

1 **Gene markers for exon capture and phylogenomics in ray-finned fishes**

2 Jiamei Jiang, Hao Yuan, Xin Zheng, Qian Wang, Ting Kuang, Jingyan Li, Junning Liu,  
3 Shuli Song, Weicai Wang, Fangyuan Cheng, Hongjie Li, Junman Huang, and Chenhong  
4 Li\*

5 Shanghai Universities Key Laboratory of Marine Animal Taxonomy and Evolution

6 Key Laboratory of Exploration and Utilization of Aquatic Genetic Resources (Shanghai  
7 Ocean University), Ministry of Education, Shanghai 201306, China

8 National Demonstration Center for Experimental Fisheries Science Education (Shanghai  
9 Ocean University)

10 \*Corresponding author: Chenhong Li, e-mail: [chli@shou.edu.cn](mailto:chli@shou.edu.cn)

11

11 **Abstract**

12 Gene capture coupled with the next generation sequencing has become one of the favorable  
13 methods in subsampling genomes for phylogenomic studies. Many target gene markers  
14 have been developed in plants, sharks, frogs, reptiles and others, but few have been  
15 reported in the ray-finned fishes. Here, we identified a suite of “single-copy” protein coding  
16 sequence (CDS) markers through comparing eight fish genomes, and tested them  
17 empirically in 83 species (33 families and 11 orders) of ray-finned fishes. Sorting through  
18 the markers according to their completeness and phylogenetic decisiveness in taxa tested  
19 resulted in a selection of 4,434 markers, which were proven to be useful in reconstructing  
20 phylogenies of the ray-finned fishes at different taxonomic level. We also proposed a  
21 strategy of refining baits (probes) design a posteriori based on empirical data. The markers  
22 that we have developed may fill a gap in the tool kit of phylogenomic study in vertebrates.

23 *Key words:* Actinopterygii, target enrichment, nuclear gene markers, phylogenomics,  
24 population genomics, baits design.

25

## 25 **Introduction**

26 Next-generation sequencing (NGS) drastically reduced the cost of genome sequencing, so  
27 that reconstructing phylogenetic relationships using whole genomes became feasible (Jarvis,  
28 et al. 2014). However, sequencing whole genomes is still costly and sometime unnecessary.  
29 Subsampling genome sequences has gained popularity in phylogenomics and population  
30 genomics in recent years (Emerson, et al. 2010; Faircloth, et al. 2012; Lemmon, et al. 2012;  
31 Peterson, et al. 2012; Li, et al. 2013). There are two camps that prefer different genome  
32 subsampling tools. One is associated with restriction site related markers, such as restriction  
33 site associated DNA (RAD) (Baird, et al. 2008) and double digest RADseq (ddRAD)  
34 markers (Peterson, et al. 2012), which could be used to produce sequences from a  
35 tremendous number of anonymous loci, particularly useful in studying population genomics  
36 or species-level phylogeny (Davey and Blaxter 2010). The other camp uses methods of  
37 gene capture, also known as target enrichment to capture and sequence target loci, which  
38 often result in less missing data than the restriction site related methods does (Collins and  
39 Hrbek 2015), and the target loci can be applied across highly divergent taxonomic groups  
40 (Faircloth, et al. 2012; Lemmon, et al. 2012; Li, et al. 2013).

41 Gene capture is based on hybridizing RNA/DNA baits (probes) to DNA library of  
42 targeted species and pulling out sequences similar to the baits for subsequent  
43 high-throughput sequencing. Two popular methods, Ultraconserved Element Captures  
44 (UCE) (Faircloth, et al. 2012) and Anchored Hybrid Enrichment (AHE) (Lemmon, et al.  
45 2012) were developed to pull out highly conserved elements in the genome along with  
46 variable flanking regions. Both UCE and AHE methods were designed to anchor highly  
47 conserved regions of the genome and make use of variation in flanking sequences. A third  
48 method, exon capture was designed explicitly to capture single-copy coding sequences  
49 across moderate to highly divergent species (Bi, et al. 2012; Hedtke, et al. 2013; Li, et al.  
50 2013). The advantage of exon capture is that exon sequences are easier to align and better  
51 studied for phylogenetics than anonymous non-coding regions. Furthermore, lowered

52 stringency in hybridization and washing steps of exon capture can generate data from more  
53 loci than methods focused only on highly conserved elements.

