

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); zhouyf@ffrc.cn(YZ);

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

1 26 zhangmy@ffrc.cn(MZ); youxinxin@genomics.cn (XY); youy@ffrc.cn (YY);  
2 27 chenjieming@genomics.cn (JC); yuhui@genomics.cn (HY); xugc@ffrc.cn(GX);  
3 28 fangda@ffrc.cn(DF); qiangj@ffrc.cn(JQ); 420219380@qq.com(SJ); hej@ffrc.cn(JH);  
4 29 xujunmin@genomics.cn (JX); shiqiong@genomics.cn (QS); 13906292412@139.com  
5 30 (ZZ); xup@ffrc.cn (PX)  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15

## 16 33 **Abstract**

17 34 **Background:** Chinese clearhead icefish, *Protosalanx hyalocranius*, is a  
18 35 representative icefish species with economic importance and special appearance. Due  
19 36 to its great economic values in China, the fish was introduced into Lake Dianchi and  
20 37 several other lakes from the Lake Taihu half a century ago. Similar to the  
21 38 *Sinocyclocheilus* cavefish, the clearhead icefish has certain cavefish-like traits, such  
22 39 as transparent body and nearly scaleless skin. Here, we provide the whole genome  
23 40 sequence of this surface-dwelling fish and generated a draft genome assembly, aiming  
24 41 at exploring molecular mechanisms for the biological interests.  
25  
26  
27  
28  
29  
30  
31  
32  
33

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

47 47 **Conclusion:** We report the first draft genome of the Chinese clearhead icefish. The  
48 48 genome assembly will provide a solid foundation for further molecular breeding and  
49 49 germplasm resource protection in Chinese clearhead icefish, as well as other icefishes.  
50 50 It is also a valuable genetic resource for revealing the molecular mechanisms for the  
51 51 cavefish-like characters.  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

56 52 **Keywords:** Clearhead icefish; *Protosalanx hyalocranius*; Whole genome  
57 53 sequencing; Genome assembly; Gene prediction; Repetitive sequences  
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*, Figure 1), a diadromous fish, mainly inhabits in coastal  
60 areas and adjacent freshwaters [4-6]. As an economically important fish in China, the  
61 clearhead icefish was widely introduced into some lakes from the original Lake Taihu  
62 half a century ago, and it has developed a resident life history in these water areas [2,  
63 7, 8]. 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

5 83

### 6 84 *Genome size estimation and genome assembly*

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  
12  
13 87 Ns and low-quality bases) and PCR-replicates as well as adaptor sequences. In total,  
14  
15 88 we obtained 169.0 Gb of clean reads. Subsequently, we estimated the genome size  
16  
17 89 based on the 17-mer depth frequency distribution method [12]. We applied the  
18  
19 90 following formula to calculate the genome size:  $G = k\_num / k\_depth = b\_num / b\_depth$   
20  
21 91 ( $k\_num$  is the total number of K-mers from the sequencing data,  $k\_depth$  is the  
22  
23 92 expected coverage depth for k-mers,  $b\_num$  is the total number of bases,  $b\_depth$  is  
24  
25 93 the expected coverage depth of bases; As one read with length  $L$  generates  $L - K + 1$   
26  
27 94 k-mers,  $k\_num / b\_num = (L - K + 1) / L$ ). In our current study, the  $K\_num$  was  
28  
29 95 10,500,000,000 and the  $K\_depth$  was 20. Hence, we estimated that the genome size of  
30  
31 96 Chinese clearhead icefish is 525 Mb.  
32  
33

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

114

### 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).

127

### 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. Finally, we identified 23,132 and 21,379 pro-coding  
136 genes by using the AUGUSTUS and GENSCAN software (Table 3).

