

piClusterBuster: Software for Automated Classification and Characterization of piRNA Cluster Loci

Patrick Schreiner^{1*} and Peter W. Atkinson^{1,2}

¹Interdepartmental Graduate Program in Genetics, Genomics & Bioinformatics, University of California, Riverside, CA 92521, USA, ²Department of Entomology and Institute for Integrative Genome Biology, University of California, Riverside, CA 92521, USA

*corresponding author

Abstract

Background

Piwi-interacting RNAs (piRNAs) are sRNAs that have a distinct biogenesis and molecular function from siRNAs and miRNAs. The piRNA pathway is well-conserved and shown to play an important role in the regulatory capacity of germline cells in Metazoans. Significant subsets of piRNAs are generated from discrete genomic loci referred to as piRNA clusters. Given that the contents of piRNA clusters dictate the target specificity of primary piRNAs, and therefore the generation of secondary piRNAs, they are of great significance when considering transcriptional and post-transcriptional regulation on a genomic scale. A quantitative comparison of top piRNA cluster composition can provide further insight into piRNA cluster biogenesis and function.

Results

We have developed software for general use, piClusterBuster, which performs nested annotation of piRNA cluster contents to ensure high-quality characterization, provides a quantitative representation of piRNA cluster composition by feature, and makes available annotated and unannotated piRNA cluster sequences that can be utilized for downstream analysis. The data necessary to run piClusterBuster and the skills necessary to execute this software on any species of interest are not overly burdensome for biological researchers.

piClusterBuster has been utilized to compare the composition of top piRNA generating loci amongst 13 Metazoan species. Characterization and quantification of cluster composition allows for comparison within piRNA clusters of the same species and between piRNA clusters of different species.

Conclusions

We have developed a tool that accurately, automatically, and efficiently describes the contents of piRNA clusters in any biological system that utilizes the piRNA pathway. The results from piClusterBuster have provided an in-depth description and comparison of the architecture of top

piRNA clusters within and between 13 species, as well as a description of annotated and unannotated sequences from top piRNA cluster loci in these Metazoans.

piClusterBuster is available for download on GitHub:

<https://github.com/pschreiner/piClusterBuster>

Background

P-element induced wimpy testis (PIWI) proteins and the utilization of the PIWI-interacting RNA (piRNA) pathway has been conserved in a diverse range of Metazoans, including sponges, roundworms, fruit flies, and humans [1]. The importance of the role of piRNAs in fertility was demonstrated in Metazoans by the observation of crosses after exposure of *Drosophila melanogaster* to the *P* transposable element (TE) [2]. When female flies containing the *P* element were mated with males that were lacking it, the progeny were fertile. However, when males containing the *P* element were mated with females lacking it, hybrid dysgenesis occurred, leading to sterility of the progeny [3]. It was later discovered that exposure to the *P* element prompted maternal deposition of piRNAs to effectively silence the *P* element and allow for fertile progeny [4]. Perturbations to the piRNA pathway have also demonstrated gametogenic defects in *M. musculus* and *D. rerio* [5, 6]. Although, piRNAs are notably absent in plant and fungal species [1].

piRNAs are a subset of sRNAs between 24-31 nucleotides in length, although the range of the piRNA size distribution varies across species [7,8]. piRNAs are generated via a primary or secondary mechanism of biogenesis [7].

Primary piRNAs derive from discrete genomic loci that are referred to as piRNA clusters. These loci can vastly range in size from under one thousand nucleotides to over one hundred thousand nucleotides in length. Transcription of piRNA clusters can occur in several distinct manners depending on the nature of the piRNA clusters [7].

piRNA clusters are characterized as unidirectional, bidirectional, or dual-stranded based on the direction transcription at the locus [9,10]. The transcripts generated from piRNA clusters serve as precursor molecules for piRNAs, undergoing dicer-independent slicing and modification at their 3' end [7,11]. The processing of primary piRNAs has been shown to demonstrate a bias of U at position 1 of the piRNAs [7]. When post-transcriptional processing of piRNAs is complete, the molecules are referred to as mature piRNAs.

Mature piRNAs then associate with an Argonaute family, PIWI protein to form a RNA-induced silencing complex (RISC). The RISC complex has the capability to facilitate both transcriptional and post-transcriptional regulation. RISC-mediated transcriptional regulation occurs via piRNA association and guiding of a PIWI protein which facilitates epigenetic modification in *Drosophila* [12]. RISC-mediated post-transcriptional regulation occurs via piRNA association with PIWI or AGO3, which leads to piRNA-directed cleavage of mRNAs in *Drosophila* [7].

piRNAs have also been implicated in post-transcriptional silencing of mRNAs via poly(A) deadenylation [13,14].

The number of PIWI proteins can differ in Metazoan species. While three PIWI proteins have been identified in *D. melanogaster*, *H. sapiens*, and *M. musculus*, as few as two PIWI proteins have been identified in *D. rerio* and as many as seven PIWI proteins have been identified in *Ae. aegypti* [7,15-18]. It has not yet been determined whether the variation in the number of PIWI proteins between these species is a result of redundant, compensatory, or additional functionality.

Secondary piRNAs are generated by the slicing mechanism of RISC regulation, resulting in what is referred to as the amplification loop, or ping-pong pathway [11]. The amplification loop functions by primary piRNA targeting of mRNA via sequence complementarity, followed by PIWI-mediated slicing of the target mRNA. The remaining fragment of the mRNA can be processed into a secondary piRNA. A secondary, mature piRNA can then associate with AGO3, and slice other mRNA targets via sequence complementarity in *Drosophila*. The overlap of complementarity between the piRNAs and their mRNA targets is generally ten base pairs in the opposite orientation, leaving an A10 bias in secondary piRNAs [11]. piRNAs have been known to target TEs, genic mRNAs, viral mRNAs, and even rRNA molecules [11,12,19,20].

Given that the contents of piRNA clusters dictate target specificity, finding the origin of these sequences is of great importance in understanding the biogenesis and function of piRNA clusters. We have developed software, piClusterBuster, to be capable of automatically, consistently, and efficiently detecting top piRNA cluster loci and thoroughly describing the contents of those loci on a large scale. The capability that piClusterBuster has to quantify piRNA cluster composition and describe annotated, as well as unannotated sequences in diverse Metazoan species allows for meaningful comparisons that can aid in facilitating a better understanding of top piRNA cluster biogenesis and function across species. Exploring piRNA cluster composition on a large scale can provide insight into conserved piRNA cluster architecture that dictates piRNA cluster biogenesis and function.

Implementation

piClusterBuster is a series of integrated R and bash scripts that interact along with other standalone bioinformatics programs to perform piRNA cluster characterization and annotation. The tool supports a variety of user input data, customization of the analyses and computational resources to be used in executing the program. piClusterBuster is intended to be executed in a Unix environment and has a series of required software dependencies (Table 2).

Flag	Default	Description
Data Input		
-fa <file>		Indicates a FASTA input file containing piRNA

Flag	Default	Description
		cluster sequences of interest
-fq <file>		Indicates a FASTQ input file containing quality trimmed sRNAs
-bed <file>		Indicates a BED input file containing the location of the piRNA clusters of interest
-gid	“Genome”	Name of the piClusterBuster Run
Databases (provide in FASTA format)		
-x <file>		Reference Genome
-gndb <file>		Organism-specific Gene Set
-tedb <file>		Transposable Element (TE) Set
Additional Analysis		
-n	5	Indicates the number of piRNA clusters to be analyzed
-ncbidb <file>	None	NCBI Nucleotide Database
--verbose	FALSE	Retain intermediate results
--go	FALSE	Perform a Gene Ontology enrichment analysis on sequence of genic origin
--all-srna	FALSE	Observe all sRNA, not just piRNAs
Performance Enhancement		
--qsub	FALSE	Submit jobs via Torque/Maui Resource Allocation
--srun	FALSE	Submit jobs via Slurm Resource Allocation
-p	1	Number of processors to utilize

Table 1 | Program Parameters. A list of the options, corresponding runtime flags, and default values available for use in piClusterBuster. A flag is an indicator to specify the type of input information to the application. A blank in the “Default” column constitutes a required parameter.

Database	Website	Reference
----------	---------	-----------

Database	Website	Reference
NCBI (nt)	http://www.ncbi.nlm.nih.gov/	[21]
RepBase	http://www.girinst.org/repbase/	[22]

Standalone Programs	Website	Reference
BLAST+ (v2.2.30+)	http://blast.ncbi.nlm.nih.gov/	[23]
CENSOR	http://www.girinst.org/censor/	[24]
proTRAC	http://www.smallrnagroup.uni-mainz.de/	[25]
RepeatMasker	http://www.repeatmasker.org/	[26]

R Packages	Reference
Biostrings	[27]
doMC	[28]
GenomicRanges	[29]
gProfileR	[30]
Plyr	[31]
Qcc	[32]
Seqinr	[33]
systemPipeR	[34]

Table 2 | List of software and databases utilized in the piClusterBusterR

Depending on the user-specified analyses to be performed, piClusterBusterR may require these **(A)** standalone software and **(B)** R libraries.

