# GigaScience

## Genome sequence of rock bream, Oplegnathus fasciatus (Temminck & Schlegel, 1884): the first draft genome in family Oplegnathidae
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-18-00300 |
| Full Title: | Genome sequence of rock bream, Oplegnathus fasciatus (Temminck & Schlegel, 1884): the first draft genome in family Oplegnathidae |
| Article Type: | Data Note |
| Funding Information: | Shandong Province Key Research and Invention Program (2017GHY15102, 2017GHY15106) — Dr. Yongshuang Xiao<br>Young Scientists Fund (41506170) — Dr. Yongshuang Xiao |

**Abstract:**

Background
The rock bream (Oplegnathus fasciatus), a member of the Oplegnathidae family of the Perciformes, is a commerically important rocky reef fish native to East Asia. O. fasciatus has become an important fishery resource for offshore cage aquaculture and fish stocking of marine ranching in China, Japan and Korea. Recently, growth of sexual dimorphism with neo-sex chromosome and widespread biotic diseases in O. fasciatus has received increasing concern. However, the adequate genome resources to make insight into sex-determining mechanism and to establish genetically basing resistant breeding systems for O. fasciatus have been lacking. Here, we performed whole genome of female fish for O. fasciatus using long-read sequencing and Hi-C data to generate chromosome-length scaffolds with highly contiguous genome assembly.

Findings
We assembled the O. fasciatus with a total of 245.0 Gb of raw reads, which were generated using both of PacBio Sequel and Illumina Hiseq 2000 platforms. The final draft genome assembly was approximately 778.7 Mb, which reached a remarkable high level of continuity with contig N50 of 2.1 Mb. The genome size was consistent with the estimated genome size (808.9 Mb) based on k-mer analysis. The identified repeat sequences account for 32.2% of the whole genome and 24 003 protein-coding genes with an average of 10.1 exons per gene were annotated using de novo method and with RNA-seq data and homologies to other teleosts. We combined Hi-C data with draft genome assembly to generate chromosome-length scaffolds. Twenty-four scaffolds corresponding to the twenty-four chromosomes were assembled to a final size of 768.8 Mb with contig N50 2.1 Mb and scaffold N50 of 33.5 Mb using 1372 contigs. According to the phylogenetic analysis using protein-coding genes, the O. fasciatus was close related to Larimichthys crocea and O. fasciatus diverged from their ancestor was at about 70.3-87.3 million years ago.

Conclusions
We generated high-quality draft genome and chromosomes assembly for O. fasciatus using long reads generated using PacBio sequencing technologies, which is the first reference genome for Oplegnathidae species. The genome assembly will provide insight into sex-determining mechanism and serve as a resource for accelerating the genome-assisted improvement of resistant breeding systems.

| | |
|---|---|
| Corresponding Author: | Yongshuang Xiao<br><br>CHINA |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Yongshuang Xiao |

| First Author Secondary Information: | |
|---|---|
| Order of Authors: | Yongshuang Xiao |
| | Zhizhong Xiao |
| | Jing Liu |
| | Daoyuan Ma |
| | Jun Li |
| Order of Authors Secondary Information: | |
| Additional Information: | |
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |
| Experimental design and statistics<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| Resources<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| Availability of data and materials<br><br>All datasets and code on which the | Yes |

conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.

Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?

# Genome sequence of rock bream, *Oplegnathus fasciatus* (Temminck & Schlegel, 1884): the first draft genome in family Oplegnathidae

Yongshuang Xiao[1,2,3,†], Zhizhong Xiao[1,2,3,†], Jing Liu[2,4*], Daoyuan Ma[1,2,3*], Jun Li[1,2,3*]

[1]CAS Key Laboratory of Experimental Marine Biology, Institute of Oceanology, Chinese Academy of Sciences, 7 Nanhai Road, Qingdao, 266071, China, [2]Laboratory for Marine Biology and Biotechnology, Qingdao National Laboratory for Marine Science and Technology, 7 Nanhai Road, Qingdao, 266071, China, [3]Center for Ocean Mega-Science, Chinese Academy of Sciences, 7 Nanhai Road, Qingdao, 266071, China.

∗Correspondence address: Jing Liu, Institute of Oceanology, Chinese Academy of Sciences, 7 Nanhai Road, Qingdao, 266071, China; Tel: +86-053282898790; E-mail: jliu@qdio.ac.cn; Daoyuan Ma, Mega-Science, Chinese Academy of Sciences, 7 Nanhai Road, Qingdao, 266071, China; