

1 **Title: *FishPi*: a bioinformatic prediction tool to link piRNA and transposable elements**  
2 **in zebrafish**

3

4 Alice M. Godden<sup>1</sup>, Benjamin Rix<sup>1</sup>, Simone Immler<sup>1</sup>

5

6 <sup>1</sup> School of Biological Sciences, University of East Anglia, Norwich Research Park, Norwich,  
7 NR4 7TJ, United Kingdom

8

9 Corresponding author:

10 **Alice M. Godden** Email: [alice.godden@uea.ac.uk](mailto:alice.godden@uea.ac.uk)

11

12

13

14

15

16 **Abstract**

17 **Background**

18 Piwi-interacting RNAs (piRNA)s are non-coding small RNAs that post-transcriptionally  
19 affect gene expression and regulation. Through complementary seed region binding with  
20 transposable elements (TEs), piRNAs protect the genome from transposition, and  
21 therefore a tool to link piRNAs with complementary TE targets is needed. Tools like  
22 *TEsmall* can process sRNA-seq datasets to produce differentially expressed piRNAs and  
23 *piRScan* developed for nematodes can link piRNAs and TEs but it requires the user to  
24 know the target region of interest and work backwards.

25 **Results**

26 We have therefore developed *FishPi* to predict the pairings between piRNA and TEs.  
27 *FishPi* works with individual piRNAs or a list of piRNA sequences in fasta format. The  
28 software focuses on the piRNA:TE seed region and analyses reference TEs for piRNA  
29 complementarity. TE type is examined, counted and stored to a dictionary, with genomic  
30 loci recorded. Any updates to piRNA-TE binding rules, can easily be incorporated by  
31 updating the code underlying *FishPi*. *FishPi* provides a graphic interface, using *tkinter*,  
32 that requires the user to input piRNA sequences to generate comprehensive reports on  
33 piRNA:TE dynamics. *FishPi* can easily be adapted to other genomes opening the  
34 interpretation of piRNA functionality to a wide community.

35 **Conclusions**

36 Users will gain insight into genome age and *FishPi* will help further our understanding of  
37 the biological role of piRNAs and their interaction with TEs in a similar way that public  
38 databases have improved the access to and the understanding of the role of small  
39 RNAs.

40

41 **Keywords**

42 PiRNAs, transposons, evolution, zebrafish

## 43 Background

44 Environmental changes and related biotic and abiotic stressors can disrupt the  
45 balance and expression of small RNAs and transposable elements (TEs) in the genome <sup>1</sup>.  
46 TEs are parasitic DNA that can cause insertions in the genome, as well as deletions and  
47 genomic rearrangements <sup>2</sup>. To counter this activity, Piwi-interacting RNAs (piRNAs), have  
48 evolved to repress the action of these repetitive elements <sup>3</sup>. piRNAs constitute a substantial  
49 fraction of small RNAs and are predominantly found in animal germ cells <sup>4</sup>. The main  
50 function of piRNAs is to silence TEs and repetitive elements in animal germ cells to preserve  
51 the integrity of the germline genome <sup>5</sup> and loss of TE silencing by piRNAs can lead to sterility  
52 for example in the nematode *Caenorhabditis elegans* <sup>6</sup>. TEs are parasitic and are not  
53 randomly distributed in the genome but have preferred regions for insertion in the genome,  
54 including distorted and bent DNA <sup>7,8</sup>, or open chromatin regions <sup>8,9</sup>. Large populations of  
55 piRNAs are complementary to TE sequences and messenger RNAs (mRNA) <sup>10</sup>. In *C.*  
56 *elegans*, the 2<sup>nd</sup> to the 7<sup>th</sup> nucleotide at the 5' end of the piRNA <sup>11</sup> or the 2<sup>nd</sup> to the 8<sup>th</sup>  
57 nucleotide <sup>12</sup> have both been suggested to be the piRNA:TE seed region. Similarly, in *Mus*  
58 *musculus*, the 2<sup>nd</sup> to 8<sup>th</sup> nucleotide was a predicted seed region <sup>13</sup> and in *Drosophila*  
59 *melanogaster*, the first 10 nucleotides at the 5' end of the piRNA are predicted to be the  
60 piRNA:TE seed binding site <sup>14</sup>. Finally, in zebrafish and other teleost fish, the first 10  
61 nucleotides at the 5' end of the piRNA are loaded into the Argonaute machinery to produce  
62 Ago2-bound piRNAs complementary to TEs <sup>3,15</sup>.