1 137 **2) Homology annotation.** We aligned the protein sequences from six published  
2 138 genomes, including *Danio rerio* [24], *Oryzias latipes* [25], *Takifugu rubripes* [26],  
3  
4 139 *Tetraodon nigroviridis* [27], *Esox lucius* [28] and *Gasterosteus aculeatus* [29], against  
5  
6 140 our assembly to predict homology-based genes. The potential homology-based genes  
7  
8 141 were searched by TblastN [30] with an e-value of  $10^{-5}$ . The TblastN results were then  
9  
10 142 processed by SOLAR (Sorting Out Local Alignment Result [31]) to obtain the best hit  
11  
12 143 of each alignment. Subsequently, GeneWise v2.2.0 [32] was performed to detect the  
13  
14 144 possible gene structure for the best hit of each alignment. The low-quality genes were  
15  
16 145 also removed as described in the above-mentioned *de novo* annotation.  
17  
18  
19

20 146 **3) Integration of annotation results.** We employed the GLEAN [33] to generate a  
21  
22 147 non-redundant and comprehensive gene set. Finally, the best hit of each protein was  
23  
24 148 obtained through all protein sequences from the GLEAN results aligned to the  
25  
26 149 databases of the SwissProt and TrEMBL [34] (Uniprot release 2011.06) by BlastP  
27  
28 150 with an e-value of  $10^{-5}$ . Overall, we generated a final gene set with 19,884 genes for  
29  
30 151 the Chinese clearhead icefish (Table 3).  
31  
32

33  
34 152 CEGMA was performed again to evaluate the coverage rate between KOG  
35  
36 153 (EuKaryotic Orthologous Groups) genes predicted by CEGMA and the predicted total  
37  
38 154 gene set. It demonstrates that the predicted gene set mapped 96.4% of the KOGs.  
39  
40 155 Simultaneously, the BUSCO was implemented again to assess completeness of the  
41  
42 156 predicted gene set. The BUSCO values were calculated as follows: C: 79% [D: 16%],  
43  
44 157 F: 9.8%, M: 10%, n: 843 (C: complete [D: duplicated], F: fragmented, M: missed, n:  
45  
46 158 genes). The assessment values from both CEGMA and BUSCO proved high accuracy  
47  
48 159 of the annotation.  
49  
50

51 160 **4) Function annotation.** The predicted protein sequences of the clearhead icefish  
52  
53 161 were aligned against several public databases (Pfam [35], PRINTS [36], ProDom [37]  
54  
55 162 and SMART [38]) for detection of functional motifs and domains. Finally, we found  
56  
57 163 that 96.2% of the predicted total gene set had been annotated with at least one  
58  
59  
60  
61  
62  
63  
64  
65

1 164 functional assignment from other public databases (Swiss-Prot [39], Interpro [40],  
2 165 TrEMBL [41] and KEGG [42]).  
3  
4

5 166

### 6 167 *Genome evolution*

7  
8  
9 168 We performed phylogenomic analyses with orthologues from representative species  
10  
11 169 for each clade. We used the Ensembl BioMart ([www.ensembl.org/biomart](http://www.ensembl.org/biomart); Ensembl  
12  
13 170 version 76) to extract orthologues for zebrafish [24], fugu [26], stickleback [29],  
14  
15 171 medaka [25] and spotted gar [43]. This generated orthologue dataset from six species  
16  
17 172 was filtered out to retain only one-to-one orthologues. Meanwhile, a new Asian  
18  
19 173 arowana gene set stem from our recent work [44]. In order to extrapolate the Biomart  
20  
21 174 orthologues to the arowana and clearhead icefish gene sets, we used zebrafish as the  
22  
23 175 reference. We ran InParanoid [45] for the three species pairs (zebrafish-arowana and  
24  
25 176 zebrafish-clearhead icefish) at default settings (i.e., a minimum BLASTP score of 40  
26  
27 177 bits, minimum 50% alignment span, minimum 25% alignment coverage, and  
28  
29 178 minimum inparalog confidence level of 0.05). By comparing the three InParanoid  
30  
31 179 outputs, we narrowed down the list of one-to-one orthologues, presented in all the  
32  
33 180 seven species, to 454 genes. Multiple alignments were subsequently performed on  
34  
35 181 proteins of each selected family using MUSCLE (version 3.8.31) [46] and protein  
36  
37 182 alignments were converted to their corresponding CDS alignments using an in-house  
38  
39 183 perl script (see supporting data). All the translated CDS sequences were linked into  
40  
41 184 one “supergene” for each species. Non-degenerated sites extracted from the  
42  
43 185 supergenes were subsequently joined into new sequence of each species to construct a  
44  
45 186 phylogenetic tree (Figure 2) using MrBayes [47] (GTR+gamma model, Version 3.2).  
46  
47 187 Our phylogenetic data demonstrate the phylogenetic position of the clearhead icefish  
48  
49 188 (Figure 2).  
50  
51  
52  
53

