

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

2 ***Protosalanx hyalocranius***

3

4 Kai Liu^{1†}, Dongpo Xu^{1†}, Jia Li^{2†}, Chao Bian^{2†}, Jinrong Duan^{1†}, Yanfeng Zhou^{1†},

5 Mingying Zhang¹, Xinxin You², Yang You¹, Jieming Chen², Hui Yu², Gangchun Xu¹,

6 Di-an Fang¹, Jun Qiang¹, Shulun Jiang¹, Jie He¹, Junmin Xu^{2,4,5}, Qiong Shi^{2,4,5,6*},

7 Zhiyong Zhang^{3*}, Pao Xu^{1,5*}

8

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

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

13 ³Institute of Oceanology & Marine Fisheries, Jiangsu 226007, China

14 ⁴BGI Zhenjiang Institute of Hydrobiology, Zhenjiang 212000, China

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

17 ⁶Laboratory 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);

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
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 54
59 55
60 56
61 57
62 58
63 59
64 60
65 61

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

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 CEGMA [15] and BUSCO
46
47 103 [16]. The CEGMA program (Core Eukaryotic Genes Mapping Approach; version 2.4)
48
49 104 assessment with 248 conserved Core Eukaryotic Genes (CEGs) was performed for
50
51 105 evaluation of the gene space completeness. Our results revealed that the assembled
52
53 106 genome had a CEGMA completeness score at 90.32% and 98.39%, which was
54
55 107 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; 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., minimum 50% alignment span,
26
27 177 minimum 25% alignment coverage, minimum BLASTP score of 40 bits, minimum
28
29 178 inparalog confidence level of 0.05). By comparing the three InParanoid outputs, we
30
31 179 narrowed down the list of one-to-one orthologues, presented in all the seven species,
32
33 180 to 454 genes. Subsequently, multiple alignments were performed on proteins of each
34
35 181 selected family by MUSCLE (version 3.8.31) [46] and protein alignments were
36
37 182 converted to their corresponding CDS alignments using an in-house perl script. All
38
39 183 the translated CDS sequences were linked into one “supergene” for each species.
40
41 184 Non-degenerated sites extracted from the supergenes were then joined into new
42
43 185 sequence of each species to construct a phylogenetic tree (Figure 1) using MrBayes
44
45 186 [47] (Version 3.2, GTR+gamma model). Our phylogenetic data demonstrate the
46
47 187 phylogenetic position of the clearhead icefish (Figure 1).
48
49
50
51

52 188

53 189 ***Synten blocks and genome duplication***

54
55 190 Genomic homology between the clearhead icefish and Nile tilapia [48] was examined
56
57 191 using i-ADHoRe 3.0 [49] using the following settings: alignment method gg2, gap
58
59
60
61
62
63
64
65

192 size 30, tandem gap 30, cluster gap 35, q value 0.85, prob cutoff 0.01, anchor points 5
193 and multiple hypothesis correction FDR. The output was processed by the pipeline
194 and included in a relational database to which visualization programs can connect and
195 on which additional statistical analysis can be performed. For synteny detection, the
196 cloud mode was enabled (cluster_type = cloud) and appropriate settings were selected
197 as follows: cloud_gap_size 20, cloud_cluster_gap 20, cloud_filter_method binomial,
198 prob cutoff 0.01, anchor points 5, multiple hypothesis correction FDR and
199 level_2_only true. Finally, we identified 771 synteny blocks containing 7,057 genes
200 between the clearhead icefish and Nile tilapia.
201 Subsequently, Protein sequences of homologous gene pairs in the identified syntenic
202 regions were aligned using MUSCLE [46], and the protein alignments were then
203 converted to the CDS alignments. Finally, four-fold degenerative third-codon
204 transversion (4DTV) values were calculated on these CDS alignments and corrected
205 using the HKY model in the PAML package [50]. These data indicate that the
206 clearhead icefish also experienced the teleost-specific whole genome duplication
207 (WGD) (Figure 2).