54 Exon capture markers have been developed in plants (Mandel, et al. 2014; Weitemier, et  
55 al. 2014; Chamala, et al. 2015), invertebrates (Hugall, et al. 2016; Mayer, et al. 2016;  
56 Teasdale, et al. 2016; Yuan, et al. 2016), and many vertebrate groups, including sharks and  
57 skates (Li, et al. 2013), frogs (Hedtke, et al. 2013; Portik, et al. 2016), skink lizards (Bragg,  
58 et al. 2016) and others, yet few exon markers have been reported in the ray-finned fishes  
59 (Actinopterygii), the most diverse group of vertebrates with more than 30,000 described  
60 species (Nelson, et al. 2016). Ilves and Lopez-Fernandez (2014) developed 923 exon  
61 markers for cichlids based on genome sequence of tilapia, but those markers probably are  
62 too specialized to be used on other ray-finned fishes. We also developed 17,817 single-copy  
63 nuclear coding (CDS) markers and applied those in the siniperid fish, but those markers  
64 have not been tested in other ray-finned fishes (Song, et al. 2017).

65 Selecting target markers and designing baits that are effective across a wide range of  
66 species is the first major challenge when applying the gene capture method. Many  
67 considerations have been taken into baits design, such as uniqueness and conserveness of  
68 markers, length and complexity of markers, and genetic distance between baits and target  
69 sequences (Bi, et al. 2012; Faircloth, et al. 2012; Lemmon, et al. 2012; Li, et al. 2013;  
70 Mayer, et al. 2016). However, all these measures were taken a priori, and nothing has been  
71 done to refine baits design after gene capture to improve the baits set for future  
72 experiments.

73 In this study, we tested the 17,817 CDS markers that we have developed in a previous  
74 study (Song, et al. 2017), and screened for the best markers for all major ray-finned fish  
75 clades. We chose the best markers according to results of pilot experiments and refined the  
76 baits design to improve evenness of reads coverage in different loci. Finally, we tested  
77 phylogenetic usefulness of selected markers in ray-finned fishes at both high taxonomic  
78 level and species level. Our goal is to provide a set of common exon markers for gene  
79 capture and phylogenomic studies in the ray-finned fishes.

## 80 **New Approaches**

### 81 *Testing the targeted gene markers in different groups of ray-finned fishes*

82 We tested the single-copy CDS markers identified from our previous study (Song, et al.  
83 2017). The markers were identified through comparing eight fish genomes (Fig. 1A) using  
84 a bioinformatics tool, EvolMarkers (Li, et al. 2012) (supplementary materials Fig. S1).  
85 Baits designing steps can be found in detailed materials and methods of supplementary  
86 materials. Thousands of the candidate CDS markers were tested empirically in 83  
87 actinopterygian species (99 individuals, 33 families of 11 orders), covering major clades of  
88 ray-finned fishes (supplementary materials Table S1). The species captured were part of  
89 five different research projects conducted in the authors' laboratory, including works on  
90 basal actinopterygians (Basal), acipenseriforms (Acipen), ostarioclupeomorphs (Ostario),  
91 gobioids (Goby) and sinipercids (Sini) (supplementary materials Fig. S2).

### 92 *Selecting the best markers and refining the baits design based on gene capture results*

93 Based on results of the pilot experiments, target gene markers and baits were evaluated and  
94 redesigned to improve their efficacy. There were two major considerations: 1) to select for  
95 markers which resulted in less missing data and were phylogenetically decisive, and 2) to  
96 identify regions with extraordinarily high read depth and mask those regions for future baits  
97 design (Fig. 2). The assembled sequences from different projects were merged (*merge.pl*).  
98 Taxa had more than 3,000 genes captured were kept (*select.pl*). Subsequently, a Perl script  
99 *deci.pl* was used to pick phylogenetically decisive loci. Phylogenetic decisiveness means  
100 that the data sets should contain all taxa whose relationships are addressed (Dell'Ampio, et  
101 al. 2014). In our case, the decisive taxonomic groups included eight major clades of the  
102 ray-finned fishes: Acipenseriformes, Lepisosteiformes, Elopomorpha, Osteoglossomorpha,  
103 Ostarioclupeomorpha, Gobiomorpharia, Ovalentariae and Percomorpharia. The  
104 Polypteridae was excluded in bait design, because both species of the polypterids sampled  
105 had less than 3000 targets captured.

106 From our previous experience, we found that partial regions of some target loci had

107 extraordinarily high number of reads mapped, which consumed a large proportion of the  
108 total data collected. Those regions escaped RepeatMasker (Smit, et al. 1996-2004 )  
109 checking in original baits design, wasted a lot of sequencing reads and are better to be  
110 excluded from future baits design. To find those problematic regions, the selected decisive  
111 data were parsed to different files by species name (*parsefast.pl*), then, the raw reads of  
112 each species were mapped to the assembled reference sequences of each species using *BWA*  
113 (Li and Durbin 2009). The reads depth data were extracted from the mapping results using  
114 *SAMtools* (Li, et al. 2009) and a custom Perl script (*mapdepth.pl*). Regions with  
115 extraordinary high read depth, i.e., 100 times than adjacent regions were identified  
116 (*pickbaits.pl*), and manually checked and masked for future baits design. All custom Perl  
117 scripts can be found at <http://www.lmse.org/markersandtools.html>.