Workflow

piClusterBusterR only requires four input parameters from the user on the command line: (1) input data, (2) a reference genome, (3) a species-specific gene set, and (4) a set of known TEs.

Additional options are available to increase the efficiency of the software and to customize the program output.

piClusterBuster allows for data input in the form of sRNA reads, piRNA cluster sequences, or piRNA cluster chromosomal loci. When sRNA reads are provided as the data input, piClusterBuster must perform additional steps in order to assign piRNA cluster loci. First, all of the reads are filtered in order to analyze only those that are 24 nucleotides in length or greater. The piRNAs from the filtered FASTQ file are then mapped to the user-provided reference genome using proTRAC's sRNA mapping tool. The piClusterBuster-generated map file is then utilized to define the top piRNA cluster loci using proTRAC [25].

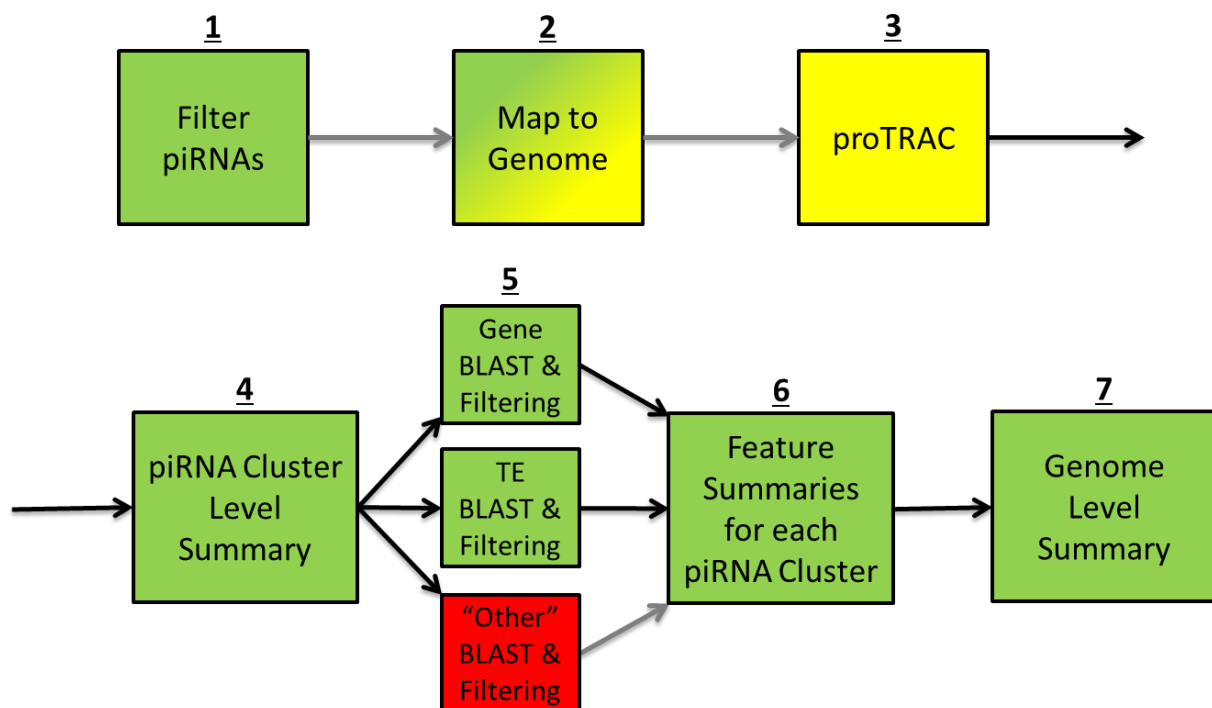


Figure 1 | Algorithm Overview. A workflow of the steps taken by piClusterBuster to annotate and characterize piRNA clusters. The step number is indicated above the boxes that describe the analysis in each step. The relative time requirement of each step is indicated by green, yellow, and red from fastest to slowest. The gray arrows indicate that the previous step may be skipped if it is unnecessary.

proTRAC is an standalone tool designed for the definition of piRNA clusters [25]. proTRAC considers features of small RNA sequence reads features that are indicative of piRNAs such as read length and a U1 or A10 bias. proTRAC uses a density-based approach to identify genomic regions that have piRNA accumulation, as defined by a significant deviation from a hypothetical uniform distribution, which then defines the degree of confidence of the piRNA cluster call. proTRAC has demonstrated efficacy in piRNA cluster definition relative to previously

established methods of piRNA cluster detection [25]. The proTRAC output is then processed to identify the top piRNA cluster loci, as defined by the number of normalized reads per piRNA cluster, and converted to a BED file of piRNA cluster loci. The BED file is then utilized to analyze the contents of the individual top piRNA cluster loci.

In the individual piRNA cluster level analysis, piClusterBuster performs a detailed characterization and quantification regarding the contents of each individual piRNA cluster. The user has the option to analyze piRNA clusters sequentially (default), or in parallel for each piRNA cluster of interest.

piClusterBuster first extracts the sequence using the chromosomal coordinates and reference genome that was provided by the user. piClusterBuster then attempts to identify the origin of the sequences within the piRNA clusters of interest.

In order to best infer the origin of the sequences within a given piRNA cluster, piClusterBuster utilizes what we refer to as nested annotation using RepeatMasker, CENSOR, and BLAST [23-24,26]. Nested annotation allows for sequential and non-redundant definition of known sequences with the piRNA cluster sequences under observation. RepeatMasker is run initially on the piRNA cluster of interest using the TE database and organism-specific gene set provided by the user [26]. TE and organism-specific data sets were extracted from RepBase and NCBI non-redundant nucleotide databases, respectively [21-22]. Any of the unannotated sequence remaining in the piRNA cluster of interest is extracted and subjected to TE and genic analysis via CENSOR [24]. Finally, the remainder of the unannotated piRNA cluster sequence is subjected to a blastn search, with a word size of 7 and maximum E-value of 1e-3, against the NCBI nucleotide database [21,23]. Any of the hits returned in the BLAST of sequences within the NCBI non-redundant (nt) database are classified as “Other,” in comparison to sequence originating from known TEs or genes (Figure 2). Regions of the piRNA cluster loci that have not been defined with a known sequence origin are then extracted and reported. In doing so, piRNA cluster sequence of unknown origin can be easily accessed for downstream analysis.

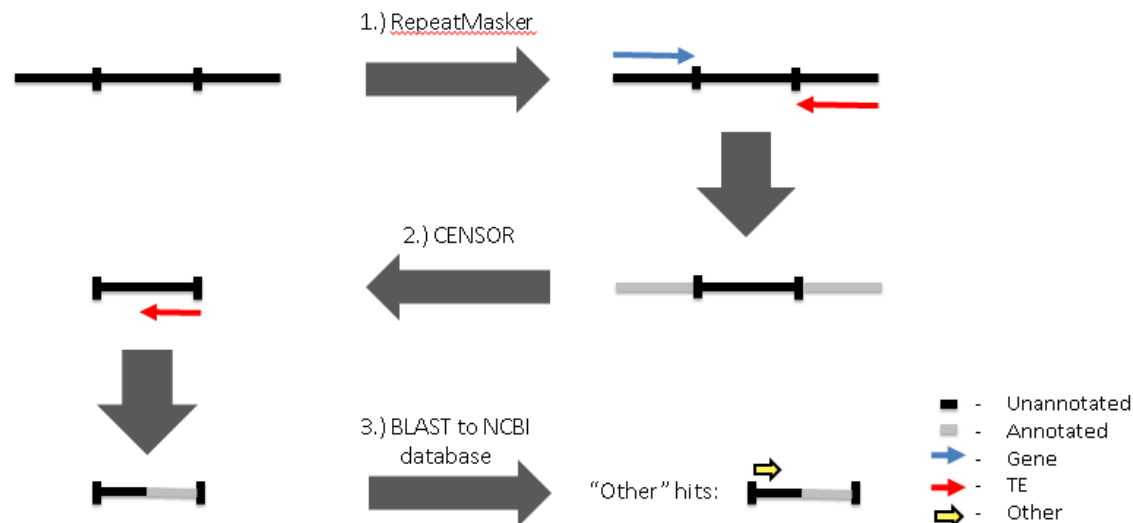


Figure 2 | Nested annotation. Workflow regarding the characterization of unannotated loci are using RepeatMasker, CENSOR, and BLAST [23-24,26]. Sequence characterization in the former steps excludes sequences from being passed to the latter.

To ensure that there is no redundancy in the sequence characterization, feature filtering is performed to only retain the best available annotation for a given piRNA cluster sequence. The best available annotation is defined as the hit with the longest available alignment length and highest similarity percentage to a known feature.

Non-redundant TE, genic, and “other” annotation is then summarized and plotted. A directory containing all of the intermediate annotation files and summary files is output in an individual directory to represent the analysis of each individual piRNA cluster. The results of the piRNA cluster level analyses of each piRNA cluster are stored so that they can be used in the genome level analysis.

In the genome level analysis, annotation is graphically compared between individual piRNA clusters. The piRNA clusters are compared in terms of their length, contents, degree of strand specificity, and percent genome occupancy. Top piRNA cluster loci can then be compared between piRNA clusters within the same species and between species on a genomic level (Figure 1).

