

[Click here to view linked References](#)

1 **Whole genome sequencing of Chinese clearhead icefish,**

2 ***Protosalanx hyalocranius***

3

4 Kai Liu<sup>1†</sup>, Dongpo Xu<sup>1†</sup>, Jia Li<sup>2†</sup>, Chao Bian<sup>2†</sup>, Jinrong Duan<sup>1†</sup>, Yanfeng Zhou<sup>1†</sup>,

5 Mingying Zhang<sup>1</sup>, Xinxin You<sup>2</sup>, Yang You<sup>1</sup>, Jieming Chen<sup>2</sup>, Hui Yu<sup>2</sup>, Gangchun Xu<sup>1</sup>,

6 Di-an Fang<sup>1</sup>, Jun Qiang<sup>1</sup>, Shulun Jiang<sup>1</sup>, Jie He<sup>1</sup>, Junmin Xu<sup>2,4,5</sup>, Qiong Shi<sup>2,4,5,6\*</sup>,

7 Zhiyong Zhang<sup>3\*</sup>, Pao Xu<sup>1,5\*</sup>

8

9 <sup>1</sup>Freshwater Fisheries Research Center, Chinese Academy of Fishery Sciences, Wuxi  
10 214081, China

11 <sup>2</sup>Shenzhen Key Lab of Marine Genomics, Guangdong Provincial Key Lab of  
12 Molecular Breeding in Marine Economic Animals, BGI, Shenzhen 518083, China

13 <sup>3</sup>Institute of Oceanology & Marine Fisheries, Jiangsu 226007, China

14 <sup>4</sup>BGI Zhenjiang Institute of Hydrobiology, Zhenjiang 212000, China

15 <sup>5</sup>BGI Research Center for Aquatic Genomics, Chinese Academy of Fishery Sciences,  
16 Shenzhen 518083, China

17 <sup>6</sup>Laboratory of Aquatic Genomics, College of Ecology and Evolution, School of Life  
18 Sciences, Sun Yat-Sen University, Guangzhou 510275, China

19

20 † Equal contributors

21 \* Correspondence: xup@ffrc.cn (PX); shiqiong@genomics.cn (QS);

22 13906292412@139.com (ZZ)

23

24 Email addresses: liuk@ffrc.cn (KL); xudp@ffrc.cn(DX); lijial@genomics.cn (JL);

25 bianchao@genomics.cn (CB); duanjr@ffrc.cn(JD); zhoyuf@ffrc.cn(YZ);

60  
61  
62  
63  
64  
65

1 26 zhangmy@ffrc.cn(MZ); youxinxin@genomics.cn (XY); youy@ffrc.cn (YY);  
2  
3 27 chenjieming@genomics.cn (JC); yuhui@genomics.cn (HY); xugc@ffrc.cn(GX);  
4  
5 28 fangda@ffrc.cn(DF); qiangj@ffrc.cn(JQ); 420219380@qq.com(SJ); hej@ffrc.cn(JH);  
6  
7 29 xujunmin@genomics.cn (JX); shiqiong@genomics.cn (QS); 13906292412@139.com  
8  
9 30 (ZZ); xup@ffrc.cn (PX)

10  
11 31

12  
13 32

### 14 15 33 **Abstract**

16  
17 34 **Background:** Chinese clearhead icefish, *Protosalanx hyalocranius*, is a  
18  
19 35 representative icefish species with economic importance and special appearance. Due  
20  
21 36 to its great economic values in China, the fish was introduced into Lake Dianchi and  
22  
23 37 several other lakes from the Lake Taihu half a century ago. Similar to the  
24  
25 38 *Sinocyclocheilus* cavefish, the clearhead icefish has certain cavefish-like traits, such  
26  
27 39 as transparent body and nearly scaleless skin. Here, we provide the whole genome  
28  
29 40 sequence of this surface-dwelling fish and generated a draft genome assembly, aiming  
30  
31 41 at exploring molecular mechanisms for the biological interests.

32  
33  
34 42 **Findings:** A total of 252.1 gigabases (Gb) of raw reads were sequenced. Subsequently,  
35  
36 43 a novel draft genome assembly was generated, with the scaffold N50 reaching 1.163  
37  
38 44 Mb. The genome completeness was estimated to be 98.39% by using the CEGMA  
39  
40 45 evaluation. Finally, we annotated 19,884 protein-coding genes and observed that  
41  
42 46 repeat sequences account for 24.43% of the genome assembly.

43  
44 47 **Conclusion:** We report the first draft genome of the Chinese clearhead icefish. The  
45  
46 48 genome assembly will provide a solid foundation for further molecular breeding and  
47  
48 49 germplasm resource protection in Chinese clearhead icefish, as well as other icefishes.  
49  
50 50 It is also a valuable genetic resource for revealing the molecular mechanisms for the  
51  
52 51 cavefish-like characters.

