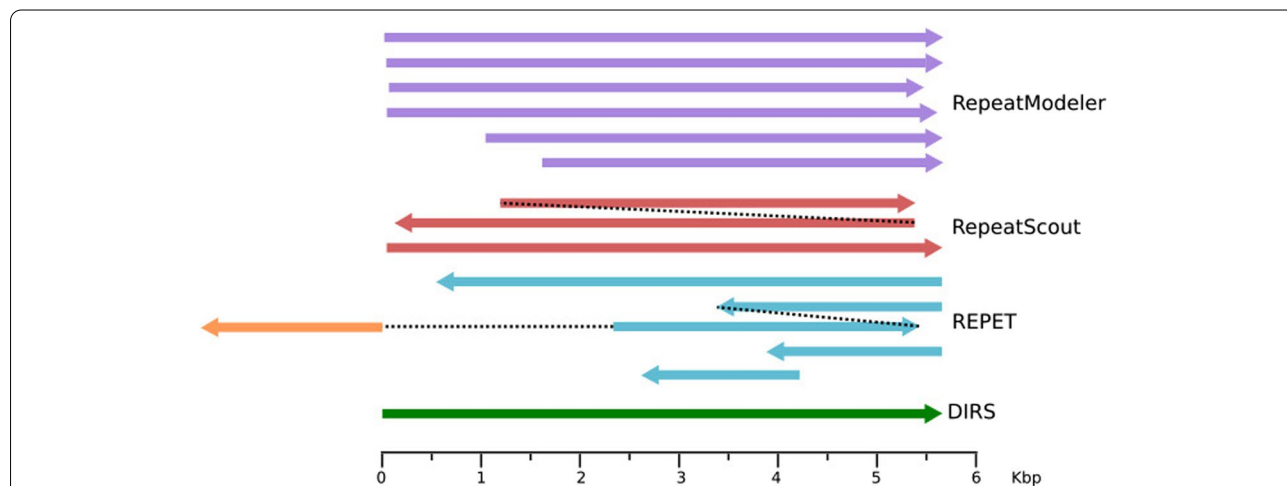


Fig. 5 Consensus sequences generated by each software for a single TE from the DIRS family. The length and position reflect the mapping to the TE, black dotted lines show the continuation of the same model, the orange segment represents a fragment coming from another TE (TRANSIB)

number of models generated for human data might be linked to the smaller sequence data compared to the two other datasets. However, based on the TE-annotation, the total length of TEs in the fruit fly genome is comparable to the total length of TEs in human chromosome 21, 24 MB versus 20 Mb, respectively. In general, the real data seem to harbor more versatile repertoire of TEs than our simulated data resulting in many more TE models (compare Table 2 and Table 3).

To investigate this matter further, we compared TE models created by the software analyzed with TE families annotated in the sequences used for the benchmarking. We simply run RepeatMasker with the created TE libraries as a query. Interestingly, in all cases there were many consensus sequences (TE models) that didn't produce any significant hits in the RepeatMasker analysis. This may suggest that there might still be undiscovered repeats in the analyzed data, although some of them could be false positive results (see Table 4). Interestingly, recent de novo analysis of the repetitive elements of the T2T-CHM13 assembly [20] revealed 49 novel repeat types [13] showing that even well annotated genomes may still harbor undefined repeats.

Individual repeats annotation

To get a better idea of the different results of the de novo annotation obtained by the six tools used, we plotted the coordinates of each one in tracks along with the reference annotation, as shown in Fig. 6. Simply by visualizing the results it is quite evident that there is a tendency to get a fragmented annotation when using k-mer counting tools, particularly P-Clouds and phRAIDER. Red uses a smoothing function to merge nearby high frequency k-mers, giving less fragmented results, as shown in Figs. 6 and 7. As it was expected due to the methodology used, the best results for detecting transposons were obtained by software that calculate TE-models, but as it is shown more in detail in Fig. 7, most of the predictions are fragmented TEs, annotations without clear borders, or missing some smaller or incomplete elements.