209 **Conclusion**

210 We generated a draft genome assembly of the Chinese clearhead icefish. The novel
211 genome data were deposited in publicly accessible repositories to promote further
212 biological research, molecular breeding and resource protection of this representative
213 and valuable icefish.

215 **Availability of supporting data**

216 Supporting data are available in the GigaDB database, and the raw genome sequences
217 are deposited in the SRA under the bioproject number PRJNA328051.

219 **Competing interests**

1 220 The authors declare that they have no competing interests.

2 221

3
4
5 222 **Funding**

6
7 223 This study was supported by a grant from Natural Science Foundation of Jiangsu
8
9 224 Province (No.BK2012093), fish investigation in Taihu Lake (No.TH2016WT007),
10
11 225 National Infrastructure of Fishery Germplasm Resources (No.2016DKA30470), Basic
12
13 226 Research Funds from Freshwater Fisheries Research Center (No. 2013JBFM07),
14
15 227 Special Project on the Integration of Industry, Education and Research of Guangdong
16
17 228 Province (No. 2013B090800017), Shenzhen Special Program for Future Industrial
18
19 229 Development (N 192 o. JSGG20141020113728803), and Zhenjiang Leading Talent
20
21 230 Program for Innovation and Entrepreneurship.

22
23
24
25 231 **Author's Contributions**

26
27 232 KL, PX, QS, DX, JX, CB and ZZ conceived the project. MZ, XY, HY, JC, GX, DF,
28
29 233 JQ, SJ and JH collected the samples and extracted the genomic DNA. JL, CB and HY
30
31 234 performed the genome assembly and data analysis. JL, CB, QS, KL, XP, KL, YY and
32
33 235 ZZ wrote the paper.