53  
54 52 **Keywords:** Clearhead icefish; *Protosalanx hyalocranius*; Whole genome  
55  
56 53 sequencing; Genome assembly; Gene prediction; Repetitive sequences  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 55 **Data description**

### 56 ***Background***

57 Icefishes (Osmeriformes, Salangidae) are widely distributed in freshwater, coastal and  
58 estuarine habitats in East Asian countries [1-3]. Chinese clearhead icefish  
59 (*Protosalanx hyalocranius*), a diadromous fish, mainly inhabits in coastal areas and  
60 adjacent freshwaters [4-6]. As an economically important fish in China, the clearhead  
61 icefish was widely introduced into some lakes from the original Lake Taihu half a  
62 century ago, and it has developed a resident life history in these water areas [2, 7, 8].  
63 Because of its transparent body and nearly scaleless skin, similar to the  
64 *Sinocyclocheilus* cavefishes [9], we are very interested in this surface-dwelling fish  
65 and are performing comparative genomics studies to explore the mechanisms for  
66 these biological phenotypes. However, with the rapid development of the Chinese  
67 economy in recent decades, population size of the clearhead icefish has been seriously  
68 declining because of overfishing, construction of water conservancy facilities and  
69 water pollution in the ecological systems [10]. To maintain its sustainable  
70 development in China, here we performed whole genome sequencing of Chinese  
71 clearhead icefish to support its biological and economic importance.

72

### 73 ***Sample and Sequencing***

74 In this study, we applied Illumina whole genome sequencing (WGS) strategy to  
75 sequence the genome of Chinese clearhead icefish (NCBI Taxonomy ID: 418454;  
76 Fishbase ID: 12236). Genomic DNA was isolated from the muscle tissue of an  
77 individual collected from the Lake Taihu of Jiangsu Province in China. We  
78 constructed seven paired-end libraries with three short-insert libraries (250, 500 and  
79 800 bp) and four long-insert libraries (2, 5, 10 and 20 kb) using the standard protocol  
80 provided by Illumina (San Diego, USA). Subsequent paired-end sequencing was

1 81 performed by the Illumina HiSeq 2000 platform for each library. Finally, we obtained  
2 82 252.1 Gb of raw reads for further analysis.  
3

4 83

### 5 84 *Genome size estimation and genome assembly*

6  
7  
8  
9 85 The SOAPfilter v2.2 software [11] with optimized parameters (-y -p -g 1 -o clean -M  
10 86 2 -f 0) was utilized to remove low-quality raw reads (including reads with 10 or more  
11 87 Ns and low-quality bases) and PCR-replicates as well as adaptor sequences. In total,  
12  
13 88 we obtained 169.0 Gb of clean reads. Subsequently, we estimated the genome size  
14  
15 89 based on the 17-mer depth frequency distribution method [12]. A 17-mer represents  
16  
17 90 an artificial division with 17-bp length nucleotide segment of sequencing reads,  
18  
19 91 therefore, a raw sequence read with a total length of L bp contains (L-17+1) 17-mers.  
20  
21 92 The genome size was estimated with the following formula:  $G = K\_num / K\_depth$ , in  
22  
23 93 which G is the genome size, K\_num is the total number of 17-mer, and K\_depth is the  
24  
25 94 highest frequency of 17-mer analysis. In our current study, the K\_num was  
26  
27 95 10,500,000,000 and the K\_depth was 20. Hence, we estimated that the genome size of  
28  
29 96 Chinese clearhead icefish is 525 Mb.  
30  
31  
32  
33

34 97 The filtered reads were assembled using SOAPdenovo2 v2.04.4 software [13] with  
35 98 optimized parameters (pregraph -K 79 -d 1; contig -M 1; scaff -F -b 1.5 -p 16) to  
36 99 generate contigs and original scaffolds. The gaps were filled using GapCloser v1.12  
37 100 software [14] with default parameters and -p set to 25. Finally, we generated a draft  
38 101 genome assembly of 536 Mb, with the scaffold N50 reaching 1.163 Mb (Table 1).  
39 102 The completeness of our assembly was evaluated by using CEGMA [15] and BUSCO  
40 103 [16]. The CEGMA program (Core Eukaryotic Genes Mapping Approach; version 2.4)  
41 104 assessment with 248 conserved Core Eukaryotic Genes (CEGs) was performed for  
42 105 evaluation of the gene space completeness. Our results revealed that the assembled  
43 106 genome had a CEGMA completeness score at 90.32% and 98.39%, which was  
44 107 calculated from the complete gene set and the partial gene set, respectively.  
45 108 Meanwhile, we used the representative metazoa gene set [17], which contains 843  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

109 single-copy genes that are widely present in metazoan, as a reference. The assessment  
110 demonstrated that the BUSCO values is 89%, containing [D: 10%], F: 7.7%, M: 2.9%,  
111 n: 843 (C: complete [D: duplicated], F: fragmented, M: missed, n: genes). These data  
112 from CEGMA and BUSCO indicate that the assembled genome covered majority of  
113 the gene space.