For model-building software, in most cases RepeatModeler got the best results, but RepeatScout obtained also comparable results. REPET failed in some scenarios

and particularly with short divergent fragments (Alu elements). However, it must be noticed that it failed only in one of the twenty cases and performed quite well on a real data. One of the reasons could be that this software was developed with the idea to be used with large genomes and here all the tests were run with sequences of around 100 MB. We compared the annotation results obtained using simulated sequences with known identities between TEs ranging from 60 to 100% and then compared the coverage of the annotation in relation to it, as is shown in Fig. 6. It is expected that TEs with higher identities are detected more precisely. Indeed TEs with a higher identity are better detected by all software and the differences seen are inherent to the performance of each tool. For k-mer counting software, Red performs significantly better than the rest, e.g. at 70% identity, Red detects approximately 85% of the TE regions. Meanwhile, P-Clouds detects about 60% and phRAIDER 25% (Fig. 8). The model-building software display a much better and uniform performance and are less affected by more divergent TEs (Fig. 8). When we also consider the different TE orders annotated and the proportion of coverage for the TEs of each order, we observe that there's no significant difference between them and all the software have a consistent performance (supplementary material, Fig. S1).

We also analyzed how well computed libraries can detect individual TEs and masked genomic sequences. To do that, we run RepeatMasker with genomic sequences as queries and built by different software libraries as references. In most cases RepeatScout libraries gave the results closest to the current expert annotation of cognate genomic sequence (see Table 5). Interestingly, none of the libraries seem to work well with the human chromosome 21 data masking only between 50 and 79% of the originally masked sequences. However, if we look at the number of annotated TEs, the situation is not looking as bad (see Table 6). This is probably due the fact that all the software used here produced libraries with rather fragmented models as compared to curated data (see Discussion above).

Finally we evaluated the performance of each tool against the datasets using the Matthews Correlation Coefficient (Fig. 9). Among k-mer counting software run

Table 4 General annotation of models built by analyzed software using RepeatMasker and cognate libraries. In the "Novel" column we list the number of models that didn't produce any individual matches during RepeatMasker run

Software	Zebrafish chrom. 1			Fruit fly genome			Human chrom. 21		
	TEs	SSR/Sat/rRNA	Novel	TEs	SSR/Sat/rRNA	Novel	TEs	SSR/Sat/rRNA	Novel
RepeatScout	2343	152	424	1558	89	946	397	29	38
RepeatModeler	1670	30	79	582	30	74	411	14	3
REPET	331	5	6	510	42	5	31	32	2