63 TEs are known to be activated through environmental stressors. In *C. elegans* for example,  
64 thermal stress has been shown in mutant piRNA biogenesis strains that lack the ability to  
65 generate piRNAs. In this work nematodes were heat-treated and displayed transgenerational  
66 impairments for up to three subsequent generations: with thermal stress decreasing the rate  
67 of piRNA biogenesis. <sup>16,17</sup>. In *Drosophila*, the *mariner-Mos1* transposon is enriched following  
68 heat-stress but not following ultraviolet irradiation <sup>18</sup>. TEs are expressed in a tissue-specific  
69 manner and can affect transcription and gene expression <sup>19</sup>, by generating new genes and

70 mRNAs through genomic rearrangements and *de novo* mutations, and cause alteration of  
71 regulatory networks and even trigger immune responses <sup>7</sup>. Understanding the interaction  
72 between piRNAs and their target activities is therefore key to understand the defence  
73 mechanisms in the germ line against the effects of environmental conditions.

74 In the zebrafish (*Danio rerio*) genome, small RNAs post-transcriptionally regulate gene  
75 expression and regulation <sup>20</sup>. Piwi proteins are predominantly expressed in the animal germ  
76 line and associate with piRNAs. The loss of Ziwi (zebrafish Piwi) leads to apoptosis of germ  
77 cells in embryonic development <sup>3</sup>. Ziwi and Zili zebrafish Piwi proteins provide a piRNA  
78 amplification system that is highly conserved. TEs represent over 50% of the zebrafish  
79 genome <sup>21</sup> and it is therefore key to improve our understanding of the mechanisms that  
80 evolved to allow an organism to regulate their activity. TEs have over 2000 families in the  
81 zebrafish and exploit the host's genome by dynamic transposition events. The mechanisms  
82 underlying this transposition can be divided into two categories: retrotransposons, these TEs  
83 synthesise an RNA intermediate that can insert itself into the host genome <sup>22</sup>. In contrast,  
84 DNA transposons cut-and-paste themselves into the host genome <sup>7,22,23</sup>.

85 We therefore developed the software FishPi to provide a means to map zebrafish piRNAs  
86 against complementary TEs and link them to understand their biological association <sup>10</sup>.  
87 FishPi facilitates research into the function of piRNAs and accelerates analyses in the same  
88 way as microRNA databases like Targetscan have supported small RNA research <sup>24</sup>. FishPi  
89 provides a user-friendly graphic interface (Figure 1), which is based on the latest zebrafish  
90 reference genome and TE sequences and generates a comprehensive output that can be  
91 exported as a report including complementary TE counts, family level annotation,  
92 chromosomal and counts plots, with the option to export the full list of TEs and their  
93 sequences that are complementary to the piRNA sequence fed in by the user. The user can  
94 enter individual piRNAs or lists of piRNAs in fasta format. For zebrafish and other model  
95 organisms, a list of piRNA sequences are available at piRBase

96 (<http://bigdata.ibp.ac.cn/piRBase/browse.php>),<sup>4</sup> or piRNAclusterDB  
97 (<https://www.smallnagroup.uni-mainz.de/piRNAclusterDB/>),<sup>25</sup>.

## 98 **Implementation**

99 FishPi is a python-based programme and uses the following packages and versions: Python  
100 v3.11, Tkinter (included in Python v3.11<sup>26</sup>), Pillow 10.0.1<sup>27</sup>, and Matplotlib 3.8.0<sup>28</sup>. The  
101 code for the programme was written and tested using PyCharm Community Edition 2023.2.1  
102 and is platform-independent, as it is designed to be executed from the Linux command line;  
103 however, the software can be run and used in PyCharm IDE. The software initially defines  
104 functions to assay piRNA TE complementarity, creating a dictionary to count the TEs by  
105 class. The families assigned to each class were checked using dfam.org (Table 1),<sup>29</sup>.