34
35
36
37 236
38
39 237 **References**

- 40
41
42 238 1. Wang ZS, Cui Zhang FU: **Biodiversity of Chinese Icefishes (Salangidae)**
43 239 **and their conserving strategies.** *Chinese Biodiversity* 2002, **10**(4):416-424.
44
45 240 2. Zhang J, Li M, Xu M, Takita T, Wei F: **Molecular phylogeny of icefish**
46 241 **Salangidae based on complete mtDNA cytochrome b sequences, with**
47 242 **comments on estuarine fish evolution.** *Biological Journal of the Linnean*
48 243 *Society* 2007, **91**(2):325-340.
49
50 244 3. Wang Z, Lu C, Hu H, Xu C, Lei G: **Dynamics of Icefish (Salangidae) Stocks**
51 245 **in Nanyi Lake, Eastern China: Degradation and Overfishing.** *Journal of*
52 246 *Freshwater Ecology* 2004, **19**(2):271-278.
53
54 247 4. Xia DQ, Cao Y, Ting ting WU, Yang H: **Study on lineages of Protosalanx**
55 248 **chinensis, Neosalanx taihuensis and N.oligodontis in Taihu Lake with**
56 249 **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
- 250 5. Xia DQ, Cao Y, Ting Ting WU, Yang H: **Genetic Structures of Population**
251 **of Protosalanx Chinensis, Neosalanx Taihuensis and Neosalanx**
252 **Oligodontis in Lake Taihu.** *Journal of Fisheries of China* 1999(03):254-260.
 - 253 6. Armani A, Castigliero L, Tinacci L, Gianfaldoni D, Guidi A: **Molecular**
254 **characterization of icefish, (S alangidae family), using direct sequencing**
255 **of mitochondrial cytochrome b gene.** *Food Control* 2011, **22**(6):888-895.
 - 256 7. Wang Z, Lu C, Hu H, Zhou Y, Xu C, Lei G: **Freshwater icefishes**
257 **(Salangidae) in the Yangtze River basin of China: Spatial distribution**
258 **patterns and environmental determinants.** *Environmental Biology of Fishes*
259 2005, **73**(3):253-262.
 - 260 8. Ye S, Yang J, Liu H, Oshima Y: **Use of elemental fingerprint analysis to**
261 **identify localities of collection for the large icefish protosalanx chinensis**
262 **in Taihu Lake, China.** *Journal of the Faculty of Agriculture, Kyushu*
263 *University* 2011, **56**(1):41-45.
 - 264 9. Yang J, Chen X, Jie B, Fang D, Ying Q, Jiang W, Hui Y, Chao B, Jiang L, He
265 S: **The Sinocyclocheilus cavefish genome provides insights into cave**
266 **adaptation.** *Bmc Biology* 2016, **14**(1):1-13.
 - 267 10. Xu J, Xie P, Zhang M, Zhou Q, Zhang L, Wen Z, Cao T: **Icefish (salangidae)**
268 **as an indicator of anthropogenic pollution in freshwater systems using**
269 **nitrogen isotope analysis.** *Bulletin of environmental contamination and*
270 *toxicology* 2007, **79**(3):323-326.
 - 271 11. Kar HK, Narayan R, Gautam RK, Jain RK, Doda V, Sengupta D, Bhargava
272 NC: **Mucocutaneous disorders in Hiv positive patients.** *Indian journal of*
273 *dermatology, venereology and leprology* 1996, **62**(5):283-285.
 - 274 12. Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, Li Z, Chen Y, Mu D, Fan W:
275 **Estimation of genomic characteristics by analyzing k-mer frequency in de**
276 **novo genome projects.** *Quantitative Biology* 2013, **35**(s 1–3):62-67.
 - 277 13. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y *et*
278 *al*: **SOAPdenovo2: an empirically improved memory-efficient short-read**
279 **de novo assembler.** *Gigascience* 2012, **1**:18.
 - 280 14. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J: **SOAP2: an**
281 **improved ultrafast tool for short read alignment.** *Bioinformatics* 2009,
282 **25**(15):1966-1967.
 - 283 15. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate**
284 **core genes in eukaryotic genomes.** *Bioinformatics* 2007, **23**(9):1061-1067.
 - 285 16. Sim AFO, Waterhouse MR, Ioannidis P, Kriventseva VE, Zdobnov ME:
286 **BUSCO: assessing genome assembly and annotation completeness with**
287 **single-copy orthologs.** *Bioinformatics* 2015, **31**(19):3210-3212.
 - 288 17. Kriventseva EV, Tegenfeldt F, Petty TJ, Waterhouse RM, Simão FA,
289 Pozdnyakov IA, Ioannidis P, Zdobnov EM: **OrthoDB v8: update of the**
290 **hierarchical catalog of orthologs and the underlying free software.** *Nucleic*
291 *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
- 292 18. Maziade M, Bouchard S, Gingras N, Charron L, Cardinal A, Roy MA,
293 Gauthier B, Tremblay G, Cote S, Fournier C *et al*: **Long-term stability of**
294 **diagnosis and symptom dimensions in a systematic sample of patients**
295 **with onset of schizophrenia in childhood and early adolescence. II:**
296 **Postnegative distinction and childhood predictors of adult outcome.** *The*
297 *British journal of psychiatry : the journal of mental science* 1996,
298 **169(3):371-378.**
- 299 19. Jurka, J., Kapitonov, V. V., Pavlicek, A. ,et al: **Rebase Update, a database**
300 **of eukaryotic repetitive elements.** *Cytogenetic & Genome Research* 2005,
301 **110(1-4):462-467.**
- 302 20. Chen N.: Using RepeatMasker to Identify Repetitive Elements in Genomic
303 Sequences[J]. 2004, Chapter 4(Unit 4):4.10.1-4.10.14.
- 304 21. Benson G, . **Tandem repeats finder: a program to analyze DNA sequences.**
305 *Nucleic Acids Research* 1999, **27(2):573-580.**
- 306 22. Mario S, Oliver K, Irfan G, Alec H, Stephan W, Burkhard M: **AUGUSTUS:**
307 **ab initio prediction of alternative transcripts.** *Nucleic Acids Research* 2006,
308 **34:435-439.**
- 309 23. Burge C., Karlin S., **Prediction of complete gene structures in human**
310 **genomic DNA.** *Journal of Molecular Biology* 1997, **268(1):78-94.**
- 311 24. Collins JE, White S, Searle SMJ, Stemple DL: **Incorporating RNA-seq data**
312 **into the zebrafish Ensembl genebuild.** *Genome Research* 2012,
313 **22(10):2067-2078.**
- 314 25. Kasahara M, Naruse K, Sasaki S, Nakatani Y, Qu W, Ahsan B, Yamada T,
315 Nagayasu Y, Doi K, Kasai Y: **The medaka draft genome and insights into**
316 **vertebrate genome evolution.** *Nature* 2007, **447(7145):714-719.**
- 317 26. Kesteven GL: **Whole-Genome Shotgun Assembly and Analysis of the**
318 **Genome of Fugu rubripes.** *Science (New York, NY)* 2002,
319 **297(5585):1301-1310.**
- 320 27. Jaillon O, Aury JM, Brunet F, Petit JL, Stangethomann N, Mauceli E,
321 Bouneau L, Fischer C, Ozoufcozaz C, Bernot A: **Genome duplication in the**
322 **teleost fish Tetraodon nigroviridis reveals the early vertebrate**
323 **proto-karyotype.** *Nature* 2004, **431(7011):946-957.**
- 324 28. Rondeau EB, Minkley DR, Leong JS, Messmer AM, Jantzen JR, von
325 Schalburg KR, Lemon C, Bird NH, Koop BF: **The genome and linkage map**
326 **of the northern pike (Esox lucius): conserved syntenic revealed between**
327 **the salmonid sister group and the Neoteleostei.** *PLoS ONE* 2014,
328 **9(7):e102089.**
- 329 29. Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford
330 R, Pirun M, Zody MC, White S: **The genomic basis of adaptive evolution in**
331 **threespine sticklebacks.** *Nature* 2012, **484(7392):55-61.**
- 332 30. Pevsner J: **Basic Local Alignment Search Tool (BLAST):** John Wiley &
333 Sons, Inc.; 2005.