54 189

### 55 190 *Synten blocks and genome duplication*

56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 191 Genomic homology between the clearhead icefish and Nile tilapia [48] was examined  
2  
3 192 using i-ADHoRe 3.0 [49] using the following settings: alignment method gg2, gap  
4  
5 193 size 30, tandem gap 30, cluster gap 35, q value of 0.85, prob cutoff 0.01, anchor  
6  
7 194 points 5 and using multiple hypothesis correction FDR. The output of this was  
8  
9 195 processed by the pipeline and incorporated in a relational database to which  
10  
11 196 visualization programs can connect and on which additional statistical analysis can  
12  
13 197 then be performed. For synteny detection, the cloud mode was enabled (cluster\_type  
14  
15 198 = cloud) and appropriate settings were selected as follows: cloud\_gap\_size 20,  
16  
17 199 cloud\_cluster\_gap 20, cloud\_filter\_method binomial, prob cutoff 0.01, anchor points  
18  
19 200 5, multiple hypothesis correction FDR and level\_2\_only true. Finally, we identified  
20  
21 201 771 synteny blocks containing 7,057 genes between the clearhead icefish and Nile  
22  
23 202 tilapia.

24  
25  
26 203 Subsequently, Protein sequences of homologous gene pairs in the identified syntenic  
27  
28 204 regions were aligned using MUSCLE [46], and the protein alignments were then  
29  
30 205 converted to the CDS alignments. Finally, four-fold degenerative third-codon  
31  
32 206 transversion (4DTV) values were calculated on these CDS alignments and corrected  
33  
34 207 using the HKY model in the PAML package [50]. These data indicate that the  
35  
36 208 clearhead icefish also experienced the teleost-specific whole genome duplication  
37  
38 209 (WGD) (Figure 3).

40  
41 210

## 42 211 **Conclusion**

43  
44  
45 212 We generated a draft genome assembly of the Chinese clearhead icefish. The novel  
46  
47 213 genome data were deposited in publicly accessible repositories to promote further  
48  
49 214 biological research, molecular breeding and resource protection of this representative  
50  
51 215 and valuable icefish.

52  
53 216

## 54 217 **Availability of supporting data**



1 218 Supporting data and materials are available in the *GigaScience* GigaDB database [51],  
2  
3 219 with the raw genome sequences deposited in the SRA under the bioproject number  
4  
5 220 PRJNA328051.  
6

7 221

## 9 222 **Competing interests**

11 223 The authors declare that they have no competing interests.  
12  
13 224

## 16 225 **Funding**

19 226 This study was supported by a grant from Natural Science Foundation of Jiangsu  
20  
21 227 Province (No.BK2012093), fish investigation in Taihu Lake (No.TH2016WT007),  
22  
23 228 National Infrastructure of Fishery Germplasm Resources (No.2016DKA30470), Basic  
24  
25 229 Research Funds from Freshwater Fisheries Research Center (No. 2013JBFM07),  
26  
27 230 Special Project on the Integration of Industry, Education and Research of Guangdong  
28  
29 231 Province (No. 2013B090800017), Shenzhen Special Program for Future Industrial  
30  
31 232 Development (N 192 o. JSGG20141020113728803), and Zhenjiang Leading Talent  
32  
33 233 Program for Innovation and Entrepreneurship.  
34  
35

## 37 234 **Author's Contributions**

39 235 KL, PX, QS, DX, JX, CB and ZZ conceived the project. MZ, XY, HY, JC, GX, DF,  
40  
41 236 JQ, SJ and JH collected the samples and extracted the genomic DNA. JL, CB and HY  
42  
43 237 performed the genome assembly and data analysis. JL, CB, QS, KL, XP, KL, YY and  
44  
45 238 ZZ wrote the paper.  
46  
47

48 239

## 50 240 **References**

- 54 241 1. Wang ZS, Cui Zhang FU: **Biodiversity of Chinese Icefishes (Salangidae)**  
55 242 **and their conserving strategies.** *Chinese Biodiversity* 2002, **10**(4):416-424.  
56  
57 243 2. Zhang J, Li M, Xu M, Takita T, Wei F: **Molecular phylogeny of icefish**  
58 244 **Salangidae based on complete mtDNA cytochrome b sequences, with**  
59  
60  
61  
62  
63  
64  
65

245 **comments on estuarine fish evolution.** *Biological Journal of the Linnean*  
246 *Society* 2007, **91**(2):325-340.