#### 118 *Testing phylogenetic usefulness of the markers selected and efficacy of the new baits*

119 A phylogeny of 16 orders of ray-finned fishes, including 10 species with gene captured data  
120 and 7 species with sequence data extracted from genomes were reconstructed. A phylogeny  
121 of four species of freshwater sleepers (*Odontobutis*, Gobiiformes) also was reconstructed  
122 based on gene capture data of the chosen markers, including five individuals of each  
123 species of *Odontobutis sinensis*, *O. potamophila* and *O. yaluensis* and one individual of *O.*  
124 *haifengensis*. Two individuals of *Perccottus glenii* were used as outgroup. Therefore, the  
125 phylogenetic usefulness of the chosen markers was evaluated in reconstructing phylogenies  
126 of ray-finned fishes at both high and low taxonomic levels. Additionally, we extracted  
127 single nucleotide polymorphisms (SNPs) from captured data of the *Odontobutis*, and  
128 visualized inter- and intra-specific genetic variation among individuals of the four  
129 *Odontobutis* species using the principal component analysis (PCA).

130 The new baits refined based on empirical data were compared with the baits designed a  
131 priori. Reads depth and evenness of reads coverage were summarized from the gene capture  
132 data. The comparison was done on results of capturing a goby species (*Rhinogobius*  
133 *giurinus*). Finally, to help researchers to design baits using reference species that are closer

134 to their organism of interested than the eight model fishes that we used, we developed a  
135 pipeline of retrieving sequences of the target loci from user provided genomes  
136 (supplementary materials file 3).

## 137 **Results**

### 138 *Single-copy protein coding markers for ray-finned fishes*

139 The number of loci captured ranged from 435 to 11,534 in different samples. All but four  
140 samples had more than three thousand loci captured (supplementary materials Fig. S2). The  
141 samples did the worst in gene capture experiment included two polypteriforms  
142 (*Erpetoichthys calabaricus* and *Polypterus endlicher*), one sturgeon (*Acipenser ruthenus*)  
143 and the Waigeo barramundi (*Psammoperca waigiensis*). After combining the data from all  
144 five projects, excluding taxa with less than 3,000 loci captured and selecting for  
145 phylogenetic decisive loci, we obtained 4,434 CDS markers of 2,261 genes. The information  
146 of the target loci and sequences of the eight model fish species can be found at  
147 <http://www.lmse.org/markersandtools.html>.

### 148 *Phylogenetic usefulness of selected markers*

149 The average length of coding region of the chosen markers was 236 bp (94 bp to 4,718bp).  
150 GC content ranged from 37% to 69% with an average of 55%. Average pairwise distance  
151 (p-dist) among the 17 species varied from 0.06 to 0.50 substitutions per site, with an overall  
152 average of 0.19. Average consistency index (CI) was 0.60 (0.43-0.93), and average  
153 retention index (RI) was 0.52 (0.47-0.62) (supplementary materials Fig. S3). Maximum  
154 likelihood (ML) analyses concatenating 4,434 loci resulted in a well-resolved tree of major  
155 ray-finned fish clades, and all nodes had 100 bootstrap support values (Fig. 1). The  
156 resulting phylogenetic tree is consistent with recent studies (Betancur, et al. 2013; Faircloth,  
157 et al. 2013), except that the Elopomorpha and the Osteoglossomorpha were found sister to  
158 each other.

159 There were 4,296 of 4,434 loci captured at least in one *Odontobutis* sample. A total of

160 1,630 loci were captured in all samples. The average length of target regions was 265 bp  
161 (120 bp to 5,637 bp). The average length of captured non-coding flank region was 487 bp.  
162 A concatenated ML tree was reconstructed for the four Chinese *Odontobutis* species with *P.*  
163 *glenii* as the outgroup. The species level phylogeny was well resolved with 100 bootstrap  
164 support values for each node. *Odontobutis sinensis* was sister to the rest of *Odontobutis*  
165 species. *Odontobutis yaluensis* was grouped with *O. potamophila* and *O. haifengensis* was  
166 placed as sister to them. Species tree is consistent with ML tree with a normalized quartet  
167 score 0.64 (supplementary materials Fig. S4). We extracted 36,440 single nucleotide  
168 polymorphisms (SNPs) sites from target regions (35 SNPs per kb). In PCA, axis 1 and axis  
169 2 explained 48.42% and 11.21% of the variability respectively. Individuals of *O. yaluensis*  
170 and *O. potamophila* were close to each other, whereas individuals of *O. sinensis* were apart  
171 from them and *O. haifengensis* lied in between (supplementary materials Fig S5).