### 115 ***Repeat annotation***

116 Firstly, a *de novo* repeat library was constructed by the RepeatModeller v1.05 [18]  
117 and LTR\_FINDER.x86\_64-1.0.6 [10] with default parameters. Then, the assembled  
118 genome sequences were aligned against the RepBase v21.01 [19] and the *de novo*  
119 repeat libraries to recognize the known and novel transposable elements ( TEs ) using  
120 the RepeatMasker v4.06 [20]. Meantime, the Tandem Repeat Finder v4.07 [21] with  
121 parameters “Match=2, Mismatch=7, Delta=7, PM=80, PI=10, Minscore=50, and  
122 MaxPeriod=2000” was utilized for annotation of tandem repeats. Furthermore, the  
123 RepeatProteinMask software v4.0.6 [20] was used to predict TE relevant proteins in  
124 our genome assembly. Finally, we observed that the repeat sequences account for  
125 24.43% of the assembled genome (Table 1), and the *de novo* annotation method  
126 predicted the most abundant repeat sequence among the four methods (Table 2).

### 128 ***Genome Annotation***

129 In brief, we utilized two different methods to predict total gene set of the clearhead  
130 icefish.

131 **1) *de novo* annotation.** The AUGUSTUS v2.5 [22] and GENSCAN v1.0 [23] were  
132 executed to *ab initio* predict genes within the assembled genome, with the repetitive  
133 sequences masked as “N” in order to discard pseudo gene prediction. Those  
134 low-quality genes with short length (<150 bp), premature termination or  
135 frame-shifting were removed.

1 136 **2) Homology annotation.** We aligned the protein sequences from six published  
2 137 genomes, including *Danio rerio* [24], *Oryzias latipes* [25], *Takifugu rubripes* [26],  
3  
4 138 *Tetraodon nigroviridis* [27], *Esox lucius* [28] and *Gasterosteus aculeatus* [29], against  
5  
6 139 our assembly to predict homology-based genes. The potential homology-based genes  
7  
8 140 were searched by TblastN [30] with an e-value of  $10^{-5}$ . The TblastN results were then  
9  
10 141 processed by SOLAR (Sorting Out Local Alignment Result [31]) to obtain the best hit  
11  
12 142 of each alignment. Subsequently, GeneWise v2.2.0 [32] was performed to detect the  
13  
14 143 possible gene structure for the best hit of each alignment. The low-quality genes were  
15  
16 144 also removed as described in the above-mentioned *de novo* annotation.  
17

18  
19  
20 145 **3) Integration of annotation results.** We employed the GLEAN [33] to generate a  
21  
22 146 non-redundant and comprehensive gene set. Finally, the best hit of each protein was  
23  
24 147 obtained through all protein sequences from the GLEAN results aligned to the  
25  
26 148 databases of the SwissProt and TrEMBL [34] (Uniprot release 2011.06) by BlastP  
27  
28 149 with an e-value of  $10^{-5}$ . Overall, we generated a final gene set with 19,884 genes for  
29  
30 150 the Chinese clearhead icefish.  
31

32  
33  
34 151 CEGMA was performed again to evaluate the coverage rate between KOG  
35  
36 152 (EuKaryotic Orthologous Groups) genes predicted by CEGMA and the predicted total  
37  
38 153 gene set. It demonstrates that the predicted gene set mapped 96.4% of the KOGs.  
39  
40 154 Simultaneously, the BUSCO was implemented again to assess completeness of the  
41  
42 155 predicted gene set. The BUSCO values were calculated as follows: C: 79% [D: 16%],  
43  
44 156 F: 9.8%, M: 10%, n: 843 (C: complete [D: duplicated], F: fragmented, M: missed, n:  
45  
46 157 genes). The assessment values from both CEGMA and BUSCO proved high accuracy  
47  
48 158 of the annotation.  
49

50  
51 159 **4) Function annotation.** The predicted protein sequences of the clearhead icefish  
52  
53 160 were aligned against several public databases (Pfam [35], PRINTS [36], ProDom [37]  
54  
55 161 and SMART [38]) for detection of functional motifs and domains. Finally, we found  
56  
57 162 that 96.2% of the predicted total gene set had been annotated with at least one  
58  
59  
60  
61  
62  
63  
64  
65

1 163 functional assignment from other public databases (Swiss-Prot [39], Interpro [40],  
2 164 TrEMBL [41] and KEGG [42]).