247 3. Wang Z, Lu C, Hu H, Xu C, Lei G: **Dynamics of Icefish (Salangidae) Stocks**  
248 **in Nanyi Lake, Eastern China: Degradation and Overfishing.** *Journal of*  
249 *Freshwater Ecology* 2004, **19**(2):271-278.

250 4. Xia DQ, Cao Y, Ting ting WU, Yang H: **Study on lineages of Protosalanx**  
251 **chinensis, Neosalanx taihuensis and N.oligodontis in Taihu Lake with**  
252 **RAPD technique.** *Journal of Fisheryences of China* 2000, **7**(01):12-15.

253 5. Xia DQ, Cao Y, Ting Ting WU, Yang H: **Genetic Structures of Population**  
254 **of Protosalanx Chinensis, Neosalanx Taihuensis and Neosalanx**  
255 **Oligodontis in Lake Taihu.** *Journal of Fisheries of China* 1999(03):254-260.

256 6. Armani A, Castigliengo L, Tinacci L, Gianfaldoni D, Guidi A: **Molecular**  
257 **characterization of icefish, (S alangidae family), using direct sequencing**  
258 **of mitochondrial cytochrome b gene.** *Food Control* 2011, **22**(6):888-895.

259 7. Wang Z, Lu C, Hu H, Zhou Y, Xu C, Lei G: **Freshwater icefishes**  
260 **(Salangidae) in the Yangtze River basin of China: Spatial distribution**  
261 **patterns and environmental determinants.** *Environmental Biology of Fishes*  
262 2005, **73**(3):253-262.

263 8. Ye S, Yang J, Liu H, Oshima Y: **Use of elemental fingerprint analysis to**  
264 **identify localities of collection for the large icefish protosalanx chinensis**  
265 **in Taihu Lake, China.** *Journal of the Faculty of Agriculture, Kyushu*  
266 *University* 2011, **56**(1):41-45.

267 9. Yang J, Chen X, Jie B, Fang D, Ying Q, Jiang W, Hui Y, Chao B, Jiang L, He  
268 S: **The Sinocyclocheilus cavefish genome provides insights into cave**  
269 **adaptation.** *Bmc Biology* 2016, **14**(1):1-13.

270 10. Xu J, Xie P, Zhang M, Zhou Q, Zhang L, Wen Z, Cao T: **Icefish (salangidae)**  
271 **as an indicator of anthropogenic pollution in freshwater systems using**  
272 **nitrogen isotope analysis.** *Bulletin of environmental contamination and*  
273 *toxicology* 2007, **79**(3):323-326.

274 11. Kar HK, Narayan R, Gautam RK, Jain RK, Doda V, Sengupta D, Bhargava  
275 NC: **Mucocutaneous disorders in Hiv positive patients.** *Indian journal of*  
276 *dermatology, venereology and leprology* 1996, **62**(5):283-285.

277 12. Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, Li Z, Chen Y, Mu D, Fan W:  
278 **Estimation of genomic characteristics by analyzing k-mer frequency in de**  
279 **novo genome projects.** *Quantitative Biology* 2013, **35**(s 1–3):62-67.

280 13. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y *et*  
281 *al*: **SOAPdenovo2: an empirically improved memory-efficient short-read**  
282 **de novo assembler.** *Gigascience* 2012, **1**:18.

283 14. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J: **SOAP2: an**  
284 **improved ultrafast tool for short read alignment.** *Bioinformatics* 2009,  
285 **25**(15):1966-1967.

