

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

14  
15 <sup>3</sup>Institute of Oceanology & Marine Fisheries, Jiangsu, 226007, China

16  
17 <sup>4</sup>BGI Zhenjiang Institute of Hydrobiology, Zhenjiang, 212000, China

18  
19 <sup>5</sup>BGI Research Center for Aquatic Genomics, Chinese Academy of Fishery Sciences,  
20 Shenzhen, 518083, China

21  
22 <sup>6</sup>Laboratory of Aquatic Genomics, College of Ecology and Evolution, School of Life  
23 Sciences, Sun Yat-Sen University, Guangzhou, 510275, China

24  
25 † Equal contributors

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 \* Correspondence: xup@ffrc.cn (PX); shiqiong@genomics.cn (QS);  
2  
3 27 13906292412@139.com (ZZ)  
4  
5 28  
6  
7 29 Email addresses: liuk@ffrc.cn (KL); xudp@ffrc.cn(DX); lijial@genomics.cn (JL);  
8  
9 30 bianchao@genomics.cn (CB); duanjr@ffrc.cn(JD); zhouyf@ffrc.cn(YZ);  
10  
11 31 zhangmy@ffrc.cn(MZ); youxinxin@genomics.cn (XY); youy@ffrc.cn (YY);  
12  
13 32 chenjieming@genomics.cn (JC); yuhui@genomics.cn (HY); xugc@ffrc.cn(GX);  
14  
15 33 fangda@ffrc.cn(DF); qiangj@ffrc.cn(JQ); 420219380@qq.com(SJ); hej@ffrc.cn(JH);  
16  
17 34 xujunmin@genomics.cn (JX); shiqiong@genomics.cn (QS); 13906292412@139.com  
18  
19 35 (ZZ); xup@ffrc.cn (PX)  
20  
21  
22  
23  
24  
25

## 26 38 **Abstract**

27 39 **Background:** Chinese clearhead icefish, *Protosalanx hyalocranius*, is a  
28  
29 40 representative species of icefishes with economic importance and special appearance.  
30  
31 41 Due to its great economic value in China, the fish was introduced to Lake Taihu and  
32  
33 42 several other lakes half a century ago. Similar to the *Sinocyclocheilus* cavefishes, the  
34  
35 43 clearhead icefish also has certain cavefish-like traits, such as transparent body and  
36  
37 44 nearly scaleless skin. Here, we provided the whole genome sequence of this  
38  
39 45 surface-dwelling fish and generated a high-quality genome assembly, aiming at  
40  
41 46 exploring molecular mechanisms for these biological characteristics.  
42  
43  
44  
45

46 47 **Findings:** A total of 252.1 gigabases (Gb) of raw reads were sequenced. Subsequently,  
47  
48 48 a novel high-quality genome assembly was generated, with the scaffold N50 reaching  
49  
50 49 1.163 Mb. The genome completeness was estimated to be 98.39% by using CEGMA  
51  
52 50 and BUSCO evaluation. Finally, we annotated 19,884 protein-coding genes and  
53  
54 51 observed that repeat sequences account for 24.43% of the genome assembly.  
55

56 52 **Conclusion:** We report the whole genome sequencing of the Chinese clearhead  
57  
58 53 icefish. The assembled genome will provide a significant foundation for further  
59  
60  
61  
62  
63  
64  
65

1 54 molecular breeding and germplasm resource protection in the clearhead icefish, as  
2 55 well as other icefishes. It is also a valuable genetic resource for revealing the  
3  
4 56 molecular mechanisms for the cavefish-like characteristics.  
5  
6  
7 57

8  
9 58 **Keywords:** Icefish; *Protosalanx hyalocranius*; Whole genome sequencing;  
10  
11 59 Genome assembly; Gene prediction; Repetitive sequences  
12  
13 60