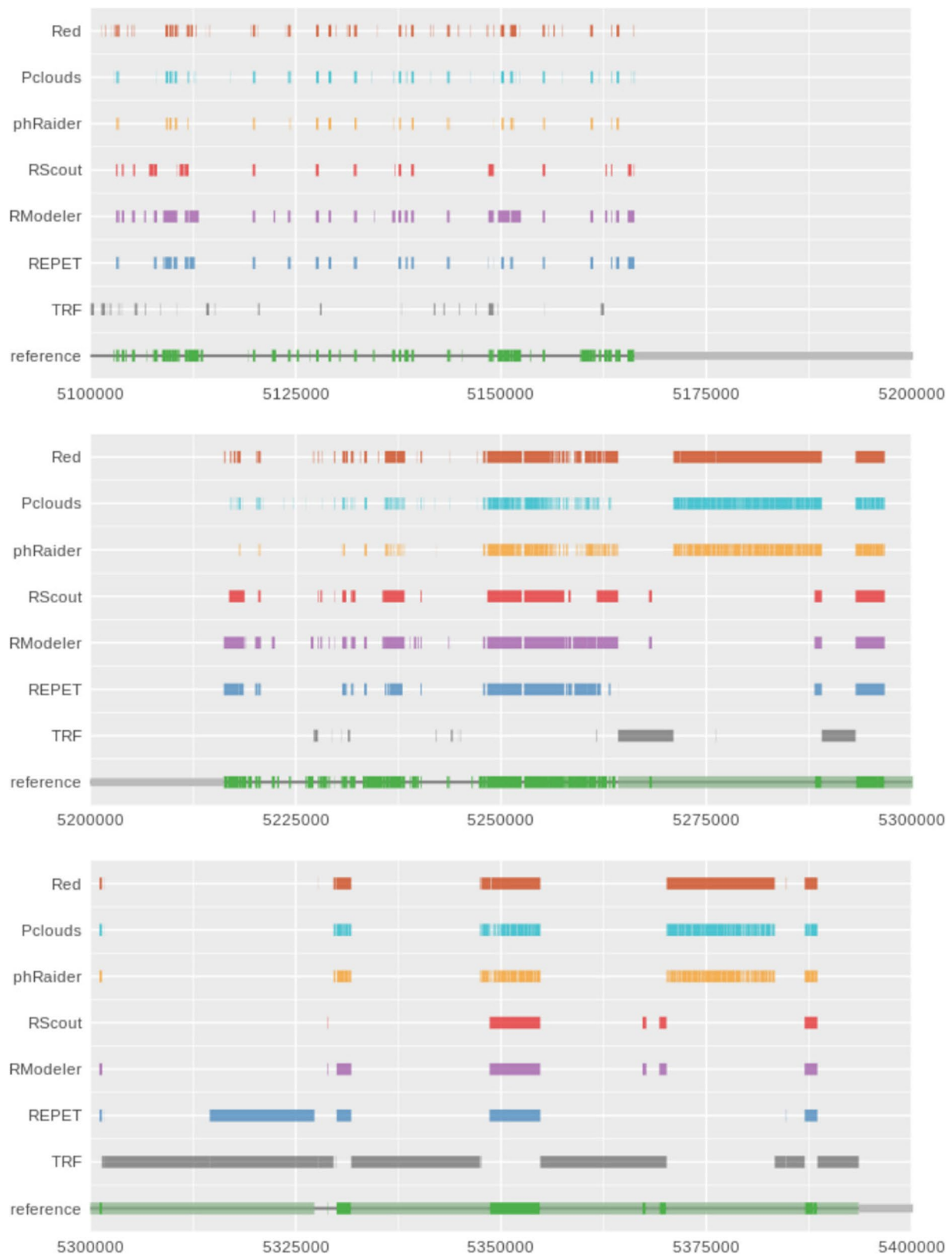


Fig. 6 Different tracks of the coordinates obtained from a de novo identification of transposons using six different software tools for detecting interspersed repeats. In the reference track, green blocks are transposons

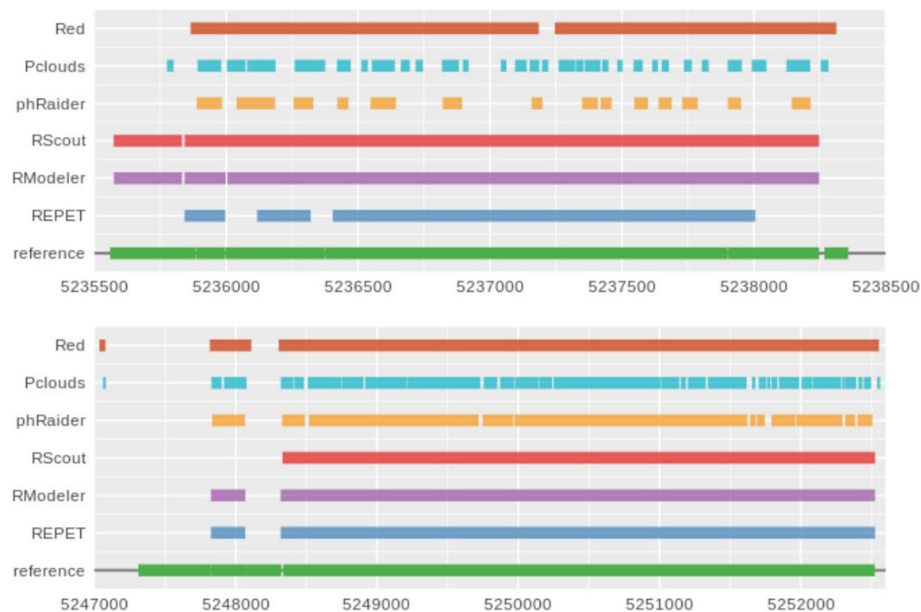


Fig. 7 Comparison of the fragmentation of results in a region of the human chromosome 21. In each track there are the predictions obtained from each tool and in green the reference. Notice how most of the results are usually incomplete or fragmented