106 To generate a fasta file of the TE sequences to use in place of the zebrafish reference  
107 genome GRCz11 TE sequences, a TE.bed file with clade was created by using the UCSC  
108 table genome browser<sup>30</sup> as follows: Vertebrate, group: Variation and Repeats, genome:  
109 zebrafish, assembly: May 2017 GRCz11, track: RepeatMasker, table: rmsk. To extract DNA  
110 sequences from the reference genome based on the co-ordinates supplied in the TE.bed  
111 file, the following command lines were entered:

```
112 bedtools getfasta -s -name -fi Danio_rerio.GRCz11.dna.primary_assembly.fa -fo  
113 GRCz11.teseqs.use.fasta -bed GRCz11.teannotation.bed
```

114 The following command line serves to mask the sequences in the reference genome:

```
115 bedtools maskfasta -fi Danio_rerio.GRCz11.dna.primary_assembly.fa -fo  
116 GRCz11.masked.fasta -bed GRCz11.teannotation.bed
```

117 To create the TE-merged reference genome<sup>31</sup>:

```
118 cat GRCz11.masked.fasta GRCz11.teseqs.1.fasta > GRCz11.teseqs.fishpi.fasta
```

119

120 This preparatory work was based on the pipeline described for PopoolationTE2<sup>32</sup>. The  
121 Ensembl reference genome GRCz11 v107 was used in this project and can be retrieved  
122 here: [https://ftp.ensembl.org/pub/release-107/fasta/danio\\_erio/dna/](https://ftp.ensembl.org/pub/release-107/fasta/danio_erio/dna/\)  
123 [\"Danio\\_erio.GRCz11.dna.primary\\_assembly.fa.gz\"](https://ftp.ensembl.org/pub/release-107/fasta/danio_erio/dna/Danio_erio.GRCz11.dna.primary_assembly.fa.gz)<sup>33</sup>. These instructions can be adjusted  
124 for other reference genomes and allow the user to run the analyses on any available  
125 genome version.

## 126 **Results & Discussion**

127 FishPi successfully links piRNAs and TEs and is a helpful tool for anyone interested in  
128 studying the role of TEs and sRNAs in shaping the germline genome. FishPi is designed to  
129 be used by those with minimal coding experience. A detailed walk-through use and  
130 installation is listed in the readme on the git repository here:  
131 <https://github.com/alicegodden/fishpi>. The simple graphic user interface offers a point and  
132 click tool to rapidly gain insight into many piRNA sequences. Each time a new piRNA  
133 sequence is provided, the programme resets itself before providing the new results. As more  
134 research into piRNA:TE seed sites evolve, FishPi can be adapted and updated to reflect  
135 changes and updates to the reference genome. Alternatively, users can generate and use  
136 other reference TE genome files and use FishPi for other model systems. This makes FishPi  
137 particularly powerful for users looking at non-model systems. The software has a graphic  
138 user interface for FishPi (Figure 1), and the outputs can be exported as hi-resolution bar-  
139 charts and chromosomal plots presenting raw count data of the different target TE  
140 categories for presentation (Figure 1, Supplementary Files 2 & 3). Additionally, users can  
141 export a .csv file with all complementary TE mRNA sequences, along with genomic co-  
142 ordinates, allowing the user to conduct further downstream analyses and explore each class  
143 of TE individually (Supplementary File 1). Benchmarking data can be found in  
144 Supplementary File 4 for optimal running of the software. FishPi's advantage over other tools  
145 is that it can be easily updated and modified as more research into piRNA:TE seeds are

146 validated and will be fully customisable for future genome releases for other understudied  
147 species including fish and other model systems.

## 148 **Conclusions**

149 By analysing piRNA complementarity to a TE, we can infer recent genomic expansion but  
150 also adaptive evolution. For DNA class TEs, we can find footprints of evolutionary activity, as  
151 DNA TEs are often associated with past transposition events. In the zebrafish genome, DNA  
152 TEs comprise 40%, with retrotransposons comprising 10% zebrafish genome<sup>23</sup>. Currently  
153 available tools to link piRNA and TEs are limited and require knowledge of known targets or  
154 are restricted to a pre-defined genome and species due to piRNA binding rules<sup>11-13,34</sup>.  
155 Further analysis with RNA-seq data and tools such as RepeatMasker and RepeatLandscape  
156 would complement this by Kimura distance analysis<sup>35</sup>. Additionally, use of the Retroseq  
157 pipeline can highlight non-reference insertion mutations by TEs from high-coverage DNA-  
158 seq data and help predict TE activity<sup>36</sup>. In the future, other adaptations of the software could  
159 be made to look at other small RNAs. The software allows users to rapidly confirm  
160 predictions of complementary TEs from specific piRNAs. FishPi software is freely available  
161 and is distributed under GPL-3.0 licence.