Output

The main directory represents the outcome of the genome-level analysis. Four output files are generated in the genome-level analysis: (1) a BED file containing the piRNA cluster coordinates, (2) an aggregate file describing the total occupancy of piRNA clusters relative to the size of the organism’s genome, as well as the final data necessary to make the genome summary plots in a (3) graphical and (4) text format [36]. The genome-level graphical output contains a comparison of piRNA cluster size, piRNAs associated with each piRNA cluster, feature composition, and

strandedness of feature calls, followed by the average feature content composition across all piRNA cluster loci analyzed (Figure 3).

The genome level analysis also provides an individual directory for each piRNA cluster of interest in the order specified within the BED file of piRNA cluster loci. Within each piRNA cluster directory resides intermediate data files and summary files that are necessary to produce the piRNA cluster-level graphical output. The intermediate data files that were used in the data collection are available in the respective program output format defaults for each utilized tool (Table 1). The unfiltered BLAST output for each piRNA cluster, however, can often be large in size and is therefore removed by default. The piRNA cluster-level summary is also available in a text and graphical output.

The piRNA cluster-level graphical output contains a representation of the number of each feature that was characterized within the piRNA cluster, the nucleotide occupancy of each feature called, the nucleotide occupancy of all feature calls in both orientations, a representation of the prominent TE superfamilies within the piRNA cluster, the prominent specific TEs called within the piRNA cluster, and optionally, the most significant GO terms associated with genic hits within the piRNA cluster, a GO enrichment analysis of genic hits within the piRNA cluster, and stranded sRNA coverage plot with annotated features (Figure 4).

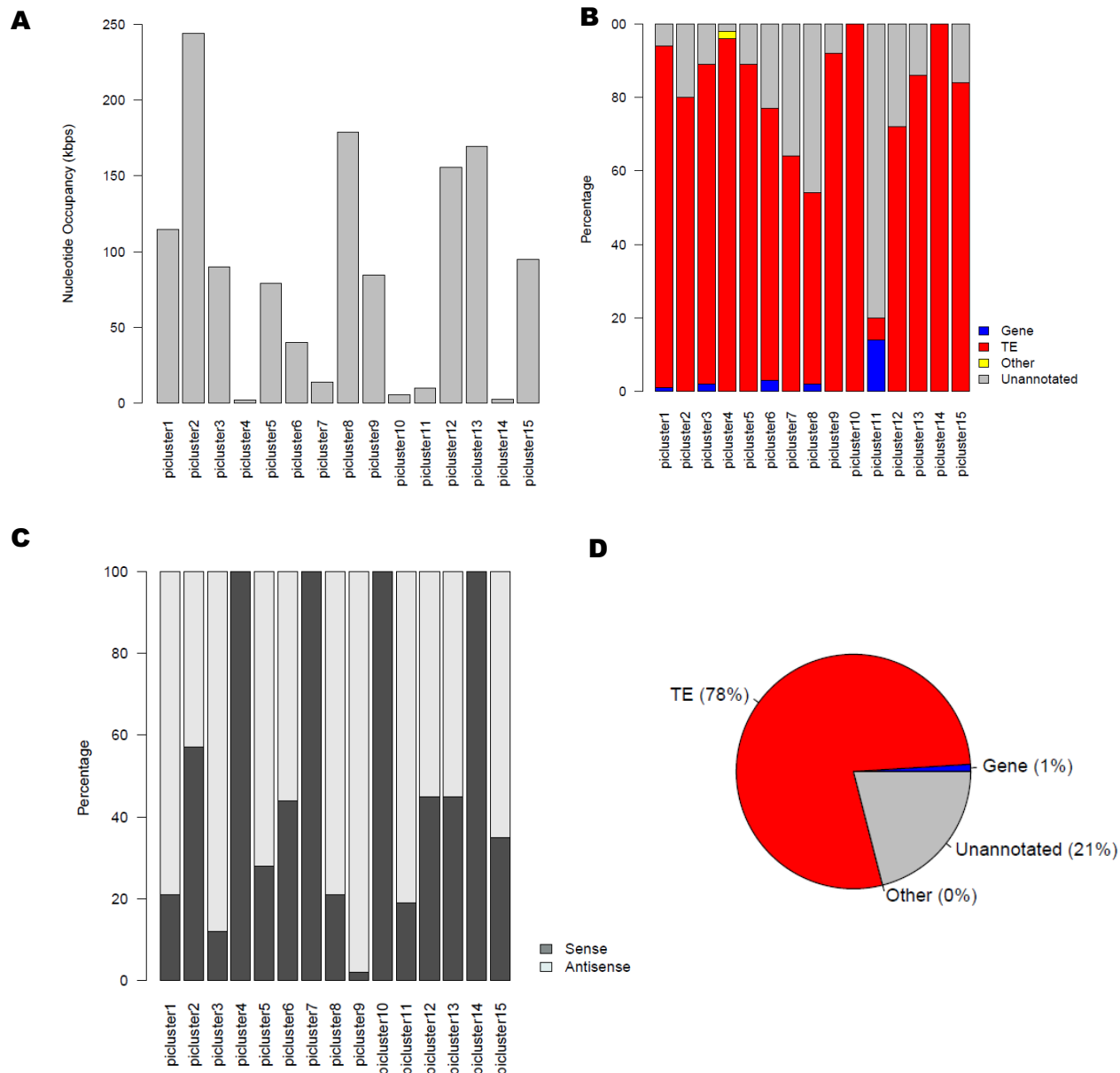
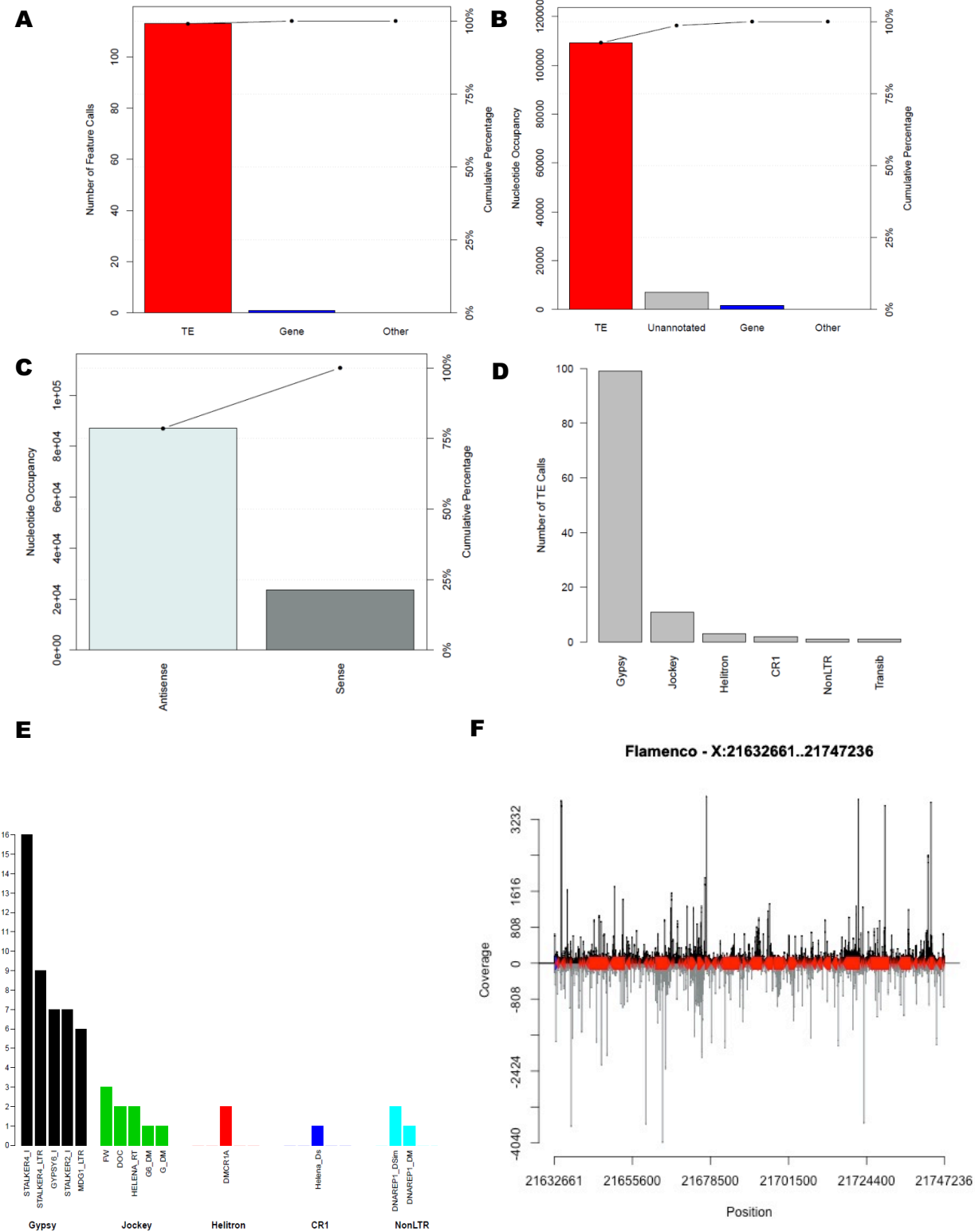


Figure 3 | Genome-Level Analysis of Top 15 piRNA Cluster Contents in *Drosophila melanogaster*. (A) Comparison of piRNA clusters by the total number of nucleotides occupied (B) Relative nucleotide occupancy occupied by each feature (C) Stranded nucleotide occupancy of feature calls. Unannotated sequences are not considered in this representation (D) Average nucleotide occupancy occupied by each feature across the top 15 *D. melanogaster* piRNA clusters identified by Brennecke *et al.* 2007 [9]. The *flamenco* and 42AB loci are represented as piRNA clusters 1 and 2, respectively [9,12].