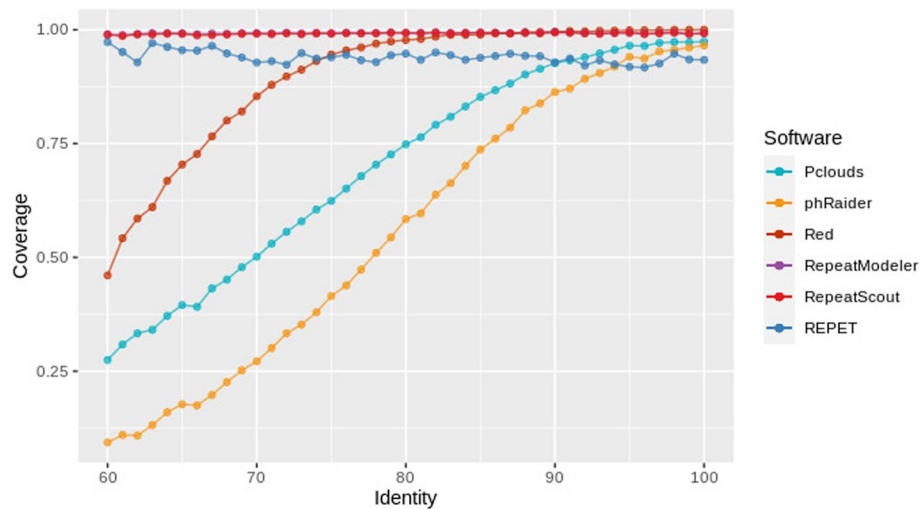


Fig. 8 Coverage of TEs in simulated sequence in relation to their average identity. In k-mer counting software there's a great drop in the detection of more divergent TEs, this behavior is not seen when using model-builders. RepeatModeler and RepeatScout results are almost identical and they are overlapped in this plot

with defaults parameters Red outperformed P-Clouds and phRAIDER and this can be explained by the fact that Red merges nearby k-mers more frequently than the others, giving less fragmented results. In model-building group RepeatModeler obtained the best results, although it is interesting to notice that RepeatScout that is part of RepeatModeler pipeline is faster than the latter but

obtained almost as good results as the pipeline. However, it should be pointed out that RepeatModeler tries to classify calculated models into repeat families, something that RepeatScout does not. Hence, execution time difference is rather expected than surprising. REPET on the other hand was probably not the best tool for this type of analysis. REPET has a default configuration but also

Table 5 Number of nucleotides detected as TEs by software analyzed. Numbers in parentheses represent fraction of the sequence covered by those TEs

Reference library	Zebrafish chromosome 1	Fruit fly genome	Human chromosome 21
RepeatScout	26,447,536 (44.39%)	23,540,405 (17.11%)	12,266,284 (26.26%)
RepeatModeler	16,561,994 (27.80%)	15,919,335 (11.57%)	13,551,826 (29.01%)
REPET	18,712,280 (31.41%)	22,412,191 (16.29%)	8,560,157 (18.33%)
Original annotation	26,465,290 (44.42%)	23,990,534 (17.44%)	17,188,977 (36.80%)
Size of analyzed sequence	59,578,282	137,547,960	46,709,983

Table 6 Number of individual TEs annotated by RepeatMasker using different reference libraries

Reference library	Zebrafish chromosome 1	Fruit fly genome	Human chromosome 21
RepeatScout	143,463	55,799	46,628
RepeatModeler	86,813	30,754	49,621
REPET	71,891	38,154	30,838
Original annotation	107,997	38,282	53,540

different tools can be added to the pipeline and each step is highly configurable, nevertheless, one of the downsides is that it can be very complex to configure and run for an unexperienced user.