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

- 376 46. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy**
377 **and high throughput**. *Nucleic acids research* 2004, **32**(5):1792-1797.
- 378 47. Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Hohna S,
379 Larget B, Liu L, Suchard MA, Huelsenbeck JP: **MrBayes 3.2: efficient**
380 **Bayesian phylogenetic inference and model choice across a large model**
381 **space**. *Systematic biology* 2012, **61**(3):539-542.
- 382 48. Brawand D, Wagner CE, Li YI, Malinsky M, Keller I, Fan S, Simakov O, Ng
383 AY, Lim ZW, Bezault E *et al*: **The genomic substrate for adaptive**
384 **radiation in African cichlid fish**. *Nature* 2014, **513**(7518):375-381.
- 385 49. Proost S, Fostier J, De Witte D, Dhoedt B, Demeester P, Van de Peer Y,
386 Vandepoele K: **i-ADHoRe 3.0--fast and sensitive detection of genomic**
387 **homology in extremely large data sets**. *Nucleic acids research* 2012,
388 **40**(2):e11.
- 389 50. Yang Z: **PAML: a program package for phylogenetic analysis by**
390 **maximum likelihood**. *Computer applications in the biosciences : CABIOS*
391 1997, **13**(5):555-556.

392
393

394 Tables

395 **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
Gap length (Mb)	122
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%

396
397
398
399
400

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

401

402

403 **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

404

405

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

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
	<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