#### 172 *Gene-capture marker refinement*

173 We examined the results of gene capture experiments using original baits. We found that 26  
174 loci of *Rhinogobius giurinus* had extreme high number of reads mapped. We manually  
175 checked those loci and found that all regions with high reads depth had low complexity. We  
176 masked those regions, redesigned the baits and carried a new round of gene capture  
177 experiment. The gene capture results from new baits had better even coverage among  
178 different loci than the results from the original baits (Fig. 3).

## 179 **Discussions**

### 180 *Exon capture*

181 Protein-coding sequences are easy to align and molecular evolution of protein sequence is  
182 better studied than non-coding flank regions, whose variation tend to increase when further  
183 apart from the conserved core region (Faircloth, et al. 2012). Our experiments showed that  
184 the markers selected and the baits designed were effective in studying phylogenetic  
185 relationship of major groups of the ray-finned fishes, and closely related species as well.



186 We notice that our exon capture protocol also produced data from flanking non-coding  
187 regions with an average length of 487 bp. We did not analyze the sequence data from the  
188 flanking region, because the non-coding flanking regions of many loci could not be aligned.  
189 Further investigation on how to process and utilize the data of flanking regions for studies  
190 at inter- and intraspecific level should be carried out.

#### 191 *A posteriori marker design*

192 The simple repeats in the markers were detected and masked using RepeatMasker by the  
193 manufacturer, MYcroarray (Ann Arbor, Michigan) before synthesizing the baits. However,  
194 repeats with some variations or complex repeats could not be detected with RepeatMasker,  
195 thus resulted in a high read depth in some regions (Fig. 3B). Extreme high read depth  
196 suggests that many reads were not from the target regions, which could cause problem in  
197 subsequent read assembly and waste sequencing resource. Based on the sequencing results,  
198 we masked these unusual regions in the following baits refinement in gobies, which has  
199 shown more even coverage for the targeted loci (Fig. 3B). If a pilot study is planned before  
200 a large-scale experiment, we recommend applying our method to refine baits design to  
201 improve the efficacy of baits.

#### 202 *Orthology checking and data filtering*

203 Problem of mistakenly using paralogous genes for phylogenetic reconstruction is  
204 exacerbated with phylogenomic data, and currently there is no ideal method to validate  
205 orthology of loci assembled from NGS data (McCormack, et al. 2013; Chakrabarty, et al.  
206 2017). The targeted loci we selected for are “single-copy” (Li, et al. 2012), which may have  
207 less chance to be paralogous than members of gene families, (Li, et al. 2007). In addition,  
208 we performed a “re-blast” step in data processing pipeline to identify and exclude potential  
209 paralogs (Yuan, et al. 2016). Nonetheless, both method cannot guarantee orthology of  
210 targeted sequences due to the third round of whole-genome duplication event in teleost and  
211 slow and steady loss of some paired genes in the subsequent 250 My (Inoue, et al. 2015).  
212 Tree based methods, such as filtering the loci a posteriori based on known monophyly of

213 taxa could be used to alleviate the problem of paralogy.

## 214 **Materials and Methods**

215 For detailed materials and methods, see supplementary materials file 4.

## 216 **Supplementary Materials**

217 Supplementary File 1: Figures S1 - S5.

218 Supplementary File 2: Tables S1.

219 Supplementary File 3: A user-friendly pipeline to retrieve target sequences of the 4,434 loci  
220 from new genome sequences or transcriptomes.

221 Supplementary File 4: Detailed materials and methods.

## 222 **Acknowledgements**

223 This work was supported by the Innovation Program of Shanghai Municipal Education  
224 Commission and the Program for Professor of Special Appointment (Eastern Scholar) at  
225 Shanghai Institutions of Higher Learning. We would like to thank Shanghai Oceanus  
226 Supercomputing Center (SOSC) for providing computational resource.

## 227 **References**

- 228 Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA,  
229 Johnson EA. 2008. Rapid SNP discovery and genetic mapping using sequenced RAD markers.  
230 *PloS one* 3:e3376.
- 231 Betancur RR, Broughton RE, Wiley EO, Carpenter K, Lopez JA, Li C, Holcroft NI, Arcila D,  
232 Sanciangco M, Cureton Ii JC, et al. 2013. The tree of life and a new classification of bony fishes.  
233 *PLoS currents* 5.
- 234 Bi K, Vanderpool D, Singhal S, Linderoth T, Moritz C, Good JM. 2012. Transcriptome-based exon  
235 capture enables highly cost-effective comparative genomic data collection at moderate  
236 evolutionary scales. *BMC Genomics* 13:403.
- 237 Bragg JG, Potter S, Bi K, Moritz C. 2016. Exon capture phylogenomics: efficacy across scales of  
238 divergence. *Molecular ecology resources* 16:1059-1068.
- 239 Chakrabarty P, Faircloth BC, Alda F, Ludt WB, McMahan CD, Near TJ, Dornburg A, Albert JS,  
240 Arroyave J, Stiassny ML, et al. 2017. Phylogenomic Systematics of Ostariophysan fishes:  
241 Ultraconserved Elements Support the Surprising Non-monophyly of Characiformes. *Systematic*  
242 *Biology*.