## 16 61 **Data description**

### 18 62 ***Background***

20 63 Icefishes (Osmeriformes, Salangidae) are widely distributed in freshwater, coastal and  
21 64 estuarine habitats in East Asian countries [1-3]. Chinese clearhead icefish  
22  
23 65 (*Protosalanx hyalocranius*), a diadromous fish, mainly inhabits in coastal areas and  
24  
25 66 adjacent freshwaters [4-6]. As a commercially important fish in China, the clearhead  
26  
27 67 icefish was widely introduced into some lakes half a century ago and has developed a  
28  
29 68 resident life history [2, 7, 8]. Because of its transparent body and nearly scaleless skin,  
30  
31 69 similar to the *Sinocyclocheilus* cavefishes [9], we are very interested in this  
32  
33 70 surface-dwelling fish and are performing comparative genomics studies to explore the  
34  
35 71 mechanisms for these biological phenotypes. However, with the rapid development of  
36  
37 72 the Chinese economy in recent decades, population size of the clearhead icefish has  
38  
39 73 been seriously declining because of overfishing, construction of water conservancy  
40  
41 74 facilities and water pollution in the ecological systems [10]. To maintain its  
42  
43 75 sustainable development in China, here we performed the genome sequencing of  
44  
45 76 Chinese clearhead icefish for its biological and economic importance.  
46  
47  
48  
49  
50  
51  
52

### 53 78 ***Sample and Sequencing***

54  
55 79 In this study, we applied Illumina whole genome sequencing strategy to generate the  
56  
57 80 genome of Chinese clearhead icefish (NCBI Taxonomy ID: 418454; Fishbase ID:  
58  
59 81 12236). Genomic DNAs were isolated from the muscle tissue of an individual  
60  
61  
62  
63  
64  
65

1 82 collected from the Lake Taihu of Jiangsu Province in China. We constructed seven  
2 83 paired-end libraries with three short-insert libraries (250, 500 and 800 bp) and four  
3  
4 84 long-insert libraries (2, 5, 10 and 20 kb) using the standard protocol provided by  
5  
6 85 Illumina (San Diego, USA). Subsequent paired-end sequencing was performed by the  
7  
8 86 Illumina HiSeq 2000 platform for each library. Finally, we obtained 252.1 Gb of raw  
9  
10 87 125-bp reads for further analysis.  
11  
12  
13  
14

### 15 89 *Genome size estimation and genome assembly*

16 90 The SOAPfilter v2.2 software [11] with optimized parameters (-y -p -g 1 -o clean -M  
17 91 2 -f 0) was performed to remove low-quality raw reads (including reads with 10 or  
18  
19 92 more Ns and low-quality bases) and PCR-replicates as well as adaptor sequences. In  
20  
21 93 total, we obtained 169.0 Gb of clean reads. Subsequently, we estimated the genome  
22  
23 94 size based on the 17-mers depth frequency distribution method [12]. A 17-mer  
24  
25 95 represents an artificial division with 17-bp length nucleotide segment of sequencing  
26  
27 96 reads, therefore, a raw sequence read with a total length of L bp contains (L-17+1)  
28  
29 97 17-mers. The genome size was estimated with the following formula:  $G =$   
30  
31 98  $N*(L-17+1)/K$ -depth, in which G is the genome size, N is the total number of reads,  
32  
33 99 and K-depth is the highest frequency of 17-mer analysis. In our current study, N was  
34  
35 100 10,500,000,000 and the K-depth was 20. Hence, we estimated that the genome size of  
36  
37 101 Chinese clearhead icefish is 525 Mb.  
38  
39  
40  
41  
42

43 102 The filtered reads were assembled using SOAPdenovo2 v2.04.4 software [13] with  
44  
45 103 optimized parameters (pregraph -K 79 -d 1; contig -M 1; scaff -F -b 1.5 -p 16) to  
46  
47 104 generate contigs and original scaffolds. The gaps were fulfilled using GapCloser  
48  
49 105 v1.12 software [14] with default parameters and -p set to 25. Finally, we generated a  
50  
51 106 high-quality genome assembly of 536 Mb, with the scaffold N50 reaching 1.163 Mb  
52  
53 107 (Table 1).  
54  
55

56 108 The completeness of our assembly was evaluated by using CEGMA [15] and BUSCO  
57  
58 109 [16]. The CEGMA program (Core Eukaryotic Genes Mapping Approach; version 2.4)  
59  
60  
61  
62  
63  
64  
65

1 110 assessment with 248 conserved Core Eukaryotic Genes (CEGs) was performed for  
2 111 evaluation of the gene space completeness. The results revealed that the assembled  
3 112 genome had a CEGMA completeness score about 90.32% and 98.39%, which was  
4 113 calculated from the complete gene set and the partial gene set respectively.