Discussion

Presented here benchmark is not very optimistic but how does it measure against other similar studies? Unfortunately, the comparison is not very easy. First of all, there were not many independent benchmarking tests perform

in past [22, 26]. Although, usually some benchmarking were presented with the original publication of a software, they cannot be completely trusted as they might be tuned to a specific software. Moreover, each study includes different set of the software and we didn't find any single paper that discussed the same six programs that we benchmarked here. Ou et al. employed five programs of which three (RepeatScout, RepeatModeler, and Red) were benchmarked by us as well. However, in this paper benchmarking was used to compare these software with EDTA, a pipeline created by the authors and as such cannot be treated as an independent analysis. We found only one, truly independent benchmarking study but it was published over a decade ago and only RepeatScout is a mutual software with our study [26]. Although RepeatModeler is most frequently compared software, surprisingly it is missing from Ou et al. analysis. Another difficulty is that each study used different data sets to evaluate de novo TE-detection software and surprisingly simulated data usually was not included. Similarly to our approach the datasets varied in size but none of the tools was tested on a sequence of gigabases scale. Nevertheless,

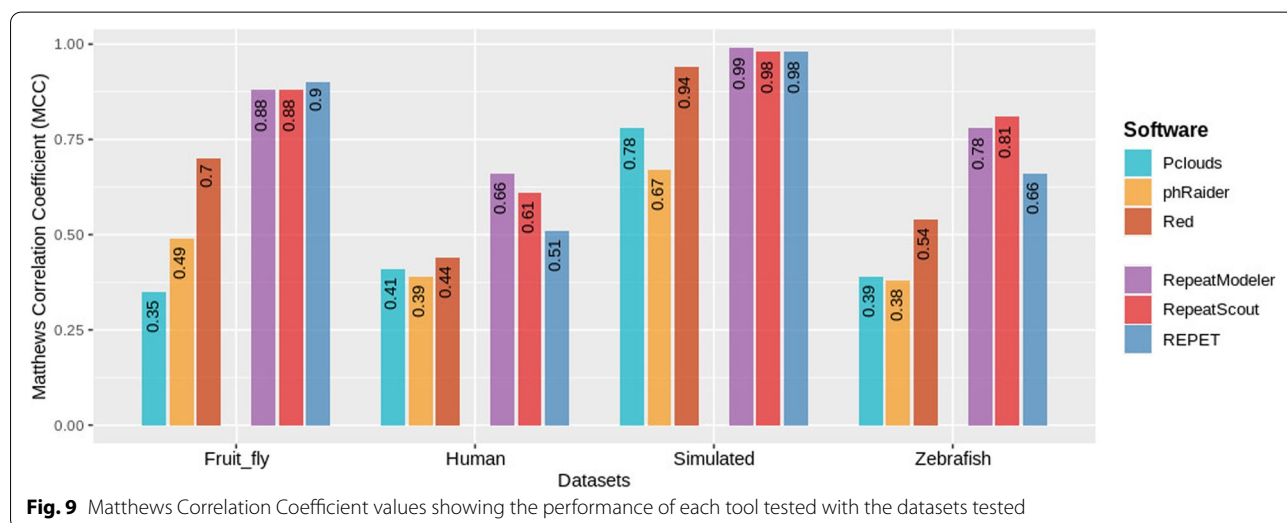


Fig. 9 Matthews Correlation Coefficient values showing the performance of each tool tested with the datasets tested

below we try to compare our findings with those previously published.

First of all, it is clear that in most studies RepeatModeler gives the best results. In the original paper Flynn et al. divided computed models into perfect, good, present, and not found [7]. Although this classification might be a bit misleading, e.g. “perfect” models might only 95% identical with the reference consensus, we looked at our results the same way (see Fig. 10). Interestingly, although in both studies *D. melanogaster* genome was subjected to the analysis, the results were not the same. Most likely, it is due to the fact that we used default parameters of RepeatModeler.