- 243 Chamala S, Garcia N, Godden GT, Krishnakumar V, Jordon-Thaden IE, De Smet R, Barbazuk WB,  
244 Soltis DE, Soltis PS. 2015. Markerminer 1.0: A New Application for Phylogenetic Marker  
245 Development Using Angiosperm Transcriptomes. *Applications in Plant Sciences* 3:1400115.
- 246 Collins RA, Hrbek T. 2015. An in silico comparison of reduced-representation and  
247 sequence-capture protocols for phylogenomics. bioRxiv preprint first posted online Nov. 21,  
248 2015; doi: <http://dx.doi.org/10.1101/032565>.
- 249 Davey JW, Blaxter ML. 2010. RADSeq: next-generation population genetics. *Briefings in*  
250 *functional genomics* 9:416-423.
- 251 Dell'Ampio E, Meusemann K, Szucsich NU, Peters RS, Meyer B, Borner J, Petersen M, Aberer AJ,  
252 Stamatakis A, Walz MG, et al. 2014. Decisive data sets in phylogenomics: lessons from studies  
253 on the phylogenetic relationships of primarily wingless insects. *Molecular biology and*  
254 *evolution* 31:239-249.
- 255 Emerson KJ, Merz CR, Catchen JM, Hohenlohe PA, Cresko WA, Bradshaw WE, Holzapfel CM.  
256 2010. Resolving postglacial phylogeography using high-throughput sequencing. *Proceedings of*  
257 *the National Academy of Sciences of the United States of America* 107:16196-16200.
- 258 Faircloth BC, McCormack JE, Crawford NG, Harvey MG, Brumfield RT, Glenn TC. 2012.  
259 Ultraconserved Elements Anchor Thousands of Genetic Markers Spanning Multiple  
260 Evolutionary Timescales. *Systematic Biology* 61:717-726.
- 261 Faircloth BC, Sorenson L, Santini F, Alfaro ME. 2013. A Phylogenomic Perspective on the  
262 Radiation of Ray-Finned Fishes Based upon Targeted Sequencing of Ultraconserved Elements  
263 (UCEs). *PloS one* 8:e65923.
- 264 Hedtke SM, Morgan MJ, Cannatella DC, Hillis DM. 2013. Targeted enrichment: maximizing  
265 orthologous gene comparisons across deep evolutionary time. *PloS one* 8:e67908.
- 266 Hugall AF, O'Hara TD, Hunjan S, Nilsen R, Moussalli A. 2016. An Exon-Capture System for the  
267 Entire Class Ophiuroidea. *Molecular biology and evolution* 33:281-294.
- 268 Ilves KL, Lopez-Fernandez H. 2014. A targeted next-generation sequencing toolkit for exon-based  
269 cichlid phylogenomics. *Molecular ecology resources* 14:802-811.
- 270 Inoue J, Sato Y, Sinclair R, Tsukamoto K, Nishida M. 2015. Rapid genome reshaping by  
271 multiple-gene loss after whole-genome duplication in teleost fish suggested by mathematical  
272 modeling. *Proceedings of the National Academy of Sciences of the United States of America*  
273 112:14918-14923.
- 274 Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SY, Faircloth BC, Nabholz B, Howard JT,  
275 et al. 2014. Whole-genome analyses resolve early branches in the tree of life of modern birds.  
276 *Science* 346:1320-1331.
- 277 Lemmon AR, Emme SA, Lemmon EM. 2012. Anchored hybrid enrichment for massively  
278 high-throughput phylogenomics. *Syst Biol* 61:727-744.
- 279 Li C, Hofreiter M, Straube N, Corrigan S, Naylor GJ. 2013. Capturing protein-coding genes across  
280 highly divergent species. *BioTechniques* 54:321-326.
- 281 Li C, Ortí G, Zhang G, Lu G. 2007. A practical approach to phylogenomics: the phylogeny of  
282 ray-finned fish (Actinopterygii) as a case study. *BMC Evol. Biol.* 7:44.
- 283 Li C, Riethoven JJ, Naylor GJ. 2012. *EvolMarkers*: a database for mining exon and intron markers