5 114 Meanwhile, we used the representative metazoa gene set [17], which contains 843  
6 115 single-copy genes that are widely present in metazoa. The assessment demonstrated  
7 116 that the BUSCO values is 89%, containing [D: 10%], F: 7.7%, M: 2.9%, n: 843 (C:  
8 117 complete [D: duplicated], F: fragmented, M: missed, n: genes). These data from  
9 118 CEGMA and BUSCO indicate that the assembled genome covered majority of the  
10 119 gene space.

11 120

### 12 121 ***Repeat annotation***

13 122 Firstly, a *de novo* repeat library was constructed by the RepeatModeller v1.05 [18]  
14 123 and LTR\_FINDER.x86\_64-1.0.6 [10] with default parameters. Then, our assembly  
15 124 genome sequences were aligned against the ReBase v21.01 [19] and the *de novo*  
16 125 repeat libraries to recognize the known and novel TEs (transposable elements) using  
17 126 the RepeatMasker v4.06 [20]. Meantime, the Tandem Repeat Finder v4.07 [21] with  
18 127 parameters “Match=2, Mismatch=7, Delta=7, PM=80, PI=10, Minscore=50, and  
19 128 MaxPerid=2000” was utilized to annotate tandem repeats. Furthermore, the  
20 129 RepeatProteinMask software v4.0.6 [20] was used to predict TE relevant proteins in  
21 130 our genome assembly. Finally, we observed that the repeat sequences account for  
22 131 24.43% of the assembled genome (Table 1).

23 132

### 24 133 ***Genome Annotation***

25 134 In brief, we utilized two different methods to predict total gene set of the clearhead  
26 135 icefish.

27 136 **1) *de novo* annotation.** The AUGUSTUS v2.5 [22] and GENSCAN v1.0 [23] were  
28 137 executed to *ab initio* predict genes within the assembled genome, with the repetitive

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

138 sequences masked as “N” in order to discard pseudo gene prediction. Those  
139 low-quality genes with short length (<150 bp), premature termination or  
140 frame-shifting were removed.

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

150 **3) Integration of annotation results.** To merge all results produced from the above  
151 methods, we employed the GLEAN [33] to generate a non-redundant and  
152 comprehensive gene set. Finally, the best hit of each protein was obtained through all  
153 protein sequences from the GLEAN results aligned to the databases of the SwissProt  
154 and TrEMBL [34] (Uniprot release 2011.06) by BlastP with an e-value of  $10^{-5}$ .  
155 Overall, we generated a final gene set with 19,884 genes for the Chinese clearhead  
156 icefish.

157 CEGMA was performed again to evaluate the coverage rate between KOG  
158 (EuKaryotic Orthologous Groups) genes predicted by CEGMA and the predicted total  
159 gene set. It demonstrates that the predicted gene set mapped 96.4% of the KOGs.  
160 Simultaneously, the BUSCO was implemented again to assess the completeness of  
161 the predicted gene set. The BUSCO values were calculated as follows: C: 79% [D:  
162 16%], F: 9.8%, M: 10%, n: 843 (C: complete [D: duplicated], F: fragmented, M:  
163 missed, n: genes). The assessment values from both CEGMA and BUSCO proved  
164 high accuracy of the annotation.

1 165 **4) Function annotation.** The predicted protein sequences of clearhead icefish were  
2  
3 166 aligned against several public databases (Pfam [35], PRINTS [36], ProDom [37] and  
4  
5 167 SMART [38]) for detection of functional motifs and domains . Finally, we found that  
6  
7 168 96.2% of the predicted total gene set had been annotated with at least one functional  
8  
9 169 assignment from other public databases (Swiss-Prot [39], Interpro [40], TrEMBL [41]  
10  
11 170 and KEGG [42]).