## 162 **Availability and requirements**

163 **Project name:** FishPi

164 **Project home page:** <https://github.com/alicegodden/fishpi/tree/fishpi>

165 **Operating system(s):** Platform independent

166 **Programming language:** Python

167 **Other requirements:** Python v3.11

168 **License:** GPL-3.0

169 **Any restrictions to use by non-academics:** GPL-3.0 license needed

170

171

172 **Declarations**

173 **Ethics**

174 No ethics approvals were required for this research project.

175 **Consent for publication**

176 All authors consent to publication.

177 **Availability of data and materials**

178 FishPi is available on the Git repository here: <https://github.com/alicegodden/fishpi/tree/fishpi>

179 **Competing interests**

180 Authors declare no financial or competing interests.

181 **Funding**

182 This project was funded by a Consolidator Grant from the European Research Council to SI

183 (SELECTHAPLOID – 101).

184 **Author's contributions**

185 AMG conceived the idea for the project, AMG and BR wrote the software, BR tested the  
186 software, AMG wrote the manuscript and SI contributed to the writing and provided guidance  
187 throughout.

188 **Acknowledgements**

189 AMG would like to acknowledge Prof. Oliver Buckley from University of East Anglia for  
190 software publication and licensing advice.

191

192 **List of abbreviations**

193 mRNA- Messenger RNA

194 PiRNA- Piwi-Interacting RNA

195 TE- Transposable element

196

197

198 **References**

199

- 200 1 Godden, A. M. & Immler, S. The potential role of the mobile and non-coding  
201 genomes in adaptive response. *Trends Genet* (2022).  
202 <https://doi.org/10.1016/j.tig.2022.08.006>
- 203 2 Kofler, R. piRNA Clusters Need a Minimum Size to Control Transposable Element  
204 Invasions. *Genome Biol Evol* **12**, 736-749 (2020).  
205 <https://doi.org/10.1093/gbe/evaa064>
- 206 3 Houwing, S. *et al.* A role for Piwi and piRNAs in germ cell maintenance and  
207 transposon silencing in Zebrafish. *Cell* **129**, 69-82 (2007).  
208 <https://doi.org/10.1016/j.cell.2007.03.026>



- 209 4 Wang, J. *et al.* piRBase: a comprehensive database of piRNA sequences. *Nucleic*  
210 *Acids Res* **47**, D175-D180 (2019). <https://doi.org/10.1093/nar/gky1043>
- 211 5 Brennecke, J. *et al.* Discrete small RNA-generating loci as master regulators of  
212 transposon activity in *Drosophila*. *Cell* **128**, 1089-1103 (2007).  
213 <https://doi.org/10.1016/j.cell.2007.01.043>
- 214 6 Spichal, M. *et al.* Germ granule dysfunction is a hallmark and mirror of Piwi mutant  
215 sterility. *Nat Commun* **12**, 1420 (2021). <https://doi.org/10.1038/s41467-021-21635-0>
- 216 7 Bourque, G. *et al.* Ten things you should know about transposable elements. *Genome*  
217 *Biol* **19**, 199 (2018). <https://doi.org/10.1186/s13059-018-1577-z>
- 218 8 Sultana, T., Zamborlini, A., Cristofari, G. & Lesage, P. Integration site selection by  
219 retroviruses and transposable elements in eukaryotes. *Nat Rev Genet* **18**, 292-308  
220 (2017). <https://doi.org/10.1038/nrg.2017.7>
- 221 9 Gebrie, A. Transposable elements as essential elements in the control of gene  
222 expression. *Mob DNA* **14**, 9 (2023). <https://doi.org/10.1186/s13100-023-00297-3>
- 223 10 Houwing, S., Berezikov, E. & Ketting, R. F. Zili is required for germ cell  
224 differentiation and meiosis in zebrafish. *EMBO J* **27**, 2702-2711 (2008).  
225 <https://doi.org/10.1038/emboj.2008.204>
- 226 11 Zhang, D. *et al.* The piRNA targeting rules and the resistance to piRNA silencing in  
227 endogenous genes. *Science* **359**, 587-592 (2018).  
228 <https://doi.org/10.1126/science.aao2840>
- 229 12 Shen, E. Z. *et al.* Identification of piRNA Binding Sites Reveals the Argonaute  
230 Regulatory Landscape of the *C. elegans* Germline. *Cell* **172**, 937-951 e918 (2018).  
231 <https://doi.org/10.1016/j.cell.2018.02.002>
- 232 13 Gou, L. T. *et al.* Pachytene piRNAs instruct massive mRNA elimination during late  
233 spermiogenesis. *Cell Res* **24**, 680-700 (2014). <https://doi.org/10.1038/cr.2014.41>
- 234 14 Aravin, A. A., Hannon, G. J. & Brennecke, J. The Piwi-piRNA pathway provides an  
235 adaptive defense in the transposon arms race. *Science* **318**, 761-764 (2007).  
236 <https://doi.org/10.1126/science.1146484>
- 237 15 La, Y. *et al.* Identification and Characterization of Piwi-Interacting RNAs for Early  
238 Testicular Development in Yak. *Int J Mol Sci* **23** (2022).  
239 <https://doi.org/10.3390/ijms232012320>
- 240 16 Belicard, T., Jareosettasin, P. & Sarkies, P. The piRNA pathway responds to  
241 environmental signals to establish intergenerational adaptation to stress. *BMC Biol* **16**,  
242 103 (2018). <https://doi.org/10.1186/s12915-018-0571-y>
- 243 17 Kurhanewicz, N. A., Dinwiddie, D., Bush, Z. D. & Libuda, D. E. Elevated  
244 Temperatures Cause Transposon-Associated DNA Damage in *C. elegans*  
245 Spermatocytes. *Curr Biol* **30**, 5007-5017 e5004 (2020).  
246 <https://doi.org/10.1016/j.cub.2020.09.050>
- 247 18 Jardim, S. S., Schuch, A. P., Pereira, C. M. & Loreto, E. L. Effects of heat and UV  
248 radiation on the mobilization of transposon mariner-Mos1. *Cell Stress Chaperones* **20**,  
249 843-851 (2015). <https://doi.org/10.1007/s12192-015-0611-2>
- 250 19 He, J. *et al.* Identifying transposable element expression dynamics and heterogeneity  
251 during development at the single-cell level with a processing pipeline scTE. *Nat*  
252 *Commun* **12**, 1456 (2021). <https://doi.org/10.1038/s41467-021-21808-x>
- 253 20 Wei, C., Salichos, L., Wittgrove, C. M., Rokas, A. & Patton, J. G. Transcriptome-  
254 wide analysis of small RNA expression in early zebrafish development. *RNA* **18**, 915-  
255 929 (2012). <https://doi.org/10.1261/rna.029090.111>
- 256 21 Huang, C. R., Burns, K. H. & Boeke, J. D. Active transposition in genomes. *Annu Rev*  
257 *Genet* **46**, 651-675 (2012). <https://doi.org/10.1146/annurev-genet-110711-155616>