- 284 for evolution, ecology and conservation studies. *Molecular ecology resources* 12:967-971.
- 285 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
- 286 *Bioinformatics* 25:1754-1760.
- 287 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009.
- 288 The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078-2079.
- 289 Mandel JR, Dikow RB, Funk VA, Masalia RR, Staton SE, Kozik A, Michelmore RW, Rieseberg LH,
- 290 Burke JM. 2014. A target enrichment method for gathering phylogenetic information from
- 291 hundreds of loci: An example from the Compositae. *Applications in plant sciences* 2.
- 292 Mayer C, Sann M, Donath A, Meixner M, Podsiadlowski L, Peters RS, Petersen M, Meusemann K,
- 293 Liere K, Wägele JW, et al. 2016. BaitFisher: A Software Package for Multispecies Target DNA
- 294 Enrichment Probe Design. *Molecular biology and evolution* 33:1875-1886.
- 295 McCormack JE, Hird SM, Zellmer AJ, Carstens BC, Brumfield RT. 2013. Applications of
- 296 next-generation sequencing to phylogeography and phylogenetics. *Molecular phylogenetics and*
- 297 *evolution* 66:526-538.
- 298 Nelson JS, Grande TC, Wilson MVH. 2016. *Fishes of the World*: Wiley.
- 299 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. 2012. Double digest RADseq: an
- 300 inexpensive method for de novo SNP discovery and genotyping in model and non-model
- 301 species. *PloS one* 7:e37135.
- 302 Portik DM, Smith LL, Bi K. 2016. An evaluation of transcriptome-based exon capture for frog
- 303 phylogenomics across multiple scales of divergence (Class: Amphibia, Order: Anura).
- 304 *Molecular ecology resources* 16:1069-1083.
- 305 RepeatMasker Open-3.0. [Internet]. 1996-2004 Available from: <http://www.repeatmasker.org>
- 306 Song S, Zhao J, Li C. 2017. Species delimitation and phylogenetic reconstruction of the sinipercids
- 307 (Perciformes: Sinipercidae) based on target enrichment of thousands of nuclear coding
- 308 sequences. *Molecular phylogenetics and evolution* 111:44-55.
- 309 Teasdale LC, Kohler F, Murray KD, O'Hara T, Moussalli A. 2016. Identification and qualification of
- 310 500 nuclear, single-copy, orthologous genes for the Eupulmonata (Gastropoda) using
- 311 transcriptome sequencing and exon capture. *Molecular ecology resources* 16:1107-1123.
- 312 Weitemier K, Straub SC, Cronn RC, Fishbein M, Schmickl R, McDonnell A, Liston A. 2014.
- 313 Hyb-Seq: Combining target enrichment and genome skimming for plant phylogenomics.
- 314 *Applications in plant sciences* 2.
- 315 Yuan H, Jiang J, Jimenez FA, Hoberg EP, Cook JA, Galbreath KE, Li C. 2016. Target gene
- 316 enrichment in the cyclophyllidean cestodes, the most diverse group of tapeworms. *Molecular*
- 317 *ecology resources* 16:1095-1106.

318

319

319 **Figure Captions**

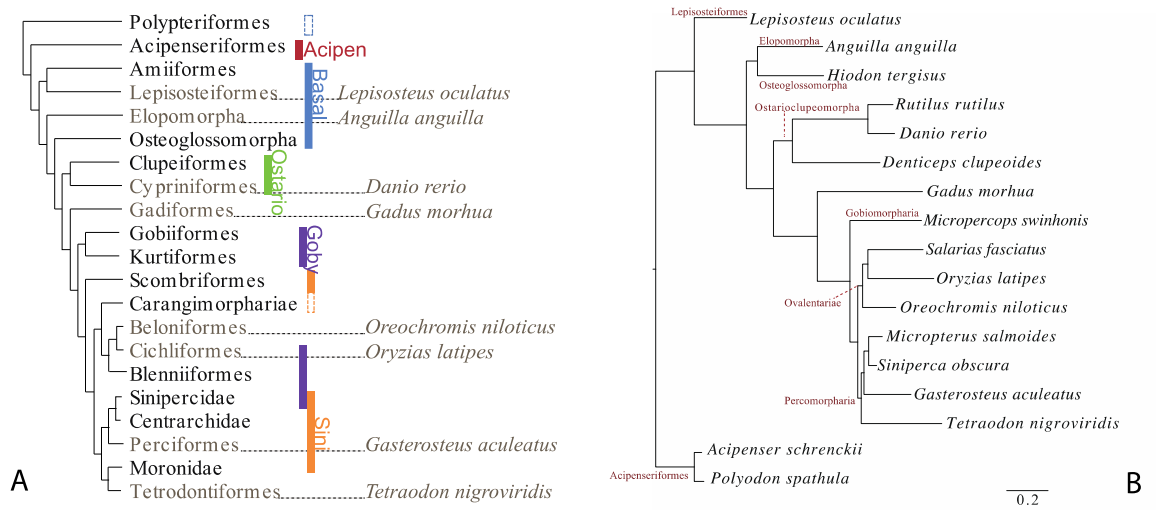
320 **FIG. 1.** A. Phylogenetic relationships among 21 groups of ray-finned fish (Betancur, et al.  
321 2013). Eight species names indicate the fishes with genome sequence available that were  
322 used in finding the target markers. The vertical bars indicate different projects carried in the  
323 author's laboratory. The unfilled vertical bars indicate groups that captured less than 3000  
324 loci. B. Maximum likelihood tree of 16 representative ray-finned fishes based on 4,434  
325 exon loci, all nodes have a 100 bootstrap value.