As mentioned above Saha et al. benchmarked six programs but only RepeatScout was mutual software between their and our study [26]. Moreover, they run tests solely on different rice genome data varying from 3 to 27 Mb with published annotation assumed as a base line (golden standard) and only sensitivity was reported. For unknown reason, different programs were run on different datasets. Nevertheless, RepeatScout performed the worst with 3 Mb sequence data with only 26.2% sensitivity but performed much better when run on a larger 27.8 Mb dataset with sensitivity increased to 84.3%. The other software used on both datasets (RepeatFinder) behaved similarly, 32.7 and 85.3% sensitivity, respectively. These results are similar to our observations that de novo repeat finding programs perform better when a larger dataset is used to build a TE library (see Table 5).

To evaluate the results it is important to consider not only the raw performance of each tool but also the difficulty to run, configurability, and speed. K-mer counting

software usually only accept a few parameters such as k-mer length, minimum frequency, and length; but these tools normally are very easy to run and require little computing power while being incredibly fast. However, one of the performance downsides can be the requirement to store large data structures in memory. In comparison, model-building software employ a strategy that requires much more computational resources. They are also more time consuming and can be more complex to install, configure, and run.

Based on our simulated data analysis it is clear that none of the analyzed software is able to compute a repeat consensus sequence perfectly. While the sequence of a repeat can be recovered with confidence, the structure of the repeat should be inspected manually and edited accordingly. This is especially important for longer transposons. In general, for a fast assessment of interspersed repeats, Red can be useful, acknowledging of course its limitations when it comes to low complexity sequences. For more in depth studies with small genomes, RepeatModeler seems to be the best option. It is also interesting to note that RepeatScout has a really good performance if we take into consideration speed and computational requirements. However, the situation may not be as bad as it sounds. For some tasks such as repeat masking this might be satisfactory. Nevertheless, to fully understand biology and evolution of transposable elements in a given genome automatic approach is not sufficient but it should be a good starting point.

The question appears if our simulated data are close enough to reality and if our results are comparable with other studies that used simulated sequences. It is very

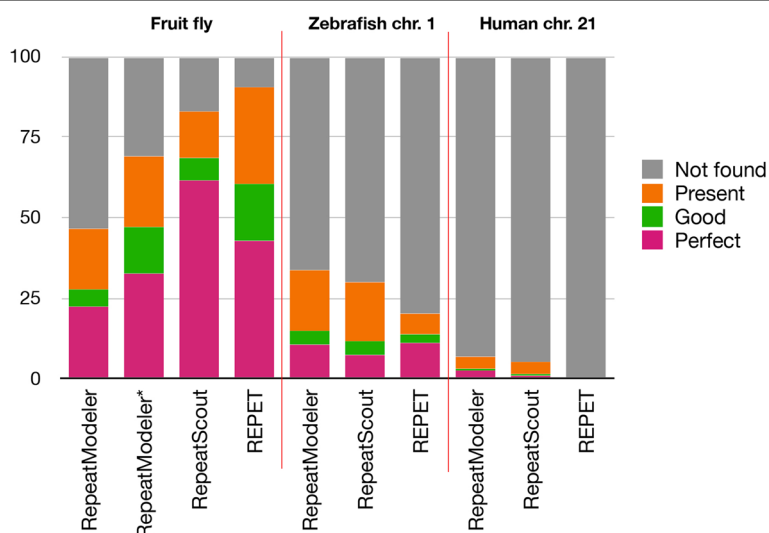


Fig. 10 Summary of models' accuracy obtained by different software and data sets. RepeatModeler* data taken from [7]. Accuracy categories as defined in Fig. 3