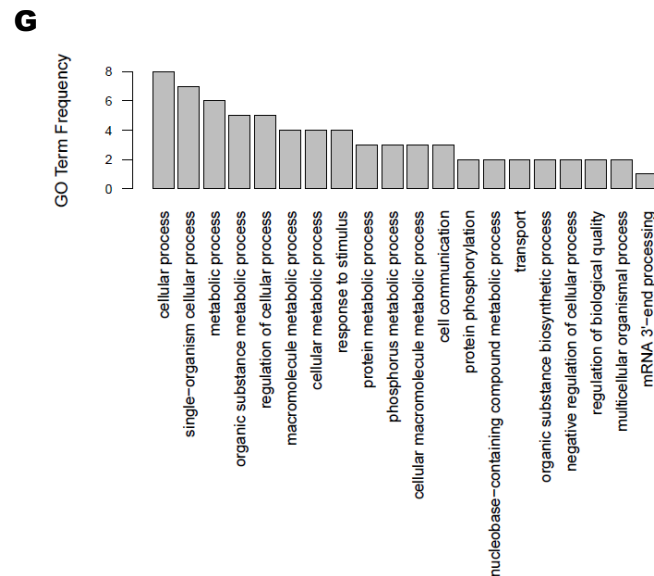


Figure 4 | piRNA Cluster-Level Analysis of the *Flamenco* Locus of *Drosophila melanogaster*. (A) Number of known TE, Gene, or “Other” feature calls (B) Nucleotide occupancy of the feature calls (C) Orientation of feature calls (D) Number TE calls within the piRNA cluster for the most represented TE superfamilies (E) Number of individual TEs called in the top 5 represented superfamilies. Additional functionality can optionally be specified by the user to prompt production of a (F) sRNA coverage plot with feature content and orientation in 0-2hr eggs libraries (G) GO term frequency plot regarding all gene hits within the piRNA cluster.

Results

Benchmarking Software Performance

piClusterBusterR was timed for the analysis of the top 5 piRNA clusters identified in *Drosophila melanogaster* ovarian samples. When running sequentially on a single Intel(R) Xeon(R) CPU E5-2683 v4 at 2.10GHz, piClusterBusterR took approximately 3 hours to complete.

When utilizing the multithreading and multitasking capability of piClusterBusterR to analyze the same 5 piRNA clusters, using 5 nodes and 6 cores per node using the same processor speed, the timing of the piClusterBusterR run took approximately 20 minutes to run. One compute node was designated per piRNA cluster and six threads were utilized on each node. This run represents the enhanced capability of piClusterBusterR if additional resources, such as a computing cluster and queue submission system, are available to the user. Output from these independent analyses was identical.

piClusterBusterR results were observed on the previous established contents of the *flamenco* locus in *Drosophila melanogaster* which were extracted from the FlyBase database [9,36].

Previous exploration with regard to the contents of the *flamenco* locus used RepeatMasker to characterize sequence content [10]. The results of this method were extracted from the UCSC

Table Browser retrieval tool [37]. The contents and strand specificity of feature calls within the *flamenco* locus identified by piClusterBuster were consistent with the previous observation (Figure 5).

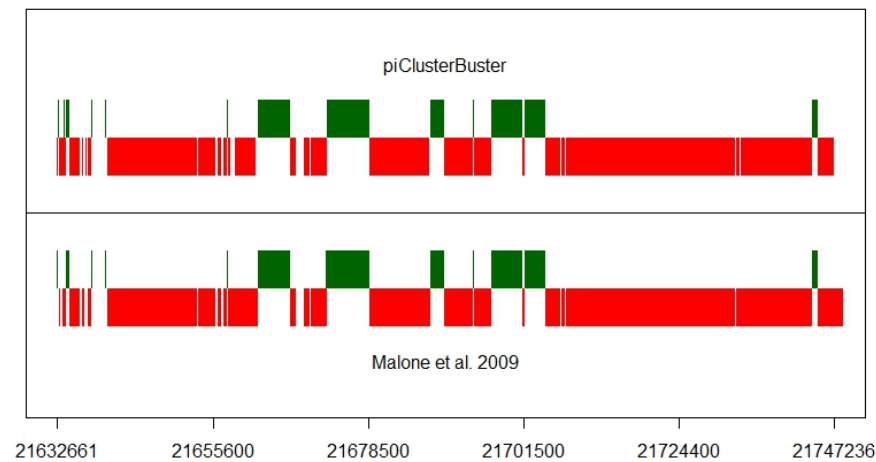
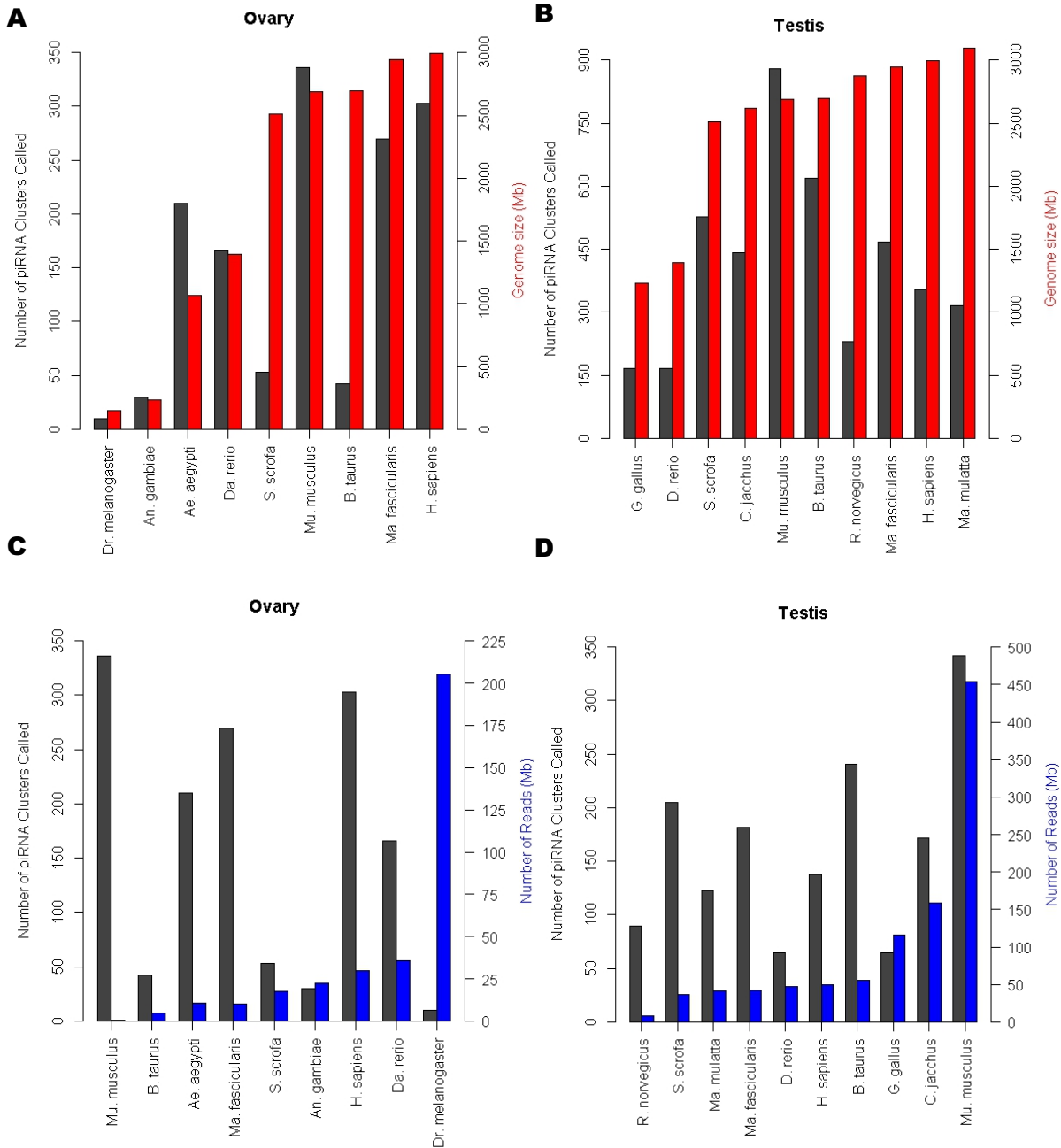


Figure 5 | Comparison of *Flamenco* TE Annotation. A depiction of the agreement of piClusterBuster characterization of the *flamenco* locus in comparison to previous reported characterization of TE contents in this locus in *Drosophila melanogaster* [10] (Figure S1). Green boxes represent a sense orientation of TE calls and the red boxes represent an antisense orientation.

piRNA Cluster Definition is Unaffected by Genome Size and Read Coverage

The 13 Metazoan species analyzed were selected based on data availability. A density-based approach of piRNA definition was implemented via use of the previous established software, proTRAC [25]. A Pearson correlation test demonstrated that piRNA cluster definition by proTRAC appears to be irrespective of genome size of the organism and the number of piRNAs available for analysis at the 1% confidence level (Figure 6A-D, Figure S6).