3  
4  
5 165

### 6 7 166 ***Genome evolution***

8  
9 167 We performed phylogenomic analyses with orthologues from representative species  
10  
11 168 for each clade. We used the Ensembl BioMart ([www.ensembl.org/biomart](http://www.ensembl.org/biomart); Ensembl  
12  
13 169 version 76) to extract orthologues for zebrafish [24], fugu [26], stickleback [29],  
14  
15 170 medaka [25] and spotted gar [43]. This generated orthologue dataset from six species  
16  
17 171 was filtered out to retain only one-to-one orthologues. Meanwhile, a new Asian  
18  
19 172 arowana gene set stem from our recent work [44]. In order to extrapolate the Biomart  
20  
21 173 orthologues to the arowana and clearhead icefish gene sets, we used zebrafish as the  
22  
23 174 reference. We ran InParanoid [45] for the three species pairs (zebrafish-arowana and  
24  
25 175 zebrafish-clearhead icefish) at default settings (i.e., minimum 50% alignment span,  
26  
27 176 minimum 25% alignment coverage, minimum BLASTP score of 40 bits, minimum  
28  
29 177 inparalog confidence level of 0.05). By comparing the three InParanoid outputs, we  
30  
31 178 narrowed down the list of one-to-one orthologues, presented in all the seven species,  
32  
33 179 to 454 genes. Subsequently, multiple alignments were performed on proteins of each  
34  
35 180 selected family by MUSCLE (version 3.8.31) [46] and protein alignments were  
36  
37 181 converted to their corresponding CDS alignments using an in-house perl script. All  
38  
39 182 the translated CDS sequences were linked into one “supergene” for each species.  
40  
41 183 Non-degenerated sites extracted from the supergenes were then joined into new  
42  
43 184 sequence of each species to construct a phylogenetic tree (Figure 1) using MrBayes  
44  
45 185 [47] (Version 3.2, GTR+gamma model). Our phylogenetic data demonstrate the close  
46  
47 186 relationship between the clearhead icefish and zebrafish & medaka (Figure 1).  
48  
49  
50  
51

52 187

### 53 54 188 **Synten blocks and genome duplication**

55  
56 189 Genomic homology between the clearhead icefish and medaka was examined using  
57  
58 190 i-ADHoRe 3.0 [48] using the following settings: alignment method gg2, gap size 30,  
59  
60  
61  
62  
63  
64  
65

1 191 tandem gap 30, cluster gap 35, q value 0.85, prob cutoff 0.01, anchor points 5 and  
2  
3 192 multiple hypothesis correction FDR. The output was processed by the pipeline and  
4  
5 193 included in a relational database to which visualization programs can connect and on  
6  
7 194 which additional statistical analysis can be performed. For synteny detection, the  
8  
9 195 cloud mode was enabled (cluster\_type = cloud) and appropriate settings were selected  
10  
11 196 as follows: cloud\_gap\_size 20, cloud\_cluster\_gap 20, cloud\_filter\_method binomial,  
12  
13 197 prob cutoff 0.01, anchor points 5, multiple hypothesis correction FDR and  
14  
15 198 level\_2\_only true. Finally, we identified 660 synteny blocks containing 6,156 genes  
16  
17 199 between the clearhead icefish and medaka.  
18  
19 200 Subsequently, Protein sequences of homologous gene pairs in the identified syntenic  
20  
21 201 regions were aligned using MUSCLE [46], and the protein alignments were then  
22  
23 202 converted to the CDS alignments. Finally, four-fold degenerative third-codon  
24  
25 203 transversion (4DTV) values were calculated on these CDS alignments and corrected  
26  
27 204 using the HKY model in the PAML package [49]. These data indicate that the  
28  
29 205 clearhead icefish also experienced the teleost-specific whole genome duplication  
30  
31 206 (WGD), and it appeared more recently than medaka (Figure 2).  
32  
33  
34  
35  
36

## 37 208 **Conclusion**

38  
39 209 We generated a draft genome assembly of the Chinese clearhead icefish. The novel  
40  
41 210 genome data were deposited in publicly accessible repositories to promote further  
42  
43 211 biological research, molecular breeding and resource protection of this representative  
44  
45 212 and valuable icefish.  
46  
47

## 48 213 49 50 214 **Availability of supporting data**

51  
52 215 Supporting data are available in the GigaDB database, and the raw genome sequences  
53  
54 216 are deposited in the SRA under the bioproject number PRJNA328051.  
55  
56

## 57 217 58 59 218 **Competing interests**

60  
61  
62  
63  
64  
65



1 219 The authors declare that they have no competing interests.

2 220

3  
4  
5 221 **Funding**

6  
7 222 This study was supported by a grant from fish investigation in Taihu Lake (No.  
8  
9 223 TH2016WT007), National Infrastructure of Fishery Germplasm Resources (No.  
10  
11 224 2016DKA30470), Basic Research Funds from Freshwater Fisheries Research Center  
12  
13 225 (No. 2013JBFM07), Special Project on the Integration of Industry, Education and  
14  
15 226 Research of Guangdong Province (No. 2013B090800017), Shenzhen Special Program  
16  
17 227 for Future Industrial Development (No. JSGG20141020113728803), and Zhenjiang  
18  
19 228 Leading Talent Program for Innovation and Entrepreneurship.