difficult to make such a comparison because not many studies were performed with simulated data. Flynn et al. [7], for instance, used simulated genome depleted of TEs to evaluate false positive discovery rate. In other studies, simulated sequencing reads were used to evaluate software for either polymorphic TE insertions [32] or TE-expression analyses [30]. Our simulations were different as we tried to evaluate software ability to detect TE families and their classification. The closest approach to our data simulation was attempted by Saha et al. [26] and Schaeffer et al. [27]. In both studies genomic sequences are semi-simulated in a way that TEs are left intact but sequences between the elements are replaced by a sequence generated by fifth-order Markov chain. Unfortunately, Schaeffer et al. provide statistics only in a relative way, i.e. phRAIDER vs. RepeatScout but we were able to compare our results with the findings of Saha et al. They simulated rice genome of two sizes: 3 Mb sequence fragment based on Chromosome 12 and full size simulated Chromosome 12. They reported only sensitivity of software tested, which was 26.2% in the case of a shorter sequence and 85.3% in the case of the whole Chromosome 12. In our study, RepeatScout sensitivity was 98.7%, a bit higher than reported by Saha et al., however we should note that our simulated data was larger (100 Mb) and apparently larger initial dataset guarantees better software performance.

Conclusions

We tested a number of tools for de novo detection of TEs. The results were compared using the MCC against a reference of annotated TEs. As expected, model-builders performed better than k-mer counting software, with RepeatModeler beating competitors in most datasets. However, even for RepeatModeler, the results are far from satisfactory based on the reference annotation. There is a tendency for most tools to identify TE-regions in a fragmented manner and it is also frequent that small TEs or fragmented TEs are not detected. We recognize that some of the results obtained may be improved by fine tuning of parameters; some tools like REPET are fully customizable and more tools can be added to the pipeline, although this can be challenging for most users. In conclusion, the contemporary tools for de novo detection of TEs benchmarked here are far from being perfect and the identification of TEs is still a challenging endeavor as it requires a significant manual curation by an experienced expert. It seems that for de novo detection of TEs extensive manual curation and using multiple tools for confirmation of the results obtained is necessary. We also found that MCC can be used as a fast and reliable test to compare the performance of these software and can give a general idea of which tool is best suited for each task.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13100-022-00266-2>.

Additional file 1.

Acknowledgements

Not applicable.

Authors' contributions

WM initiated and designed the project. MR executed the project and wrote a draft manuscript. Both authors worked on the final manuscript. The author(s) read and approved the final manuscript.

Funding

Open Access funding enabled and organized by Projekt DEAL. This work was partially funded by the DAAD Research Grants - Doctoral Programmes in Germany, 2018/19 (57381412) to MR and internal funds of the Institute of Bioinformatics.

Availability of data and materials

Scripts and simulated sequences used for the software evaluation are deposited at <https://github.com/IOB-Muenster/denovoTE-eval>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 18 December 2021 Accepted: 16 March 2022

Published online: 27 April 2022

References

- Bao WD, Kojima KK, Kohany O. Repbase update, a database of repetitive elements in eukaryotic genomes. *Mob DNA*. 2015;6:11.
- Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*. 1999;27(2):573–80.
- Biemont C. A brief history of the status of transposable elements: from junk DNA to major players in evolution. *Genetics*. 2010;186(4):1085–93.
- Boughorbel S, Jarray F, El-Anbari M. Optimal classifier for imbalanced data using Matthews correlation coefficient metric. *PLoS One*. 2017;12(6):e0177678.
- de Koning APJ, Gu WJ, Castoe TA, Batzer MA, Pollock DD. Repetitive elements may comprise over two-thirds of the human genome. *PLoS Genet*. 2011;7(12):e1002384.
- Flutre T, Duprat E, Feuillet C, Quesneville H. Considering transposable element diversification in De novo annotation approaches. *PLoS One*. 2011;6(1):e16526.
- Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, et al. RepeatModeler2 for automated genomic discovery of transposable element families. *Proc Natl Acad Sci U S A*. 2020;117(17):9451–7.
- Gao CH, Xiao ML, Ren XD, Hayward A, Yin JM, et al. Characterization and functional annotation of nested transposable elements in eukaryotic genomes. *Genomics*. 2012;100(4):222–30.
- Girgis HZ. Red: an intelligent, rapid, accurate tool for detecting repeats de-novo on the genomic scale. *Bmc Bioinformatics*. 2015;16:227.
- Gu WJ, Castoe TA, Hedges DJ, Batzer MA, Pollock DD. Identification of repeat structure in large genomes using repeat probability clouds. *Anal Biochem*. 2008;380(1):77–83.