- 258 22 Boeke, J. D., Garfinkel, D. J., Styles, C. A. & Fink, G. R. Ty elements transpose  
259 through an RNA intermediate. *Cell* **40**, 491-500 (1985). [https://doi.org/10.1016/0092-](https://doi.org/10.1016/0092-8674(85)90197-7)  
260 [8674\(85\)90197-7](https://doi.org/10.1016/0092-8674(85)90197-7)
- 261 23 Chang, N. C., Rovira, Q., Wells, J., Feschotte, C. & Vaquerizas, J. M. Zebrafish  
262 transposable elements show extensive diversification in age, genomic distribution,  
263 and developmental expression. *Genome Res* **32**, 1408-1423 (2022).  
264 <https://doi.org/10.1101/gr.275655.121>
- 265 24 Agarwal, V., Bell, G. W., Nam, J. W. & Bartel, D. P. Predicting effective microRNA  
266 target sites in mammalian mRNAs. *Elife* **4** (2015).  
267 <https://doi.org/10.7554/eLife.05005>
- 268 25 Rosenkranz, D., Zischler, H. & Gebert, D. piRNAclusterDB 2.0: update and  
269 expansion of the piRNA cluster database. *Nucleic Acids Res* **50**, D259-D264 (2022).  
270 <https://doi.org/10.1093/nar/gkab622>
- 271 26 Foundation, P. S. Python Language reference, vesion 3.11.  
272 <https://python.org/downloads/release/python-3110/> (2022).  
273 <https://doi.org/https://python.org/downloads/release/python-3110/>
- 274 27 Clark, A. <https://buildmedia.readthedocs.org/media/pdf/pillow/latest/pillow.pdf>  
275 (readthedocs, 2015).
- 276 28 Hunter, J. D. Matplotlib: A 2D graphics environment. *Computing in Science and*  
277 *Engineering* **9**, 90-95 (2007). <https://doi.org/10.1109/MCSE.2007.55>
- 278 29 Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community  
279 resource of transposable element families, sequence models, and genome annotations.  
280 *Mob DNA* **12**, 2 (2021). <https://doi.org/10.1186/s13100-020-00230-y>
- 281 30 Karolchik, D. *et al.* The UCSC Table Browser data retrieval tool. *Nucleic Acids Res*  
282 **32**, D493-496 (2004). <https://doi.org/10.1093/nar/gkh103>
- 283 31 Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing  
284 genomic features. *Bioinformatics* **26**, 841-842 (2010).  
285 <https://doi.org/10.1093/bioinformatics/btq033>
- 286 32 Kofler, R., Gomez-Sanchez, D. & Schlotterer, C. PoPoolationTE2: Comparative  
287 Population Genomics of Transposable Elements Using Pool-Seq. *Mol Biol Evol* **33**,  
288 2759-2764 (2016). <https://doi.org/10.1093/molbev/msw137>
- 289 33 Cunningham, F. *et al.* Ensembl 2022. *Nucleic Acids Res* **50**, D988-D995 (2022).  
290 <https://doi.org/10.1093/nar/gkab1049>
- 291 34 Wu, W. S. *et al.* pirScan: a webserver to predict piRNA targeting sites and to avoid  
292 transgene silencing in *C. elegans*. *Nucleic Acids Res* **46**, W43-W48 (2018).  
293 <https://doi.org/10.1093/nar/gky277>
- 294 35 Smit, A. F., Hubley, R. & Green, P. RepeatMasker software (2013).  
295 <https://doi.org/http://www.repeatmasker.org>
- 296 36 Keane, T. M., Wong, K. & Adams, D. J. RetroSeq: transposable element discovery  
297 from next-generation sequencing data. *Bioinformatics* **29**, 389-390 (2013).  
298 <https://doi.org/10.1093/bioinformatics/bts697>