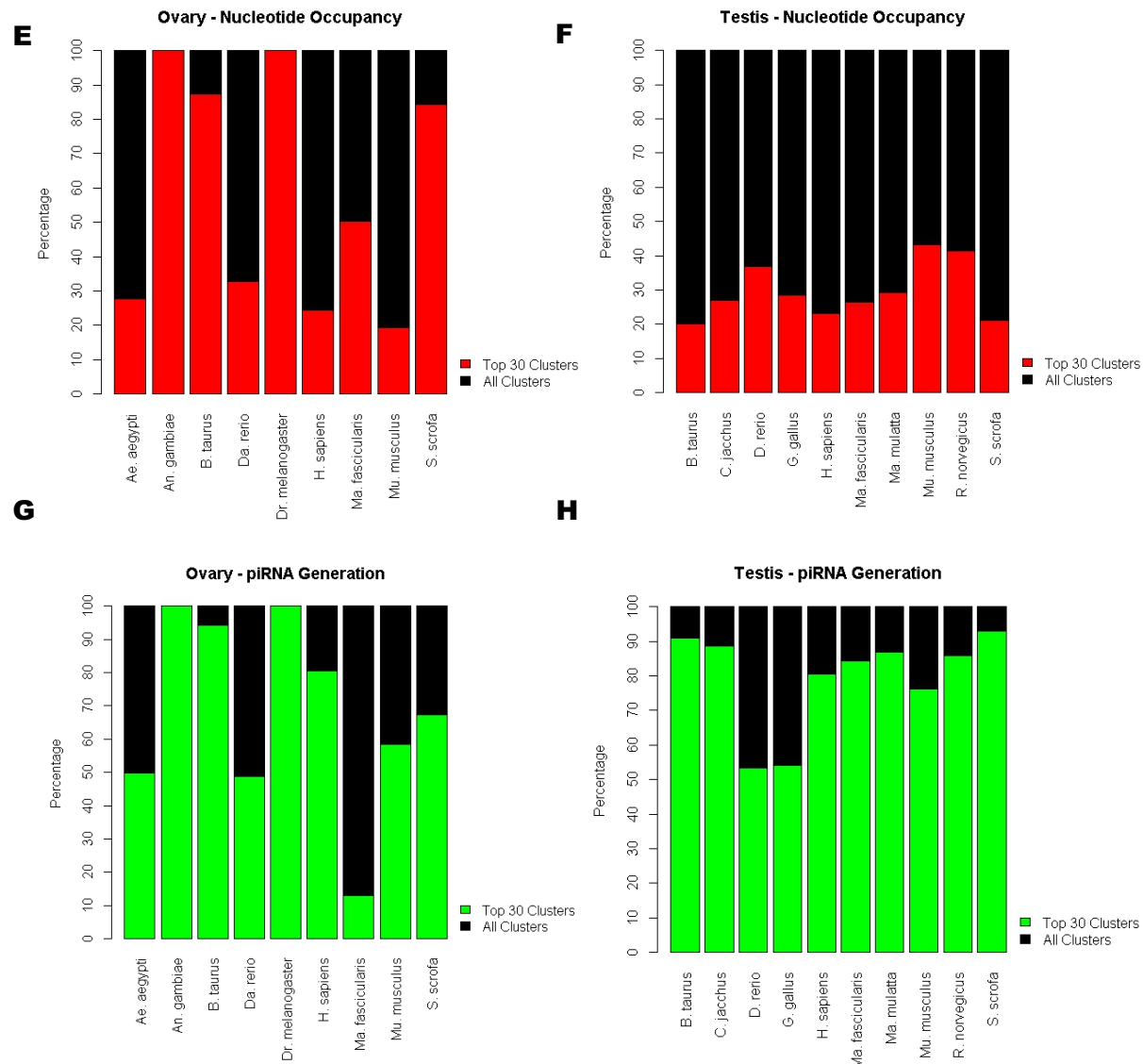


Figure 6 | Representation of piRNA Cluster Loci. Number of piRNA cluster calls relative to the genome size of the organism in (A) ovary and (B) testis. Number of piRNA cluster calls relative to the number of reads available in (C) ovary and (D) testis. The percent composition of the top 30 piRNA clusters (red/blue) relative to the full repertoire of piRNA clusters (black) called by proTRAC with regard to the percent of nucleotides occupied by piRNA clusters in (E) ovary and (F) testis. The percentage of piRNA generated from the top piRNA cluster loci in (G) ovary and (H) testis.

Given that the breadth of piRNA clusters is difficult to define in a given organism, due to the concern of false positives, we have used only the top 30 major contributing piRNA cluster loci in the between species comparisons of piRNA cluster composition. piRNA cluster definition required a length of at least five kilobases, at least 75% of the piRNAs deriving from a putative piRNA cluster with a U-1 or A-10, at least 50% of the piRNAs deriving from a putative piRNA

cluster with a U-1 and A-10, and the top 1% of piRNA sequences cannot comprise more than 90% of the piRNAs that were used to define a particular piRNA cluster.

In ovarian samples, the nucleotide occupancy of the top 30 piRNA clusters ranged from 19.3% to all of the piRNA clusters defined in a tissue with an average of 58.5% and median of 50.2% in these species (Figure 6E). The percent piRNA generation of the top 30 piRNA cluster loci ranged from 13.0% to all of the piRNAs generated in a tissue with an average of 68.0% and median of 67.3% relative to total piRNA generation (Figure 6G).

In testes samples, the nucleotide occupancy ranged from 20.1% to 43.1% relative to all of the piRNA clusters defined with an average of 29.7% and median of 27.8% in these species (Figure 6F). The percent piRNA generation of the top 30 piRNA cluster loci ranged from 53.3% to 93% with an average of 79.3% and median of 85.0% relative to total piRNA generation (Figure 6H).

Therefore, we consider the top 30 piRNA clusters to be representative of large scale architecture of genomic piRNA clusters based on the large proportion of the nucleotide occupancy and piRNA generation that is correlated with these loci.

Top piRNA Cluster Architecture is Conserved in Metazoans on a Large Scale

The analysis of piRNA cluster architecture focused on the number of piRNA clusters, piRNA cluster size, the known features within the piRNA cluster, and the orientation of the known feature.

Certain features of piRNA cluster architecture were conserved better than others. In all of the Metazoan species observed in this analysis, the majority of piRNA cluster sequence was unable to be attributed to any known origin. Unannotated sequence ranged between 18% and 70% of piRNA cluster composition. TEs were the major known contributor to piRNA cluster loci. TEs occupied up to 78% of ovarian piRNA cluster loci and 62% of testis piRNA clusters, with an average piRNA cluster occupancy of 40% to 32% in ovarian and testes libraries, respectively. Sequences of known genic origin ranged from 1 to 11%, with an average of 3% and 3.5% piRNA cluster occupancy in ovaries and testes, respectively. Non-genic, non-TE sequences within the NCBI database were the least significant contributor to piRNA cluster loci in these species, ranging from 0 to 9% of piRNA cluster composition, with an average piRNA cluster occupancy between 4 to 3.5% in ovarian and testes libraries (Figure 7).

The strand specificity of feature calls within top piRNA cluster loci was also summarized. Features were predominantly characterized on the sense strand of piRNA clusters. The nucleotide occupancy of sense features accounted for between 38.0% to 61.0% of feature calls with a piRNA cluster with an average of 50.0% and 52.8% in ovarian and testes samples, respectively, in these species.