12  
13 171

## 14 15 172 **Conclusion**

16  
17 173 We generated a high-quality genome assembly of Chinese clearhead icefish. The  
18  
19 174 novel genome data were deposited in publicly accessible repositories to promote  
20  
21 175 further biological research, molecular breeding and resource protection of this  
22  
23 176 representative icefish.

24  
25  
26 177

## 27 28 29 178 **Availability of supporting data**

30  
31 179 Supporting data are available in the GigaDB database [cite when ready], and the raw  
32  
33 180 whole genome sequences are deposited in the SRA under bioproject number  
34  
35 181 PRJNA328051.

36  
37  
38 182

## 39 40 183 **Competing interests**

41  
42 184 The authors declare that they have no competing interests.

43  
44  
45 185

## 46 47 186 **Funding**

48  
49 187 This study was supported by a grant from fish investigation in Taihu Lake (No.  
50  
51 188 TH2016WT007), National Infrastructure of Fishery Germplasm Resources (No.  
52  
53 189 2016DKA30470), Basic Research Funds from Freshwater Fisheries Research Center  
54  
55 190 (No. 2013JBFM07), Special Project on the Integration of Industry, Education and  
56  
57  
58 191 Research of Guangdong Province (No. 2013B090800017), Shenzhen Special Program

59  
60  
61  
62  
63  
64  
65

192 for Future Industrial Development (No. JSGG20141020113728803), and Zhenjiang  
193 Leading Talent Program for Innovation and Entrepreneurship.

194

## 195 **Author's Contributions**

196 KL, PX, QS, DX, JX, CB and ZZ conceived the project. MZ, XY, HY, JC, GX, DF,  
197 JQ, SJ and JH collected the samples and extracted the genomic DNA. JL, CB and HY  
198 performed the genome assembly and data analysis. KL, XP, JL, CB, QS, KL, YY and  
199 ZZ wrote the paper.

200

## 201 **References**

- 202 1. Wang ZS, Cui Zhang FU: **Biodiversity of Chinese Icefishes (Salangidae)**  
203 **and their conserving strategies**. *Chinese Biodiversity* 2002, **10**(4):416-424.
- 204 2. Zhang J, Li M, Xu M, Takita T, Wei F: **Molecular phylogeny of icefish**  
205 **Salangidae based on complete mtDNA cytochrome b sequences, with**  
206 **comments on estuarine fish evolution**. *Biological Journal of the Linnean*  
207 *Society* 2007, **91**(2):325-340.
- 208 3. Wang Z, Lu C, Hu H, Xu C, Lei G: **Dynamics of Icefish (Salangidae) Stocks**  
209 **in Nanyi Lake, Eastern China: Degradation and Overfishing**. *Journal of*  
210 *Freshwater Ecology* 2004, **19**(2):271-278.
- 211 4. Xia DQ, Cao Y, Ting ting WU, Yang H: **Study on lineages of Protosalanx**  
212 **chinensis, Neosalanx taihuensis and N.oligodontis in Taihu Lake with**  
213 **RAPD technique**. *Journal of Fisheryences of China* 2000, **7**(01):12-15.
- 214 5. Xia DQ, Cao Y, Ting Ting WU, Yang H: **Genetic Structures of Population**  
215 **of Protosalanx Chinensis, Neosalanx Taihuensis and Neosalanx**  
216 **Oligodontis in Lake Taihu**. *Journal of Fisheries of China* 1999(03):254-260.
- 217 6. Armani A, Castigliego L, Tinacci L, Gianfaldoni D, Guidi A: **Molecular**  
218 **characterization of icefish, (S alangidae family), using direct sequencing**  
219 **of mitochondrial cytochrome b gene**. *Food Control* 2011, **22**(6):888-895.
- 220 7. Wang Z, Lu C, Hu H, Zhou Y, Xu C, Lei G: **Freshwater icefishes**  
221 **(Salangidae) in the Yangtze River basin of China: Spatial distribution**  
222 **patterns and environmental determinants**. *Environmental Biology of Fishes*  
223 2005, **73**(3):253-262.
- 224 8. Ye S, Yang J, Liu H, Oshima Y: **Use of elemental fingerprint analysis to**  
225 **identify localities of collection for the large icefish protosalanx chinensis**