20 229

21  
22  
23  
24  
25 230 **Author's Contributions**

26  
27 231 KL, PX, QS, DX, JX, CB and ZZ conceived the project. MZ, XY, HY, JC, GX, DF,  
28  
29 232 JQ, SJ and JH collected the samples and extracted the genomic DNA. JL, CB and HY  
30  
31 233 performed the genome assembly and data analysis. JL, CB, QS, KL, XP, KL, YY and  
32  
33 234 ZZ wrote the paper.

34  
35  
36 235

37  
38  
39 236 **References**

- 40  
41  
42 237 1. Wang ZS, Cui Zhang FU: **Biodiversity of Chinese Icefishes (Salangidae)**  
43 238 **and their conserving strategies**. *Chinese Biodiversity* 2002, **10**(4):416-424.  
44  
45 239 2. Zhang J, Li M, Xu M, Takita T, Wei F: **Molecular phylogeny of icefish**  
46 240 **Salangidae based on complete mtDNA cytochrome b sequences, with**  
47 241 **comments on estuarine fish evolution**. *Biological Journal of the Linnean*  
48 242 *Society* 2007, **91**(2):325-340.  
49  
50 243 3. Wang Z, Lu C, Hu H, Xu C, Lei G: **Dynamics of Icefish (Salangidae) Stocks**  
51 244 **in Nanyi Lake, Eastern China: Degradation and Overfishing**. *Journal of*  
52 245 *Freshwater Ecology* 2004, **19**(2):271-278.  
53  
54 246 4. Xia DQ, Cao Y, Ting ting WU, Yang H: **Study on lineages of Protosalanx**  
55 247 **chinensis, Neosalanx taihuensis and N.oligodontis in Taihu Lake with**  
56 248 **RAPD technique**. *Journal of Fisheryences of China* 2000, **7**(01):12-15.  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 249 5. Xia DQ, Cao Y, Ting Ting WU, Yang H: **Genetic Structures of Population**  
250 **of Protosalanx Chinensis, Neosalanx Taihuensis and Neosalanx**  
251 **Oligodontis in Lake Taihu.** *Journal of Fisheries of China* 1999(03):254-260.
  - 252 6. Armani A, Castigliero L, Tinacci L, Gianfaldoni D, Guidi A: **Molecular**  
253 **characterization of icefish, (S alangidae family), using direct sequencing**  
254 **of mitochondrial cytochrome b gene.** *Food Control* 2011, **22**(6):888-895.
  - 255 7. Wang Z, Lu C, Hu H, Zhou Y, Xu C, Lei G: **Freshwater icefishes**  
256 **(Salangidae) in the Yangtze River basin of China: Spatial distribution**  
257 **patterns and environmental determinants.** *Environmental Biology of Fishes*  
258 2005, **73**(3):253-262.
  - 259 8. Ye S, Yang J, Liu H, Oshima Y: **Use of elemental fingerprint analysis to**  
260 **identify localities of collection for the large icefish protosalanx chinensis**  
261 **in Taihu Lake, China.** *Journal of the Faculty of Agriculture, Kyushu*  
262 *University* 2011, **56**(1):41-45.
  - 263 9. Yang J, Chen X, Jie B, Fang D, Ying Q, Jiang W, Hui Y, Chao B, Jiang L, He  
264 S: **The Sinocyclocheilus cavefish genome provides insights into cave**  
265 **adaptation.** *Bmc Biology* 2016, **14**(1):1-13.
  - 266 10. Xu J, Xie P, Zhang M, Zhou Q, Zhang L, Wen Z, Cao T: **Icefish (salangidae)**  
267 **as an indicator of anthropogenic pollution in freshwater systems using**  
268 **nitrogen isotope analysis.** *Bulletin of environmental contamination and*  
269 *toxicology* 2007, **79**(3):323-326.
  - 270 11. Kar HK, Narayan R, Gautam RK, Jain RK, Doda V, Sengupta D, Bhargava  
271 NC: **Mucocutaneous disorders in Hiv positive patients.** *Indian journal of*  
272 *dermatology, venereology and leprology* 1996, **62**(5):283-285.
  - 273 12. Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, Li Z, Chen Y, Mu D, Fan W:  
274 **Estimation of genomic characteristics by analyzing k-mer frequency in de**  
275 **novo genome projects.** *Quantitative Biology* 2013, **35**(s 1–3):62-67.
  - 276 13. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y *et*  
277 *al*: **SOAPdenovo2: an empirically improved memory-efficient short-read**  
278 **de novo assembler.** *Gigascience* 2012, **1**:18.
  - 279 14. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J: **SOAP2: an**  
280 **improved ultrafast tool for short read alignment.** *Bioinformatics* 2009,  
281 **25**(15):1966-1967.
  - 282 15. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate**  
283 **core genes in eukaryotic genomes.** *Bioinformatics* 2007, **23**(9):1061-1067.
  - 284 16. Sim AFO, Waterhouse MR, Ioannidis P, Kriventseva VE, Zdobnov ME:  
285 **BUSCO: assessing genome assembly and annotation completeness with**  
286 **single-copy orthologs.** *Bioinformatics* 2015, **31**(19):3210-3212.
  - 287 17. Kriventseva EV, Tegenfeldt F, Petty TJ, Waterhouse RM, Simão FA,  
288 Pozdnyakov IA, Ioannidis P, Zdobnov EM: **OrthoDB v8: update of the**  
289 **hierarchical catalog of orthologs and the underlying free software.** *Nucleic*  
290 *Acids Research* 2015, **43**(Database issue):D250-D256.

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 291 18. Maziade M, Bouchard S, Gingras N, Charron L, Cardinal A, Roy MA,  
292 Gauthier B, Tremblay G, Cote S, Fournier C *et al*: **Long-term stability of**  
293 **diagnosis and symptom dimensions in a systematic sample of patients**  
294 **with onset of schizophrenia in childhood and early adolescence. II:**  
295 **Postnegative distinction and childhood predictors of adult outcome.** *The*  
296 *British journal of psychiatry : the journal of mental science* 1996,  
297 **169(3):371-378.**
- 298 19. Jurka, J., Kapitonov, V. V., Pavlicek, A. ,et al: **Rebase Update, a database**  
299 **of eukaryotic repetitive elements.** *Cytogenetic & Genome Research* 2005,  
300 **110(1-4):462-467.**
- 301 20. Chen N.: Using RepeatMasker to Identify Repetitive Elements in Genomic  
302 Sequences[J]. 2004, Chapter 4(Unit 4):4.10.1-4.10.14.
- 303 21. Benson G, . **Tandem repeats finder: a program to analyze DNA sequences.**  
304 *Nucleic Acids Research* 1999, **27(2):573-580.**
- 305 22. Mario S, Oliver K, Irfan G, Alec H, Stephan W, Burkhard M: **AUGUSTUS:**  
306 **ab initio prediction of alternative transcripts.** *Nucleic Acids Research* 2006,  
307 **34:435-439.**
- 308 23. Burge C., Karlin S., **Prediction of complete gene structures in human**  
309 **genomic DNA.** *Journal of Molecular Biology* 1997, **268(1):78-94.**
- 310 24. Collins JE, White S, Searle SMJ, Stemple DL: **Incorporating RNA-seq data**  
311 **into the zebrafish Ensembl genebuild.** *Genome Research* 2012,  
312 **22(10):2067-2078.**
- 313 25. Kasahara M, Naruse K, Sasaki S, Nakatani Y, Qu W, Ahsan B, Yamada T,  
314 Nagayasu Y, Doi K, Kasai Y: **The medaka draft genome and insights into**  
315 **vertebrate genome evolution.** *Nature* 2007, **447(7145):714-719.**
- 316 26. Kesteven GL: **Whole-Genome Shotgun Assembly and Analysis of the**  
317 **Genome of Fugu rubripes.** *Science (New York, NY)* 2002,  
318 **297(5585):1301-1310.**
- 319 27. Jaillon O, Aury JM, Brunet F, Petit JL, Stangethomann N, Mauceli E,  
320 Bouneau L, Fischer C, Ozoufcozaz C, Bernot A: **Genome duplication in the**  
321 **teleost fish Tetraodon nigroviridis reveals the early vertebrate**  
322 **proto-karyotype.** *Nature* 2004, **431(7011):946-957.**
- 323 28. Rondeau EB, Minkley DR, Leong JS, Messmer AM, Jantzen JR, von  
324 Schalburg KR, Lemon C, Bird NH, Koop BF: **The genome and linkage map**  
325 **of the northern pike (Esox lucius): conserved syntenic revealed between**  
326 **the salmonid sister group and the Neoteleostei.** *PLoS ONE* 2014,  
327 **9(7):e102089.**
- 328 29. Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford  
329 R, Pirun M, Zody MC, White S: **The genomic basis of adaptive evolution in**  
330 **threespine sticklebacks.** *Nature* 2012, **484(7392):55-61.**
- 331 30. Pevsner J: **Basic Local Alignment Search Tool (BLAST):** John Wiley &  
332 Sons, Inc.; 2005.

- 1 333 31. Yu XJ, Zheng HK, Wang J, Wang W, Su B: **Detecting lineage-specific**  
2 334 **adaptive evolution of brain-expressed genes in human using rhesus**  
3 335 **macaque as outgroup.** *Genomics* 2006, **88**(6):745-751.
- 4 336 32. Birney E, Clamp M, Durbin R: **GeneWise and Genomewise.** *Genome*  
5 337 *Research* 2004, **14**(5):988-995.
- 6 338 33. Elsik CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM:  
7 339 **Creating a honey bee consensus gene set.** *Genome Biology* 2007,  
8 340 **8**(1):90-105.
- 9 341 34. Bairoch A, Apweiler R: **The SWISS-PROT protein sequence database and**  
10 342 **its supplement TrEMBL in 2000.** *Nucleic Acids Res* 2000, **28**(1):45-48.
- 11 343 35. Finn RD: **Pfam: the protein families database.** *Nucleic Acids Research* 2014,  
12 344 **42**(Database issue):D222-230.
- 13 345 36. Attwood TK: **The PRINTS database: A resource for identification of**  
14 346 **protein families.** *Briefings in Bioinformatics* 2002, **3**(3):252-263.
- 15 347 37. Bru C, Courcelle E, Beausse Y, Dalmar S, Kahn D: **The ProDom database of**  
16 348 **protein domain families: more emphasis on 3D.** *Nucleic Acids Research*  
17 349 2005, **33**(Database issue):212-215.
- 18 350 38. Letunic I, Copley RR, Schmidt S, Ciccarelli FD, Doerks T, Schultz J, Ponting  
19 351 CP, Bork P: **SMART 4.0: towards genomic data integration.** *Nucleic Acids*  
20 352 *Research* 2004, **32**(Database issue):D142-D144.
- 21 353 39. Boeckmann B., Bairoch A., Apweiler R., Blatter M. C., Estreicher A.. *et al*:  
22 354 **The Swiss-Prot knowledgebase and its supplement TREMBL in 2003.**  
23 355 *Nucleic Acids Research* 2003, **31**(1):365-370.
- 24 356 40. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P,  
25 357 Das U, Daugherty L, Duquenne L: **InterPro: the integrative protein**  
26 358 **signature database.** *Nucleic Acids Research* 2009, **37**(suppl 1):D211-D215.
- 27 359 41. Hingamp P, Broek AEVD, Stoesser G, Baker W: **The EMBL nucleotide**  
28 360 **sequence database.** *Molecular Biotechnology* 1999, **12**(3):255-267.
- 29 361 42. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.**  
30 362 *Nucleic Acids Research* 2000, **27**(1):29-34(26).
- 31 363 43. Braasch I, Gehrke AR, Smith JJ, Kawasaki K, Manousaki T, Pasquier J,  
32 364 Amores A, Desvignes T, Batzel P, Catchen J *et al*: **The spotted gar genome**  
33 365 **illuminates vertebrate evolution and facilitates human-teleost**  
34 366 **comparisons.** *Nature genetics* 2016, **48**(4):427-437.
- 35 367 44. Bian C, Hu Y, Ravi V, Kuznetsova IS, Shen X, Mu X, Sun Y, You X, Li J, Li  
36 368 X *et al*: **The Asian arowana (Scleropages formosus) genome provides new**  
37 369 **insights into the evolution of an early lineage of teleosts.** *Scientific reports*  
38 370 2016, **6**:24501.
- 39 371 45. Ostlund G, Schmitt T, Forslund K, Kostler T, Messina DN, Roopra S, Frings  
40 372 O, Sonnhammer EL: **InParanoid 7: new algorithms and tools for**  
41 373 **eukaryotic orthology analysis.** *Nucleic acids research* 2010, **38**(Database  
42 374 issue):D196-203.

375 46. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy**  
376 **and high throughput**. *Nucleic acids research* 2004, **32**(5):1792-1797.

377 47. Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Hohna S,  
378 Larget B, Liu L, Suchard MA, Huelsenbeck JP: **MrBayes 3.2: efficient**  
379 **Bayesian phylogenetic inference and model choice across a large model**  
380 **space**. *Systematic biology* 2012, **61**(3):539-542.

381 48. Proost S, Fostier J, De Witte D, Dhoedt B, Demeester P, Van de Peer Y,  
382 Vandepoele K: **i-ADHoRe 3.0--fast and sensitive detection of genomic**  
383 **homology in extremely large data sets**. *Nucleic acids research* 2012,  
384 **40**(2):e11.

385 49. Yang Z: **PAML: a program package for phylogenetic analysis by**  
386 **maximum likelihood**. *Computer applications in the biosciences : CABIOS*  
387 1997, **13**(5):555-556.

388  
389  
390 **Tables**

391 **Table 1.** The statistics of genome assembly and annotation for *P. hyalocranius*.

Genome assembly	
Contig N50 size (kb)	17.2
Scaffold N50 size (Mb)	1.163
Estimated genome size (Mb)	525
Assembled genome size (Mb)	536
Genome coverage (X)	315
The longest scaffold (bp)	5,398,389
Genome annotation	
Protein-coding gene number	19,884
Annotated functional gene number	19,125 (96.2%)
Unannotated functional gene number	759 (3.8%)
Repeat content	24.43%

392  
393  
394  
395

396

397 **Table 2.** Detailed classification of repeat sequences in the assembled genome.

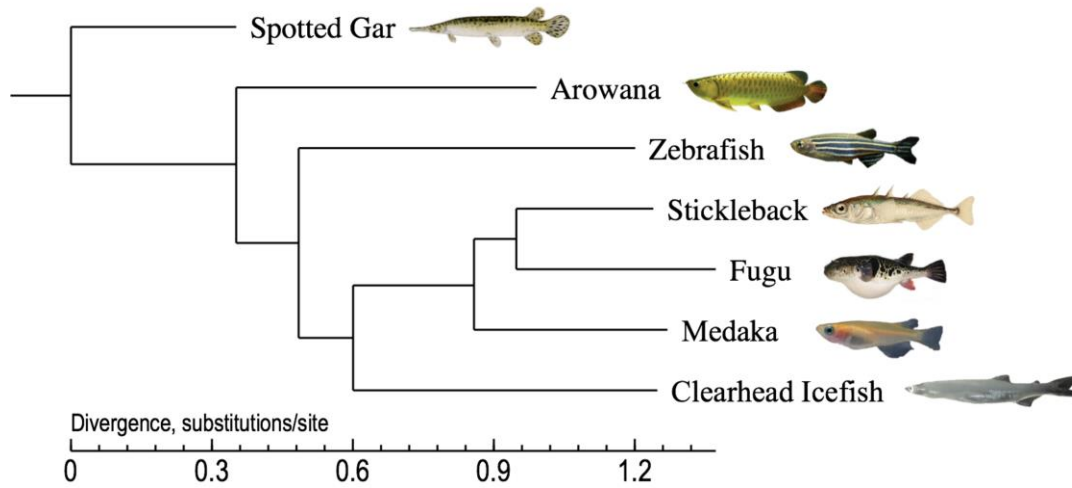
Type	Repeat Size(bp)	% of Genome
ProteinMask	9925152	1.85
RepeatMasker	5948136	1.11
TRF	66595756	12.41
De novo	93726009	17.47
Total	131090229	24.43

398

399

400

401

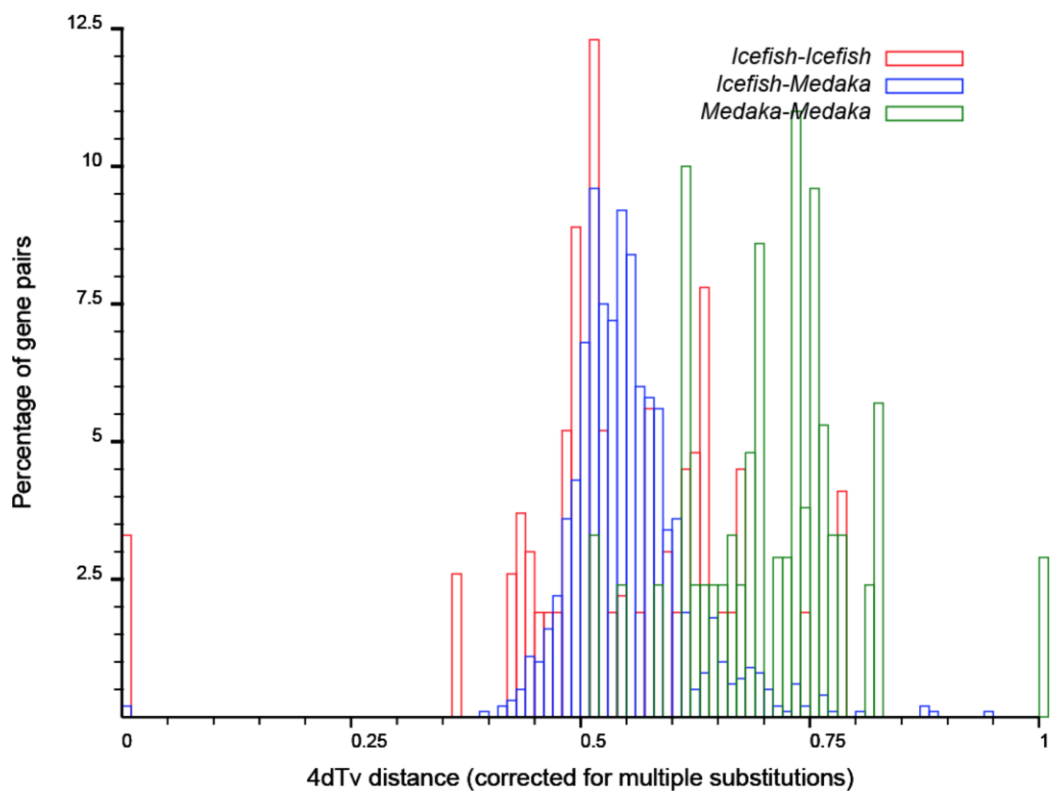


402

403 **Figure 1.** Phylogeny of seven representative ray-finned fishes. The spotted gar was used  
404 as the outgroup species.

405

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



406

407 **Figure 2.** Distribution of 4DTV distances between the clearhead icefish and medaka.

408 The horizontal axis stands for the 4DTV distance corrected using the HKY model. The

409 vertical axis represents the percentage of colinear gene pairs.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65