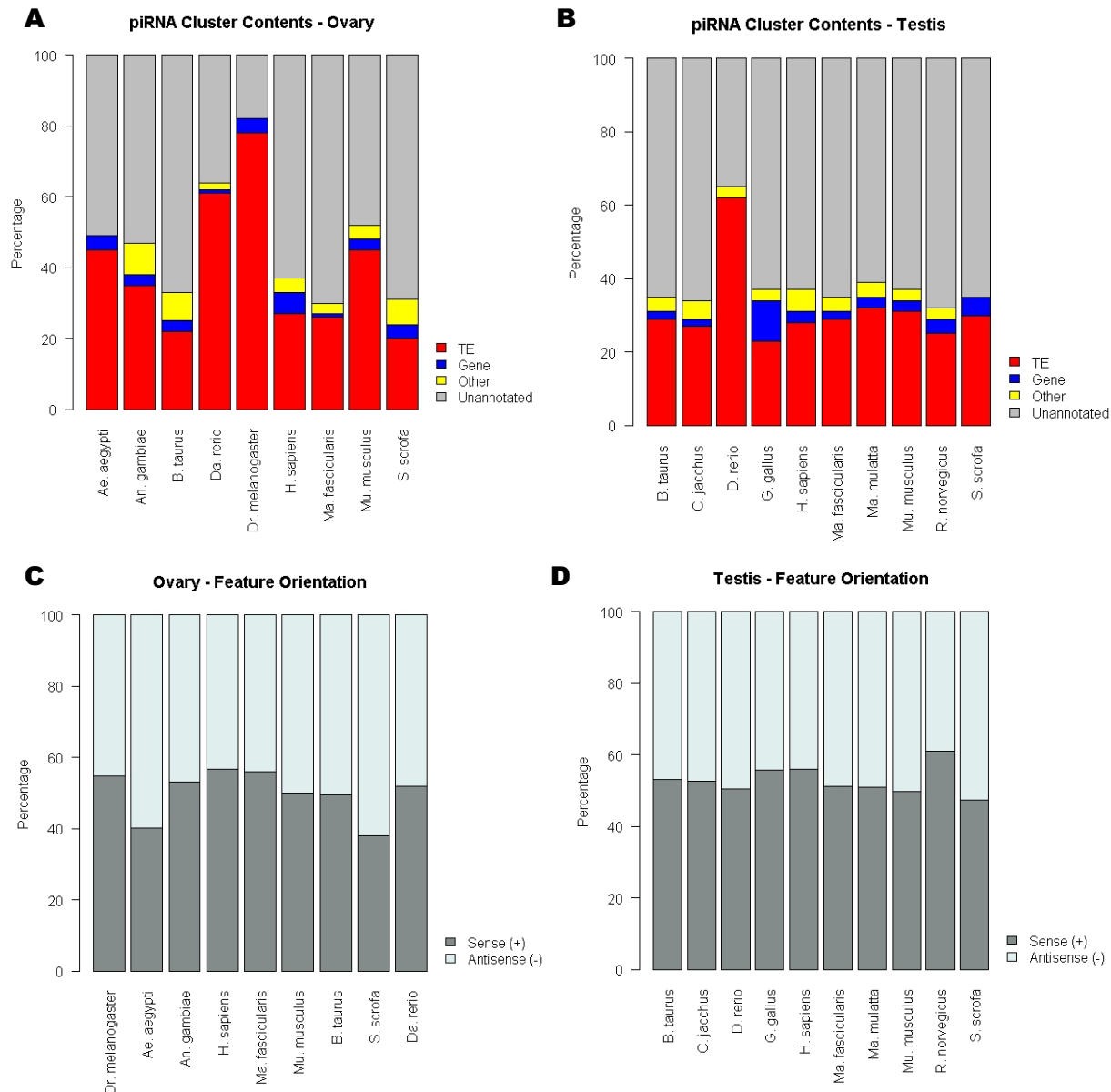


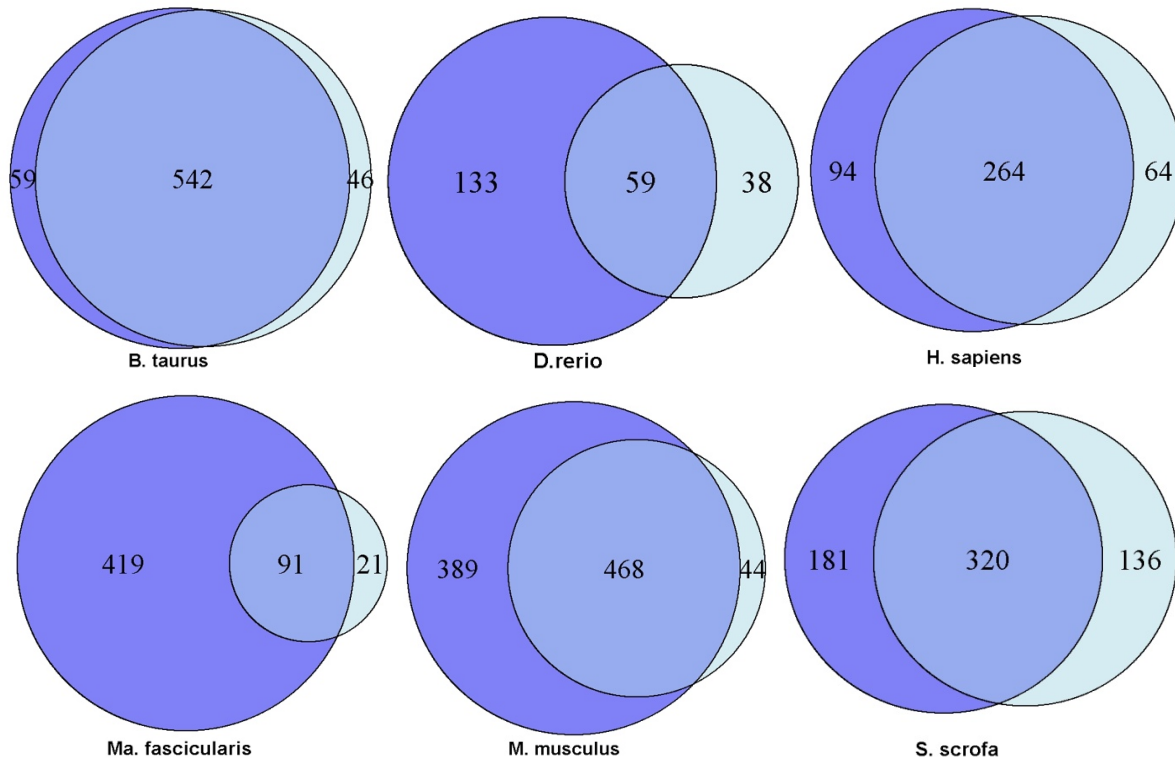
Figure 7 | Comparison of Top piRNA Cluster Composition. piRNA cluster content comparison between species in (A) ovary and (B) testis. Orientation of the feature calls in (C) ovary and (D) testis samples. These data represent an analysis of the top 30 piRNA cluster loci in each species. Available ovarian and testes datasets from the piRNA cluster database and the Short Read Archive were used to run piClusterBuster [38-39]. Only ten piRNA clusters were called in the *Dr. melanogaster* ovarian library.

Tissue-Specificity of piRNA Cluster Loci

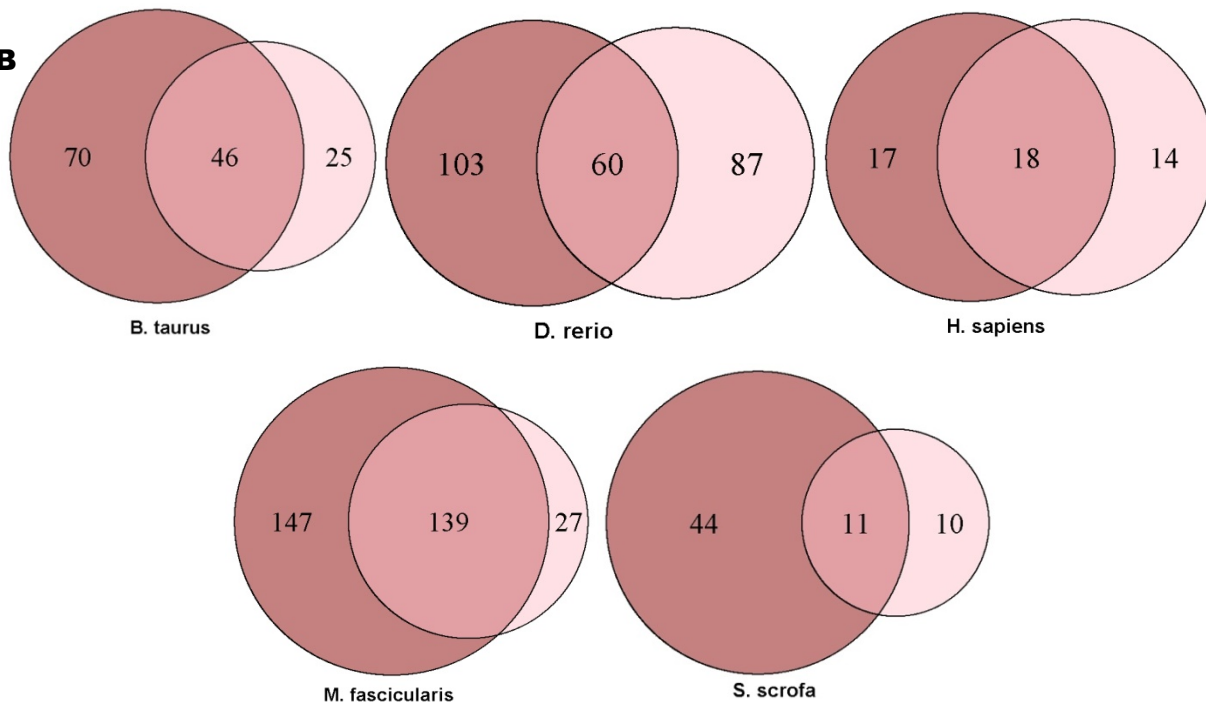
piRNA cluster definition can vary between sRNA libraries that derived from the same tissue. The number of defined piRNA clusters differed from 3 to 398 calls between two samples of the same tissue with a mean difference of 75 and median difference of 45 piRNA cluster calls.