- 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
- 286 15. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate**  
287 **core genes in eukaryotic genomes.** *Bioinformatics* 2007, **23**(9):1061-1067.
- 288 16. Sim AFO, Waterhouse MR, Ioannidis P, Kriventseva VE, Zdobnov ME:  
289 **BUSCO: assessing genome assembly and annotation completeness with**  
290 **single-copy orthologs.** *Bioinformatics* 2015, **31**(19):3210-3212.
- 291 17. Kriventseva EV, Tegenfeldt F, Petty TJ, Waterhouse RM, Simão FA,  
292 Pozdnyakov IA, Ioannidis P, Zdobnov EM: **OrthoDB v8: update of the**  
293 **hierarchical catalog of orthologs and the underlying free software.** *Nucleic*  
294 *Acids Research* 2015, **43**(Database issue):D250-D256.
- 295 18. Maziade M, Bouchard S, Gingras N, Charron L, Cardinal A, Roy MA,  
296 Gauthier B, Tremblay G, Cote S, Fournier C *et al*: **Long-term stability of**  
297 **diagnosis and symptom dimensions in a systematic sample of patients**  
298 **with onset of schizophrenia in childhood and early adolescence. II:**  
299 **Postnegative distinction and childhood predictors of adult outcome.** *The*  
300 *British journal of psychiatry : the journal of mental science* 1996,  
301 **169**(3):371-378.
- 302 19. Jurka, J., Kapitonov, V. V., Pavlicek, A. ,et al: **Repbase Update, a database**  
303 **of eukaryotic repetitive elements.** *Cytogenetic & Genome Research* 2005,  
304 **110**(1-4):462-467.
- 305 20. Chen N.: Using RepeatMasker to Identify Repetitive Elements in Genomic  
306 Sequences[J]. 2004, Chapter 4(Unit 4):4.10.1-4.10.14.
- 307 21. Benson G, . **Tandem repeats finder: a program to analyze DNA sequences.**  
308 *Nucleic Acids Research* 1999, **27**(2):573-580.
- 309 22. Mario S, Oliver K, Irfan G, Alec H, Stephan W, Burkhard M: **AUGUSTUS:**  
310 **ab initio prediction of alternative transcripts.** *Nucleic Acids Research* 2006,  
311 **34**:435-439.
- 312 23. Burge C., Karlin S., **Prediction of complete gene structures in human**  
313 **genomic DNA.** *Journal of Molecular Biology* 1997, **268**(1):78-94.
- 314 24. Collins JE, White S, Searle SMJ, Stemple DL: **Incorporating RNA-seq data**  
315 **into the zebrafish Ensembl genebuild.** *Genome Research* 2012,  
316 **22**(10):2067-2078.
- 317 25. Kasahara M, Naruse K, Sasaki S, Nakatani Y, Qu W, Ahsan B, Yamada T,  
318 Nagayasu Y, Doi K, Kasai Y: **The medaka draft genome and insights into**  
319 **vertebrate genome evolution.** *Nature* 2007, **447**(7145):714-719.
- 320 26. Kesteven GL: **Whole-Genome Shotgun Assembly and Analysis of the**  
321 **Genome of Fugu rubripes.** *Science (New York, NY)* 2002,  
322 **297**(5585):1301-1310.
- 323 27. Jaillon O, Aury JM, Brunet F, Petit JL, Stangethomann N, Mauceli E,  
324 Bouneau L, Fischer C, Ozoufcozaz C, Bernot A: **Genome duplication in the**  
325 **teleost fish Tetraodon nigroviridis reveals the early vertebrate**  
326 **proto-karyotype.** *Nature* 2004, **431**(7011):946-957.

- 1 327 28. Rondeau EB, Minkley DR, Leong JS, Messmer AM, Jantzen JR, von  
2 328 Schalburg KR, Lemon C, Bird NH, Koop BF: **The genome and linkage map**  
3 329 **of the northern pike (*Esox lucius*): conserved synteny revealed between**  
4 330 **the salmonid sister group and the Neoteleostei.** *PLoS ONE* 2014,  
5 331 **9(7):e102089.**
- 6 332 29. Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford  
7 333 R, Pirun M, Zody MC, White S: **The genomic basis of adaptive evolution in**  
8 334 **threespine sticklebacks.** *Nature* 2012, **484(7392):55-61.**
- 9 335 30. Pevsner J: **Basic Local Alignment Search Tool (BLAST):** John Wiley &  
10 336 Sons, Inc.; 2005.
- 11 337 31. Yu XJ, Zheng HK, Wang J, Wang W, Su B: **Detecting lineage-specific**  
12 338 **adaptive evolution of brain-expressed genes in human using rhesus**  
13 339 **macaque as outgroup.** *Genomics* 2006, **88(6):745-751.**
- 14 340 32. Birney E, Clamp M, Durbin R: **GeneWise and Genomewise.** *Genome*  
15 341 *Research* 2004, **14(5):988-995.**
- 16 342 33. Elsik CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM:  
17 343 **Creating a honey bee consensus gene set.** *Genome Biology* 2007,  
18 344 **8(1):90-105.**
- 19 345 34. Bairoch A, Apweiler R: **The SWISS-PROT protein sequence database and**  
20 346 **its supplement TrEMBL in 2000.** *Nucleic Acids Res* 2000, **28(1):45-48.**
- 21 347 35. Finn RD: **Pfam: the protein families database.** *Nucleic Acids Research* 2014,  
22 348 **42(Database issue):D222-230.**
- 23 349 36. Attwood TK: **The PRINTS database: A resource for identification of**  
24 350 **protein families.** *Briefings in Bioinformatics* 2002, **3(3):252-263.**
- 25 351 37. Bru C, Courcelle E, Beausse Y, Dalmar S, Kahn D: **The ProDom database of**  
26 352 **protein domain families: more emphasis on 3D.** *Nucleic Acids Research*  
27 353 2005, **33(Database issue):212-215.**
- 28 354 38. Letunic I, Copley RR, Schmidt S, Ciccarelli FD, Doerks T, Schultz J, Ponting  
29 355 CP, Bork P: **SMART 4.0: towards genomic data integration.** *Nucleic Acids*  
30 356 *Research* 2004, **32(Database issue):D142-D144.**
- 31 357 39. Boeckmann B., Bairoch A., Apweiler R., Blatter M. C., Estreicher A.. *et al*:  
32 358 **The Swiss-Prot knowledgebase and its supplement TREMBL in 2003.**  
33 359 *Nucleic Acids Research* 2003, **31(1):365-370.**
- 34 360 40. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P,  
35 361 Das U, Daugherty L, Duquenne L: **InterPro: the integrative protein**  
36 362 **signature database.** *Nucleic Acids Research* 2009, **37(suppl 1):D211-D215.**
- 37 363 41. Hingamp P, Broek AEVD, Stoesser G, Baker W: **The EMBL nucleotide**  
38 364 **sequence database.** *Molecular Biotechnology* 1999, **12(3):255-267.**
- 39 365 42. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.**  
40 366 *Nucleic Acids Research* 2000, **27(1):29-34(26).**
- 41 367 43. Braasch I, Gehrke AR, Smith JJ, Kawasaki K, Manousaki T, Pasquier J,  
42 368 Amores A, Desvignes T, Batzel P, Catchen J *et al*: **The spotted gar genome**