326 **FIG. 2.** Pipeline of screening for markers with less missing data and better phylogenetic  
327 decisiveness and posterior baits refining. I. Merge data from different project (merge.pl); II.  
328 select loci with less missing data and high phylogenetic decisiveness (select.pl; deci.pl); III.  
329 find and mask region with extraordinary read depth for bait redesign (parsefasta.pl;  
330 runbwa.pl; mapdepth.pl; pickbaits.pl). The posterior baits refining steps are optional when  
331 empirical data from pilot gene capture are available. GCMR stands for gene capture marker  
332 refinement.

333 **FIG. 3.** Comparison on evenness of read coverage between results of gene capture using the  
334 baits designed a priori (A, blue curve) and the baits refined posteriorly (A, orange curve). B  
335 and C are screenshots from visualizing the read depth of the locus  
336 *Danio\_rerio*.20.4037479.4035425 using Tablet v1.16.09.06. In this example, the result  
337 using baits designed a priori (B) is much worse than the result using refined baits (C).

338

338



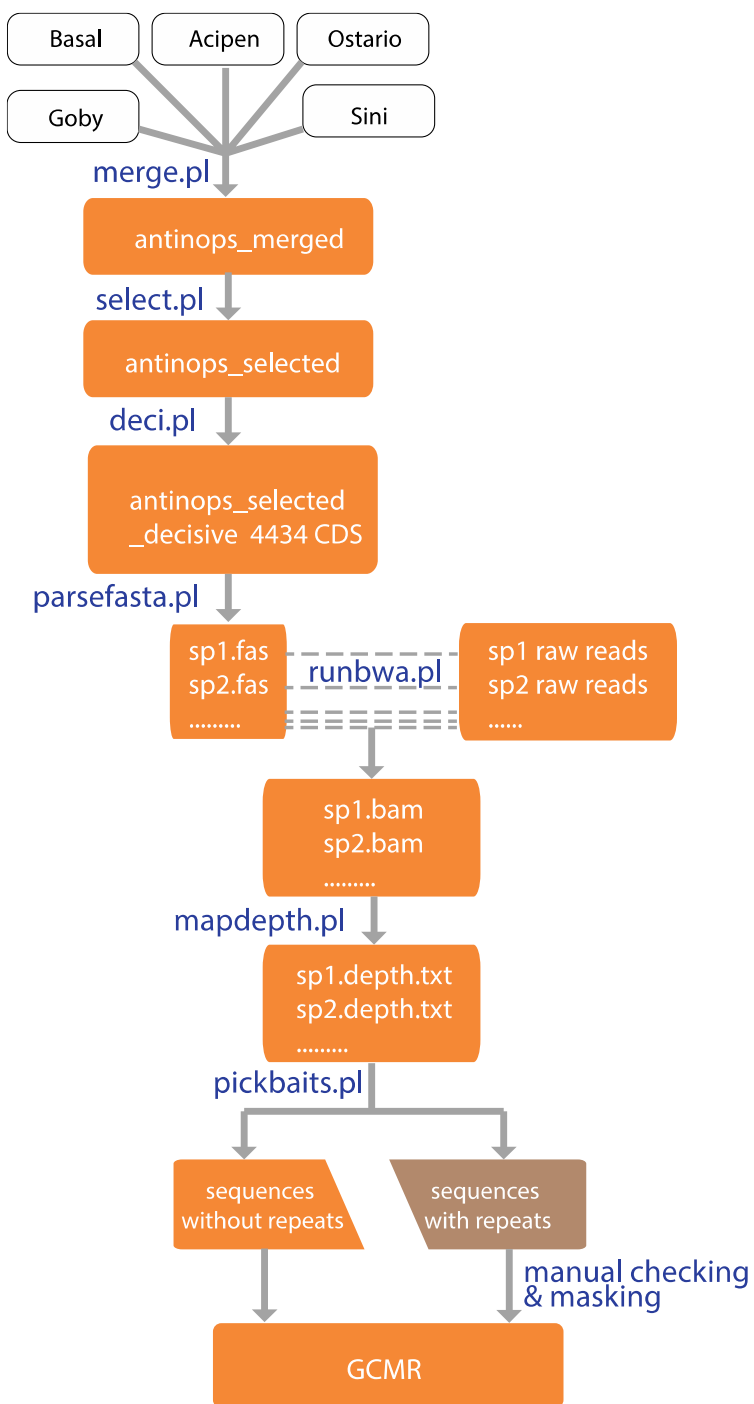
339

340

341 Figure 1

342

342

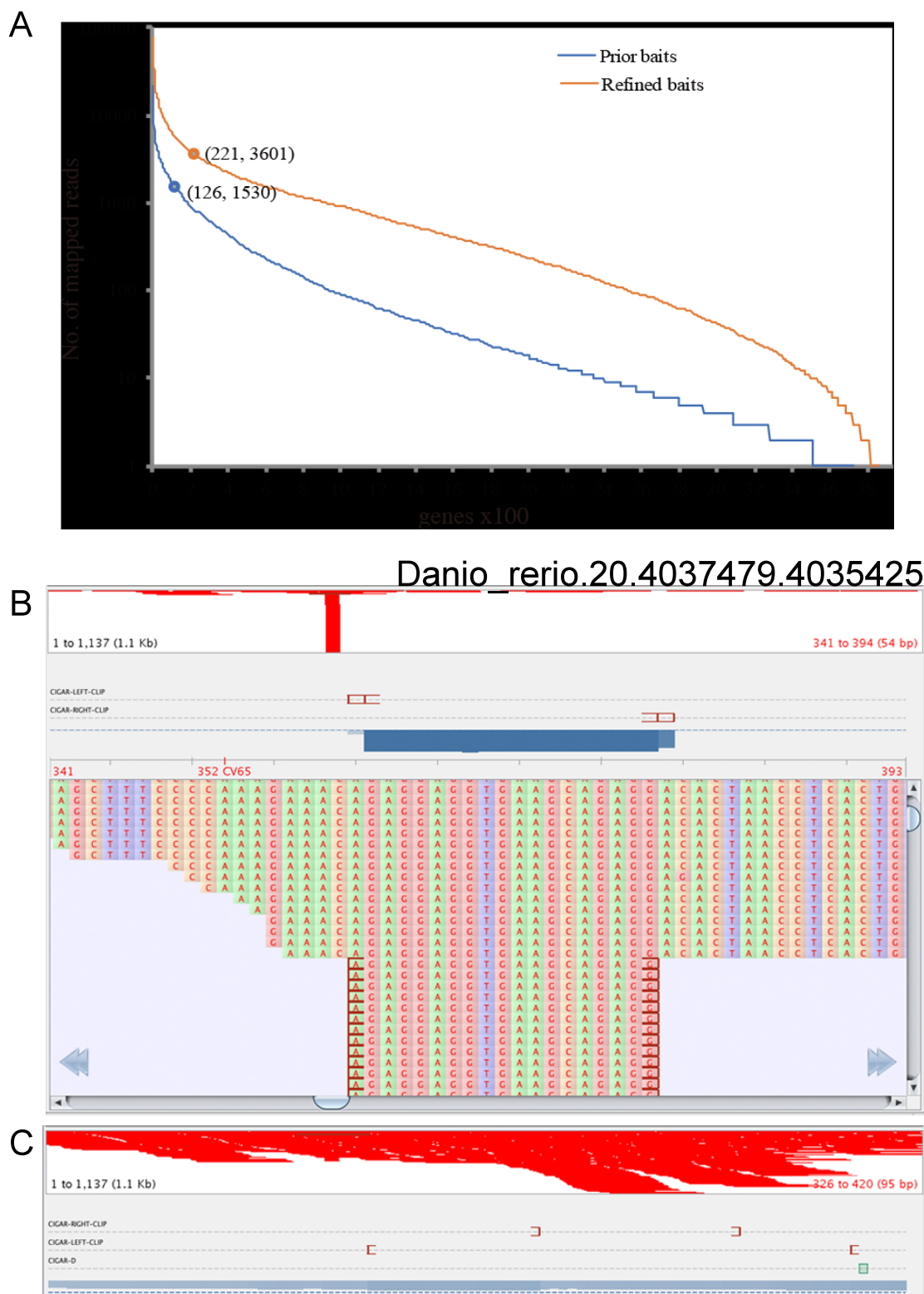


343

344

345 Figure 2

346



346

347

348 Figure 3