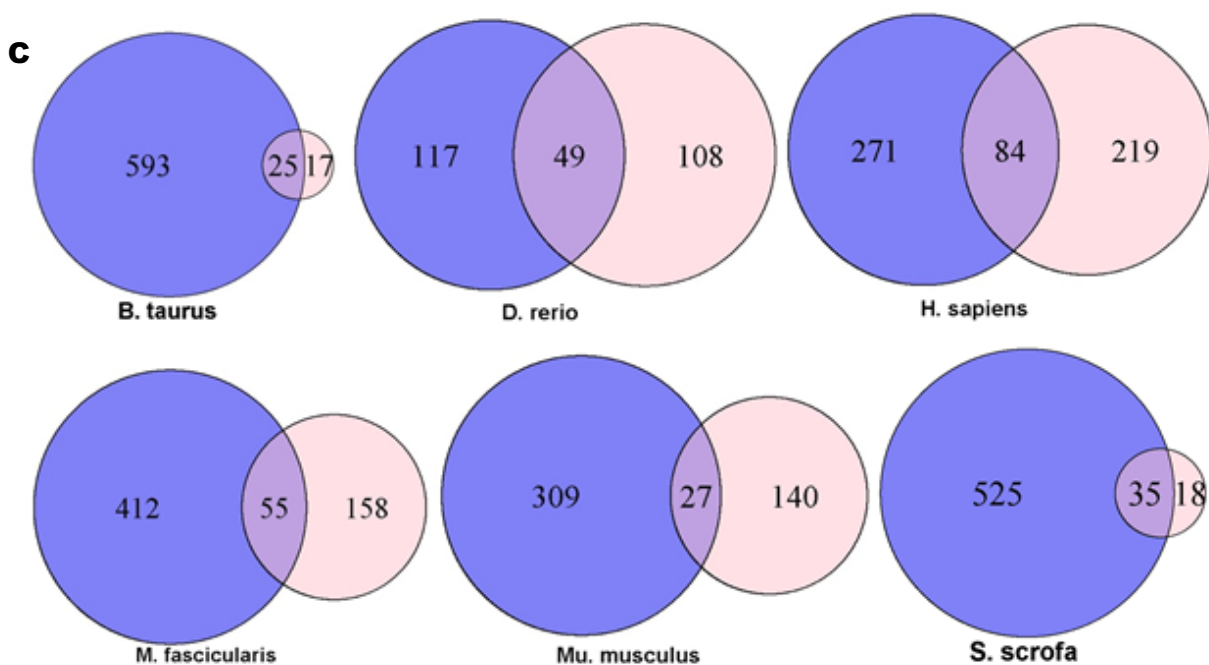
Although, at least 52.4%, and up to 92.2% of the lesser piRNA cluster definitions were also represented in the larger sample of piRNA cluster calls. piRNA cluster definition demonstrated an average of 67.8% overlap in same tissue samples (Figure 8).

A



B





D

Degree of Agreement of piRNA Cluster Definition

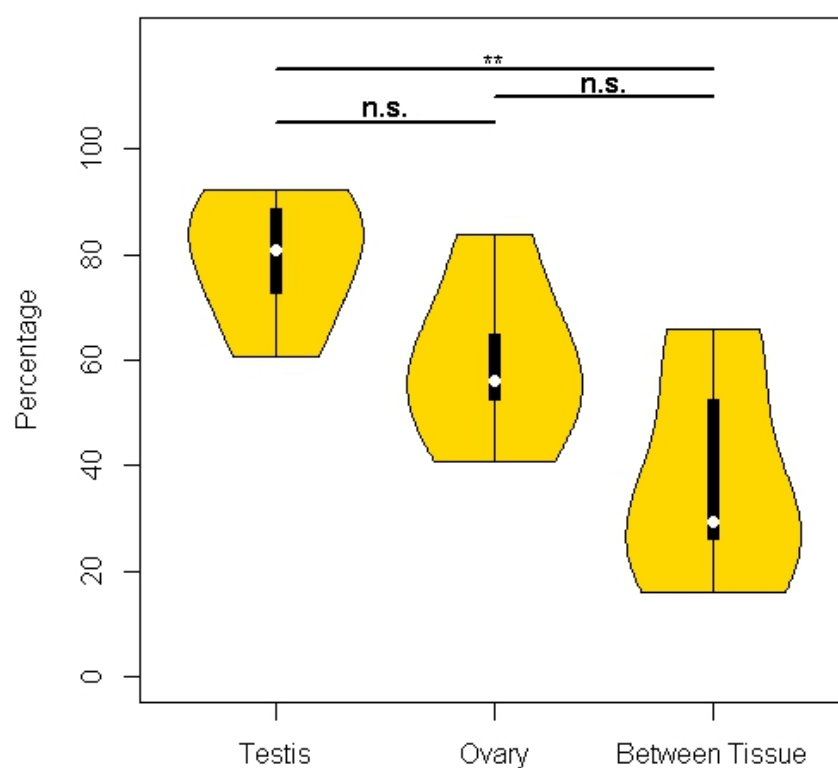


Figure 8 | Tissue-specificity of piRNA Cluster Definition. Venn Diagrams representing the degree of overlap of piRNA cluster definitions between two independent libraries of the same tissue. (A) Blue circles represent testes samples and (B) red circles represent ovarian samples. Only one ovarian sRNA library was analyzed in *M. musculus* and *D. melanogaster* testis defined no piRNA clusters. (C) Venn Diagrams representing the degree of overlap of piRNA cluster definitions between ovary and testis piRNAs. (D) Violin plot representing the agreement of piRNA cluster calls within same species testes, within same species ovaries, and a comparison between testes and ovaries piRNA cluster definition. (Figure S13)

We observed a significantly lesser degree of agreement piRNA cluster calls relative to same sample testis libraries. Same sample ovary libraries also showed a small increase in piRNA cluster agreement relative to samples between tissues (Figure 8D). The number of piRNA cluster definition differed from 9 to 576 calls between two samples of the same tissue with a mean difference of 255 and median difference of 211 piRNA cluster calls. The lesser sample of piRNA cluster definitions ranged in agreement from 16.2% to 66.0% with an average of 34.6% agreement and a median of 29.5% between samples (Figure 8C).

Discussion

The options for piClusterBuster performance enhancement allow for utilization of multitasking and multithreading. Use of multithreading allows for the execution of piClusterBuster processes by multiple nodes simultaneously. Utilization of the multitasking capability of piClusterBuster prompts independent, parallel submission for each piRNA cluster of interest to independent compute nodes. Multithreading and multitasking piClusterBuster runs allows the user the capability to significantly increase the number of piRNA clusters under observation without significantly increasing the timing of the piClusterBuster run. piClusterBuster supports both Torque/Maui or Slurm resource management software.

This comparison of piRNA cluster architecture focuses on the major genomic loci contributing to piRNA populations. By only considering only the top piRNA cluster loci in this analysis, we can be relatively confident in piRNA cluster definition relative to other piRNA-generating loci. The top 30 piRNA clusters also were a large representation of the total nucleotides occupied by piRNA clusters in these genomes, as well as disproportionally large contributors to total piRNA populations in these species (Figure 6). Taken together, the contents of the piRNA clusters in the analysis serve as the best representation of piRNA cluster architecture in these species.

Since it is difficult to determine whether RepeatMasker and CENSOR will annotate a piRNA cluster of interest more thoroughly, with higher confidence, we implemented nested annotation. A nested annotation approach allows for both of the programs that performed well in annotating sequences that are dense with repeats, RepeatMasker and CENSOR, the opportunity to characterize the sequence of interest, while only maintaining the best annotation in the

description of the contents of piRNA cluster sequence (Figure 2) [24,26]. This method allows for consistent and accurate characterization amongst diverse piRNA clusters on a large scale.

We also noted that the degree of sense or antisense orientation of feature calls within individual piRNA clusters correlated with the direction of transcription in known piRNA clusters, *flamenco* and 42AB [9,12]. Therefore, the orientation of feature calls within a piRNA cluster may be informative in the prediction of the nature of piRNA cluster transcription.

Components of piRNA architecture were strikingly similar across species. With regard to known piRNA cluster features, TEs consistently composed the majority by nucleotide occupancy and a relatively low percentage of known genic and “other” calls. The majority of informative, “other” hits within the NCBI nucleotide database were associated with mRNA that were not available in the organism-specific gene set. Other informative non-genic, non-TE sequence appeared to be of viral and rRNA origin. The orientation of feature calls within piRNA clusters were also highly conserved in these species. Taken together, these data suggest highly conserved nature, yet dynamic capacity within piRNA cluster architecture with regard to known features in Metazoans (Figure 6).