- 369 **illuminates vertebrate evolution and facilitates human-teleost**  
370 **comparisons.** *Nature genetics* 2016, **48**(4):427-437.
- 371 44. Bian C, Hu Y, Ravi V, Kuznetsova IS, Shen X, Mu X, Sun Y, You X, Li J, Li  
372 X *et al*: **The Asian arowana (*Scleropages formosus*) genome provides new**  
373 **insights into the evolution of an early lineage of teleosts.** *Scientific reports*  
374 2016, **6**:24501.
- 375 45. Ostlund G, Schmitt T, Forslund K, Kostler T, Messina DN, Roopra S, Frings  
376 O, Sonnhammer EL: **InParanoid 7: new algorithms and tools for**  
377 **eukaryotic orthology analysis.** *Nucleic acids research* 2010, **38**(Database  
378 issue):D196-203.
- 379 46. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy**  
380 **and high throughput.** *Nucleic acids research* 2004, **32**(5):1792-1797.
- 381 47. Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Hohna S,  
382 Larget B, Liu L, Suchard MA, Huelsenbeck JP: **MrBayes 3.2: efficient**  
383 **Bayesian phylogenetic inference and model choice across a large model**  
384 **space.** *Systematic biology* 2012, **61**(3):539-542.
- 385 48. Brawand D, Wagner CE, Li YI, Malinsky M, Keller I, Fan S, Simakov O, Ng  
386 AY, Lim ZW, Bezault E *et al*: **The genomic substrate for adaptive**  
387 **radiation in African cichlid fish.** *Nature* 2014, **513**(7518):375-381.
- 388 49. Proost S, Fostier J, De Witte D, Dhoedt B, Demeester P, Van de Peer Y,  
389 Vandepoele K: **i-ADHoRe 3.0--fast and sensitive detection of genomic**  
390 **homology in extremely large data sets.** *Nucleic acids research* 2012,  
391 **40**(2):e11.
- 392 50. Yang Z: **PAML: a program package for phylogenetic analysis by**  
393 **maximum likelihood.** *Computer applications in the biosciences : CABIOS*  
394 1997, **13**(5):555-556.
- 395 51. Kai Liu; Dongpo Xu; Jia Li; Chao Bian; Jinrong Duan; Yanfeng Zhou;  
396 Mingying Zhang; Xinxin You; Yang You; Jieming Chen; Hui Yu; Gangchun  
397 Xu; Di-an Fang; Jun Qiang; Shulun Jiang; Jie He; Junmin Xu; Qiong Shi;  
398 Zhiyong Zhang; Pao Xu (2016) **Supporting data for "Whole genome**  
399 **sequencing of Chinese clearhead icefish, *Protosalanx hyalocranium*".**  
400 GigaScience Database. <http://doi.org/10.5524/100262>

## 404 Tables

405 **Table 1.** The statistics of genome assembly and annotation for *P. hyalocranium*.

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
Gap length (Mb)	122
<hr/>	
Genome annotation	
<hr/>	
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%
<hr/>	

406

407

408

409

410

411

412

413 **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
Tandem Repeat Finder	66595756	12.41
De novo	93726009	17.47
Total	131090229	24.43

414

415

Method		Number	Average Transcript Length (bp)	Average CDS Length (bp)	Average Exons Per Gene	Average Exons Length (bp)	Average Intron Length (bp)
<i>De novo</i>	AUGUSTUS	23,132	4,897.24	1,264.61	5.78	218.81	760.04
	GeneScan	21,379	17,213.49	1,973.56	10.22	193.05	1,652.41
Homolog	<i>Danio rerio</i>	25,390	7,156.92	1,312.32	6.17	212.62	1,129.99
	<i>Oryzias latipes</i>	25,319	6,411.36	1,194.58	5.89	202.73	1,066.29
	<i>Takifugu rubripes</i>	16,563	7,990.91	1,759.17	11.59	151.75	588.32

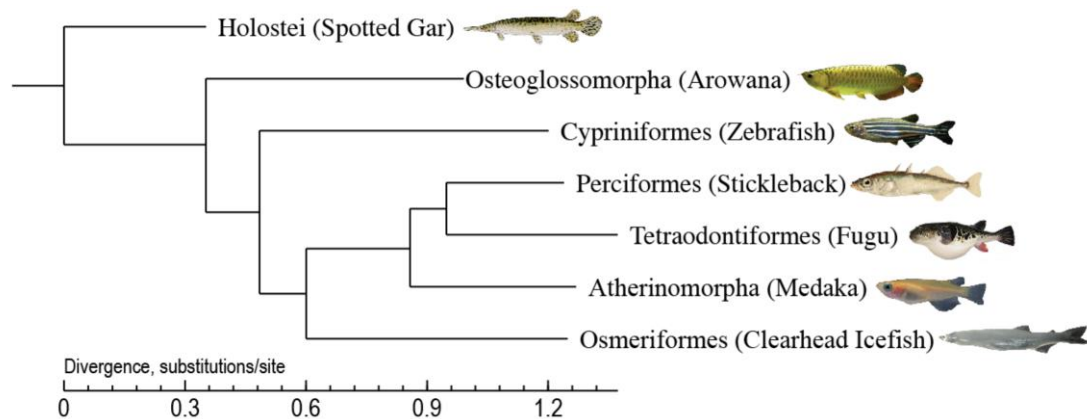
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  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