300

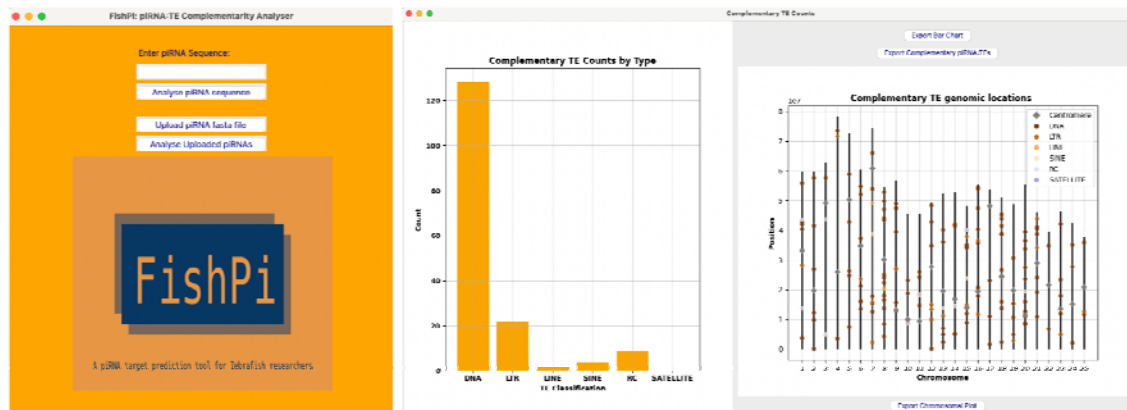
301

## Figures and Tables

302 *Table 1- Transposable element families assigned to transposable element Class (verified*  
 303 *using dfam.org).*

Class	Transposable element family
DNA	hAT, Tc1, Tc-Mar, Harbinger, Enspm, Kolobok, Merlin, Crypton, PiggyBac, Dada, Zatar, Ginger, TDR, Polinton, Maverick, Acrobat, Looper, TZF, Angel, Mariner
LTR	Gypsy, DIRS, Ngarg, ERV, Pao, Copia, BEL, HERV, Bhikari
LINE	L1, L2, L1-Tx1, Rex-Babar, RTE, Penelope, Keno, Rex
SINE	Alu, tRNA-V-RTE
RC	Helitron
Satellite	BRSATI, MOSAT

304



305

306 *Figure 1- The graphic user interface of FishPi with and pop-out window with exportable*  
 307 *results for dre-piRNA-53095 5'- CGGATAAGCATATTGCCCTAGCAACATC - 3'. The buttons*  
 308 *included allow users to export the plots at hi-resolution for presentation and publication.*  
 309 *Additionally, the user may export a .csv file list of all complementary TE mRNAs for further*  
 310 *downstream analysis.*