We observed that a significant portion of the piRNA cluster sequence was unable to be characterized in the species observed in this study (Figure 7). This observation prompts an interesting question regarding the derivation of piRNA cluster sequence whose origin is currently undetectable and its purpose within the piRNA clusters. This sequence is of particular biological interest given that these sequences occupy significant regions of piRNA clusters and may further inform scientific knowledge of piRNA cluster biogenesis and function.

Sets of piRNA clusters were differentially represented between different, and within the same, independent tissue samples. It is worth noting that differential representation of piRNA generating loci between same, independent tissue samples may be due to the previous observation that sRNA libraries represent only a subset of the complete sRNA populations within the cell, even when deep sequencing is performed [40]. However, the variability in piRNA cluster overlap was far greater between tissues than within tissues when comparing piRNA cluster definitions between libraries. Therefore, preliminary observation of these data supports a model in which different regions of the genome appear to be responsible for generating the majority of piRNAs in ovaries and testes samples in Metazoans and it may be advantageous for an organism to have a diverse, dynamic set of piRNA cluster activity in a unique cellular environment.

References

1. Grimson, A., *et al.* (2008). Early origins and evolution of microRNAs and Piwi-interacting RNAs in animals. *Nature*, 455(7217), 1193-1197.
2. Kidwell, Margaret G. "Evolution of hybrid dysgenesis determinants in *Drosophila melanogaster*." *Proceedings of the National Academy of Sciences* 80.6 (1983): 1655-1659.
3. Kidwell, Margaret G., James F. Kidwell, and John A. Sved. "Hybrid dysgenesis in *Drosophila melanogaster*: a syndrome of aberrant traits including mutation, sterility and male recombination." *Genetics* 86.4 (1977): 813-833.
4. Brennecke, J., Malone, C. D., Aravin, A. A., Sachidanandam, R., Stark, A., & Hannon, G. J. (2008). An epigenetic role for maternally inherited piRNAs in transposon silencing. *Science*, 322(5906), 1387-1392.
5. Kuramochi-Miyagawa, Satomi, *et al.* "Mili, a mammalian member of piwi family gene, is essential for spermatogenesis." *Development* 131.4 (2004): 839-849.
6. Houwing, Saskia, *et al.* "A role for Piwi and piRNAs in germ cell maintenance and transposon silencing in Zebrafish." *Cell* 129.1 (2007): 69-82.
7. Aravin, A. A., Hannon, G. J., & Brennecke, J. (2007). The Piwi-piRNA pathway provides an adaptive defense in the transposon arms race. *Science*, 318(5851), 761-764.
8. Arensburger, Peter, *et al.* "The mosquito *Aedes aegypti* has a large genome size and high transposable element load but contains a low proportion of transposon-specific piRNAs." *BMC Genomics* 12.1 (2011): 606.
9. Brennecke J, Aravin AA, Stark A, Dus M, Kellis M, Sachidanandam R, Hannon GJ. Discrete small RNA-generating loci as master regulators of transposon activity in *Drosophila*. *Cell* 2007, 128:1089–1103.
10. Malone CD, Brennecke J, Dus M, Stark A, McCombie WR, Sachidanandam R, Hannon GJ: Specialized piRNA pathways act in germline and somatic tissues of the *Drosophila* ovary. *Cell* 2009, 137:522–535.
11. Saito, K., Sakaguchi, Y., Suzuki, T., Suzuki, T., Siomi, H., & Siomi, M. C. (2007). Pimet, the *Drosophila* homolog of HEN1, mediates 2'-O-methylation of Piwi-interacting RNAs at their 3' ends. *Genes & development*, 21(13), 1603-1608.
12. Yin, H., & Lin, H. (2007). An epigenetic activation role of Piwi and a Piwi-associated piRNA in *Drosophila melanogaster*. *Nature*, 450(7167), 304-308.
13. Rouget, Christel. Maternal mRNA Deadenylation and Decay by the PiRNA Pathway in the Early *Drosophila* Embryo. *Nature* 467 (October 2010): 1128-132.
14. Barckmann, B., Pierson, S., Dufourt, J., Papin, C., Armenise, C., Port, F., ... & Curk, T. (2015). Aubergine iCLIP reveals piRNA-dependent decay of mRNAs involved in germ cell development in the early embryo. *Cell reports*, 12(7), 1205-1216.
15. Keam, Simon P., *et al.* "The human Piwi protein Hiwi2 associates with tRNA-derived piRNAs in somatic cells." *Nucleic acids research* 42.14 (2014): 8984-8995.
16. Kuramochi-Miyagawa, Satomi, *et al.* "Two mouse piwi-related genes: miwi and mili." *Mechanisms of development* 108.1 (2001): 121-133.
17. Houwing, Saskia, *et al.* "A role for Piwi and piRNAs in germ cell maintenance and transposon silencing in Zebrafish." *Cell* 129.1 (2007): 69-82.

18. Vodovar, Nicolas et al. “Arbovirus-Derived piRNAs Exhibit a Ping-Pong Signature in Mosquito Cells.” Ed. Sebastien Pfeffer. *PLoS ONE* 7.1 (2012): e30861. *PMC*. Web. 25 Jan. 2017.
19. Aravin, A. A., Sachidanandam, R., Bourc'his, D., Schaefer, C., Pezic, D., Toth, K. F., ... & Hannon, G. J. (2008). A piRNA pathway primed by individual transposons is linked to de novo DNA methylation in mice. *Molecular cell*, 31(6), 785-799.
20. García-López, J., de Dios Hourcade, J., Alonso, L., Cárdenas, D. B., & Del Mazo, J. (2014). Global characterization and target identification of piRNAs and endo-siRNAs in mouse gametes and zygotes. *Biochimica et Biophysica Acta (BBA)-Gene Regulatory Mechanisms*, 1839(6), 463-475.
21. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 1997;25:3389–402.
22. Jurka, J., Klonowski, P., Dagman, V., Pelton, P. (1996) CENSOR - a program for Identification and elimination of repetitive elements from DNA sequences. *Computers and Chemistry* Vol. 20 (No. 1): 119-122.
23. NCBI, R. C. (2013). Database resources of the National Center for Biotechnology Information. *Nucleic acids research*, 41(Database issue), D8.
24. Rosenkranz, D., & Zischler, H. (2012). proTRAC-a software for probabilistic piRNA cluster detection, visualization and analysis. *BMC bioinformatics*, 13(1), 1.
25. Jurka, J., Kapitonov, V. V., Pavlicek, A., Klonowski, P., Kohany, O., & Walichiewicz, J. (2005). Repbase Update, a database of eukaryotic repetitive elements. *Cytogenetic and genome research*, 110(1-4), 462-467.
26. Smit, AFA, Hubley, R. *RepeatModeler Open-1.0*. 2008-2015.
<<http://www.repeatmasker.org>>.
27. Pages, H., Aboyoun, P., Gentleman, R., & DebRoy, S. (2009). String objects representing biological sequences, and matching algorithms. *R package version*, 2(2).
28. Analytics, R. doMC: Foreach parallel adaptor for the multicore package, 2011a. *URL* <http://cran.r-project.org/package=doMC>. *R package version*, 1(2).
29. Lawrence, M., Huber, W., Pages, H., Aboyoun, P., Carlson, M., Gentleman, R., ... & Carey, V. J. (2013). Software for computing and annotating genomic ranges. *PLoS Comput Biol*, 9(8), e1003118.
30. Reimand, J., Kull, M., Peterson, H., Hansen, J., & Vilo, J. (2007). g: Profiler—a web-based toolset for functional profiling of gene lists from large-scale experiments. *Nucleic acids research*, 35(suppl 2), W193-W200.
31. Wickham, H. (2009). plyr: Tools for splitting, applying and combining data. *R package version 0.1*, 9, 651.
32. Scrucca, L. (2004). qcc: an R package for quality control charting and statistical process control. *dim (pistonrings)*, 1(200), 3.
33. Charif, D., & Lobry, J. R. (2007). SeqinR 1.0-2: a contributed package to the R project for statistical computing devoted to biological sequences retrieval and analysis. In *Structural approaches to sequence evolution* (pp. 207-232). Springer Berlin Heidelberg.
34. Girke, T. (2014). systemPipeR: NGS workflow and report generation environment. *UC Riverside*. <https://github.com/tgirke/systemPipeR>.
35. Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, 26(6), 841-842.

36. Attrill H, Falls K, Goodman JL, Millburn GH, Antonazzo G, Rey AJ, Marygold SJ; the FlyBase Consortium. (2016) FlyBase: establishing a Gene Group resource for *Drosophila melanogaster*. *Nucleic Acids Res.* 44(D1):D786-D792.
37. Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D, Kent WJ. The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* 2004 Jan 1;32(Database issue):D493-6.
38. Rosenkranz D. piRNA cluster database: a web resource for piRNA producing loci. *Nucleic Acids Research* 2016 44(D1):D223-D230.
39. Leinonen, R., Sugawara, H., & Shumway, M. (2010). The sequence read archive. *Nucleic acids research*, gkq1019.
40. Yamtich, Jennifer, et al. "piRNA-like small RNAs mark extended 3'UTRs present in germ and somatic cells." *BMC genomics* 16.1 (2015): 1.