11. Haeussler M, Zweig AS, Tyner C, Speir ML, Rosenbloom KR, et al. The UCSC genome browser database: 2019 update. *Nucleic Acids Res.* 2019;47(D1):D853–8.
12. Hoen DR, Hickey G, Bourque G, Casacuberta J, Cordaux R, et al. A call for benchmarking transposable element annotation methods. *Mob DNA.* 2015;6:13.
13. Hoyt SJ, Storer JM, Hartley GA, Grady PGS, Gershman A, et al. From telomere to telomere: the transcriptional and epigenetic state of human repeat elements. *Science.* 2022;376(6588):eabk3112.
14. Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, et al. The Dfam database of repetitive DNA families. *Nucleic Acids Res.* 2016;44(D1):D81–9.
15. Jurka J, Kapitonov VV, Kohany O, Jurka MV. Repetitive sequences in complex genomes: structure and evolution. *Annu Rev Genomics Hum Genet.* 2007;8:241–59.
16. Kapitonov VV, Jurka J. Self-synthesizing DNA transposons in eukaryotes. *Proc Natl Acad Sci U S A.* 2006;103(12):4540–5.
17. Kubiak MR, Makalowska I. Protein-coding Genes' Retrocopies and their functions. *Viruses.* 2017;9(4):80.
18. Makalowski W. Genomic scrap yard: how genomes utilize all that junk. *Gene.* 2000;259(1–2):61–7.
19. Makalowski W, Gotea V, Pande A, Makalowska I. Transposable elements: classification, identification, and their use as a tool for comparative genomics. *Methods Mol Biol.* 2019;1910:177–207.
20. Nurk S, Koren S, Rhie A, Rautiainen M, Bizikadze AV, et al. The complete sequence of a human genome. *Science.* 2022;376(6588):44–53.
21. Ohno, S., 1973 So much "junk" DNA in our genome, pp. 366–370 in *Evolution of Genetic Systems: Brookhaven Symposia in Biology*, edited by H. Smith. Gordon and Breach, New York.
22. Ou S, Su W, Liao Y, Chougule K, Agda JRA, et al. Benchmarking transposable element annotation methods for creation of a streamlined, comprehensive pipeline. *Genome Biol.* 2019;20(1):275.
23. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes. *Bioinformatics.* 2005;21:1351–8.
24. Quesneville H, Bergman CM, Andrieu O, Autard D, Nouaud D, et al. Combined evidence annotation of transposable elements in genome sequences. *PLoS Comput Biol.* 2005;1(2):166–75.
25. Ricker N, Qian H, Fulthorpe RR. The limitations of draft assemblies for understanding prokaryotic adaptation and evolution. *Genomics.* 2012;100(3):167–75.
26. Saha S, Bridges S, Magbanua ZV, Peterson DG. Empirical comparison of ab initio repeat finding programs. *Nucleic Acids Res.* 2008;36(7):2284–94.
27. Schaeffer CE, Figueroa ND, Liu XL, Karro JE. phRAIDER: pattern-hunter based rapid Ab initio detection of elementary repeats. *Bioinformatics.* 2016;32(12):209–15.
28. Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, et al. The B73 maize genome: complexity, diversity, and dynamics. *Science.* 2009;326(5956):1112–5.
29. Smit, A., R. Hubley, and P. Green, 2013–2015 RepeatMasker Open-4.0.
30. Teissandier A, Servant N, Barillot E, Bourc'his D. Tools and best practices for retrotransposon analysis using high-throughput sequencing data. *Mob DNA.* 2019;10:52.
31. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, et al. A unified classification system for eukaryotic transposable elements. *Nat Rev Genet.* 2007;8(12):973–82.
32. Yu T, Huang X, Dou S, Tang X, Luo S, et al. A benchmark and an algorithm for detecting germline transposon insertions and measuring de novo transposon insertion frequencies. *Nucleic Acids Res.* 2021;49(8):e44.

Publisher's Note

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

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