- 226 **in Taihu Lake, China.** *Journal of the Faculty of Agriculture, Kyushu*  
227 *University* 2011, **56**(1):41-45.
- 228 9. Yang J, Chen X, Jie B, Fang D, Ying Q, Jiang W, Hui Y, Chao B, Jiang L, He  
229 **S: The Sinocyclocheilus cavefish genome provides insights into cave**  
230 **adaptation.** *Bmc Biology* 2016, **14**(1):1-13.
- 231 10. Xu J, Xie P, Zhang M, Zhou Q, Zhang L, Wen Z, Cao T: **Icefish (salangidae)**  
232 **as an indicator of anthropogenic pollution in freshwater systems using**  
233 **nitrogen isotope analysis.** *Bulletin of environmental contamination and*  
234 *toxicology* 2007, **79**(3):323-326.
- 235 11. **SOAPfilter v2.2 software:** <http://soap.genomics.org.cn/index.html>.
- 236 12. Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, Li Z, Chen Y, Mu D, Fan W:  
237 **Estimation of genomic characteristics by analyzing k-mer frequency in de**  
238 **novo genome projects.** *Quantitative Biology* 2013, **35**(s 1–3):62-67.
- 239 13. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y *et*  
240 *al*: **SOAPdenovo2: an empirically improved memory-efficient short-read**  
241 **de novo assembler.** *Gigascience* 2012, **1**:18.
- 242 14. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J: **SOAP2: an**  
243 **improved ultrafast tool for short read alignment.** *Bioinformatics* 2009,  
244 **25**(15):1966-1967.
- 245 15. Parra G, Bradnam K, Korf I: **CEGMA: a pipeline to accurately annotate**  
246 **core genes in eukaryotic genomes.** *Bioinformatics* 2007, **23**(9):1061-1067.
- 247 16. Sim AFO, Waterhouse MR, Ioannidis P, Kriventseva VE, Zdobnov ME:  
248 **BUSCO: assessing genome assembly and annotation completeness with**  
249 **single-copy orthologs.** *Bioinformatics* 2015, **31**(19):3210-3212.
- 250 17. Kriventseva EV, Tegenfeldt F, Petty TJ, Waterhouse RM, Simão FA,  
251 Pozdnyakov IA, Ioannidis P, Zdobnov EM: **OrthoDB v8: update of the**  
252 **hierarchical catalog of orthologs and the underlying free software.** *Nucleic*  
253 *Acids Research* 2015, **43**(Database issue):D250-D256.
- 254 18. **RepeatModeller v1.05 software:**  
255 <http://www.repeatmasker.org/RepeatModeler.html>.
- 256 19. J J, VV K, A P, P K, O K, J W: **Rebase Update, a database of eukaryotic**  
257 **repetitive elements.** *Cytogenetic & Genome Research* 2005,  
258 **110**(1-4):462-467.
- 259 20. Chen N: **Using RepeatMasker to Identify Repetitive Elements in Genomic**  
260 **Sequences:** John Wiley & Sons, Inc.; 2004.
- 261 21. Benson G, . **Tandem repeats finder: a program to analyze DNA sequences.**  
262 *Nucleic Acids Research* 1999, **27**(2):573-580(578).
- 263 22. Mario S, Oliver K, Irfan G, Alec H, Stephan W, Burkhard M: **AUGUSTUS:**  
264 **ab initio prediction of alternative transcripts.** *Nucleic Acids Research* 2006,  
265 **34**(WebServerissue):435-439.
- 266 23. Burge C, ., Karlin S, . **Prediction of complete gene structures in human**  
267 **genomic DNA.** *Journal of Molecular Biology* 1997, **268**(1):78-94.