<i>Tetraodon nigroviridis</i>	19,128	8,335.40	1,351.98	7.44	181.78	1,084.78
<i>Esox lucius</i>	24,861	8,019.18	1,375.58	6.92	198.85	1,122.70
<i>Gasterosteus aculeatus</i>	25,354	6,819.62	1,183.46	6.18	191.44	1,087.68
Final gene set	19,884	12,889.35	1,821.79	9.13	199.49	1,360.92

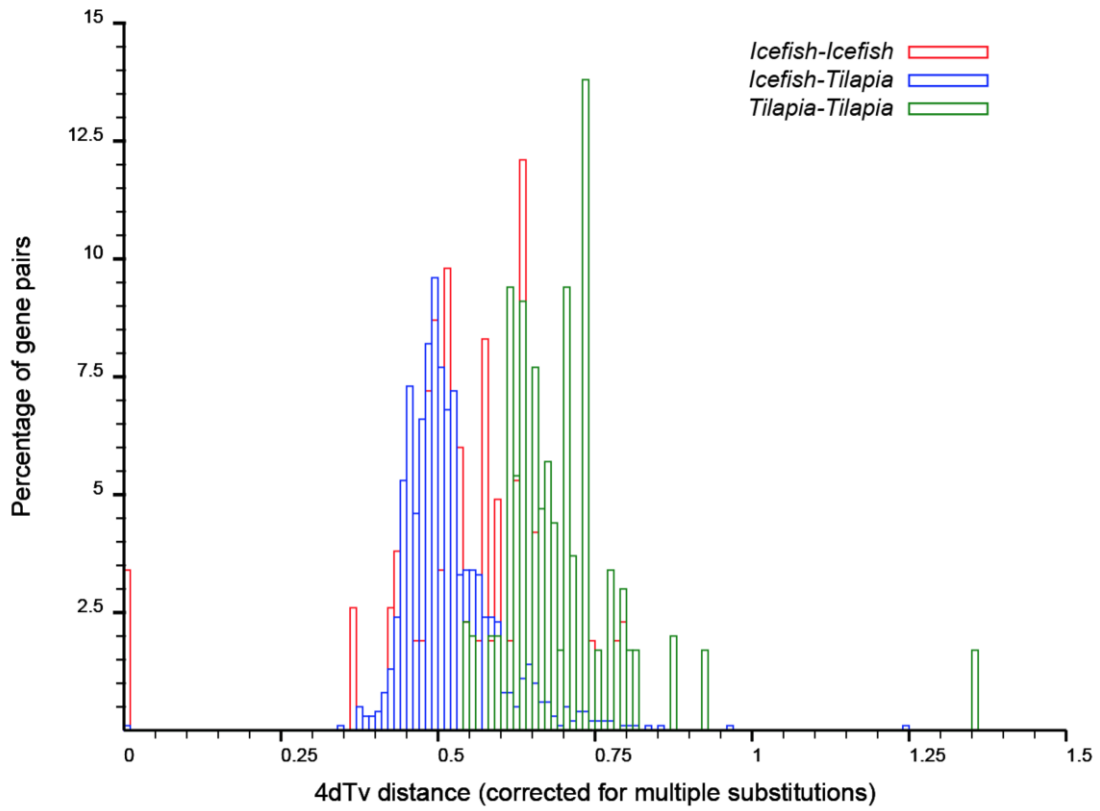
**Table 3.** Gene annotation statistics of the genome of *P. hyalocranius*.



**Figure 1.** Picture of a Chinese clearhead icefish. It was captured from the Taihu Lake of Jiangsu Province, China.



**Figure 2.** Phylogeny of seven representative ray-finned fishes. The spotted gar was used as the outgroup species.



429

430 **Figure 3.** Distribution of 4DTV distances between the clearhead icefish and tilapia. The  
 431 horizontal axis stands for the 4DTV distance corrected using the HKY model. The  
 432 vertical axis represents the percentage of collinear 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