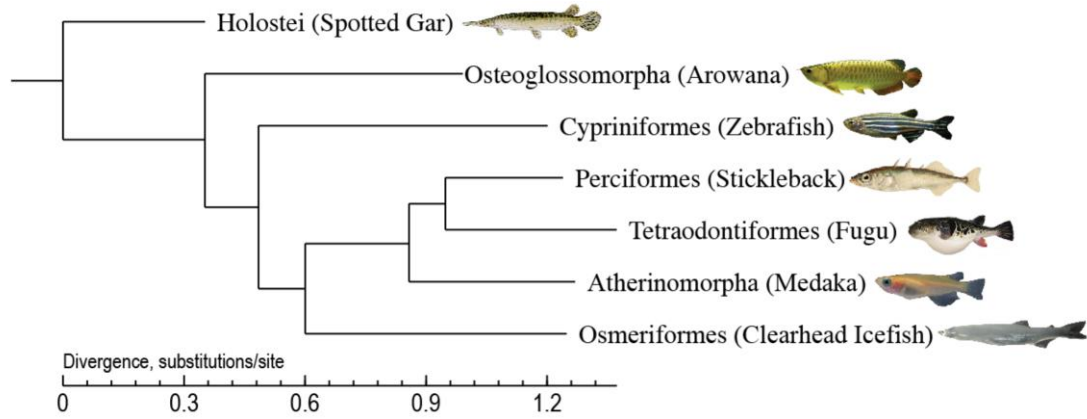
407

408

409

410

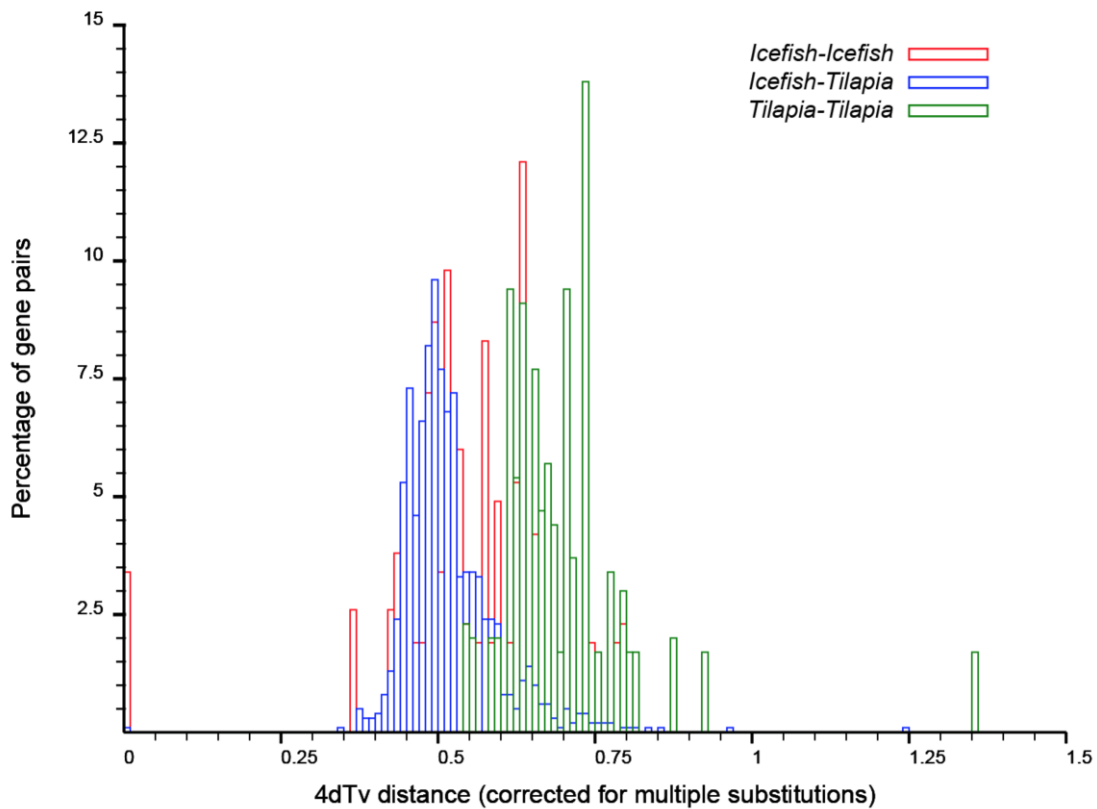
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



411

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

414



415

416 **Figure 2.** Distribution of 4DTV distances between the clearhead icefish and tilapia. The
 417 horizontal axis stands for the 4DTV distance corrected using the HKY model. The
 418 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