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

268 24. Collins JE, White S, Searle SMJ, Stemple DL: **Incorporating RNA-seq data**  
269 **into the zebrafish Ensembl genebuild.** *Genome Research* 2012,  
270 **22(10):2067-2078.**

271 25. Kasahara M, Naruse K, Sasaki S, Nakatani Y, Qu W, Ahsan B, Yamada T,  
272 Nagayasu Y, Doi K, Kasai Y: **The medaka draft genome and insights into**  
273 **vertebrate genome evolution.** *Nature* 2007, **447(7145):714-719.**

274 26. Kesteven GL: **Whole-Genome Shotgun Assembly and Analysis of the**  
275 **Genome of Fugu rubripes.** *Science (New York, NY)* 2002,  
276 **297(5585):1301-1310.**

277 27. Jaillon O, Aury JM, Brunet F, Petit JL, Stangethomann N, Mauceli E,  
278 Bouneau L, Fischer C, Ozoufcostaz C, Bernot A: **Genome duplication in the**  
279 **teleost fish Tetraodon nigroviridis reveals the early vertebrate**  
280 **proto-karyotype.** *Nature* 2004, **431(7011):946-957.**

281 28. Rondeau EB, Minkley DR, Leong JS, Messmer AM, Jantzen JR, von  
282 Schalburg KR, Lemon C, Bird NH, Koop BF: **The genome and linkage map**  
283 **of the northern pike (Esox lucius): conserved synteny revealed between**  
284 **the salmonid sister group and the Neoteleostei.** *PLoS ONE* 2014,  
285 **9(7):e102089.**

286 29. Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford  
287 R, Pirun M, Zody MC, White S: **The genomic basis of adaptive evolution in**  
288 **threespine sticklebacks.** *Nature* 2012, **484(7392):55-61.**

289 30. Pevsner J: **Basic Local Alignment Search Tool (BLAST):** John Wiley &  
290 Sons, Inc.; 2005.

291 31. Yu XJ, Zheng HK, Wang J, Wang W, Su B: **Detecting lineage-specific**  
292 **adaptive evolution of brain-expressed genes in human using rhesus**  
293 **macaque as outgroup.** *Genomics* 2006, **88(6):745-751.**

294 32. Birney E, Clamp M, Durbin R: **GeneWise and Genomewise.** *Genome*  
295 *Research* 2004, **14(5):988-995.**

296 33. Elsik CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM:  
297 **Creating a honey bee consensus gene set.** *Genome Biology* 2007,  
298 **8(1):90-105.**

299 34. Bairoch A, Apweiler R: **The SWISS-PROT protein sequence database and**  
300 **its supplement TrEMBL in 2000.** *Nucleic Acids Res* 2000, **28(1):45-48.**

301 35. Finn RD: **Pfam: the protein families database.** *Nucleic Acids Research* 2014,  
302 **42(Database issue):D222-230.**

303 36. Attwood TK: **The PRINTS database: A resource for identification of**  
304 **protein families.** *Briefings in Bioinformatics* 2002, **3(3):252-263.**

305 37. Bru C, Courcelle E, Beausse Y, Dalmar S, Kahn D: **The ProDom database of**  
306 **protein domain families: more emphasis on 3D.** *Nucleic Acids Research*  
307 2005, **33(Database issue):212-215.**

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

308 38. Letunic I, Copley RR, Schmidt S, Ciccarelli FD, Doerks T, Schultz J, Ponting  
309 CP, Bork P: **SMART 4.0: towards genomic data integration.** *Nucleic Acids*  
310 *Research* 2004, **32**(Databaseissue):D142-D144.

311 39. B B, A B, R A, MC B, A E, E G, MJ M, K M, C OD, I P *et al*: **The**  
312 **Swiss-Prot knowledgebase and its supplement TREMBL in 2003.** *Nucleic*  
313 *Acids Research* 2003, **31**(1):365-370.

314 40. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P,  
315 Das U, Daugherty L, Duquenne L: **InterPro: the integrative protein**  
316 **signature database.** *Nucleic Acids Research* 2009, **37**(suppl 1):D211-D215.

317 41. Hingamp P, Broek AEVD, Stoesser G, Baker W: **The EMBL nucleotide**  
318 **sequence database.** *Molecular Biotechnology* 1999, **12**(3):255-267.

319 42. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.**  
320 *Nucleic Acids Research* 2000, **27**(1):29-34(26).

321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340

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

341 **Tables**

342 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%

---

343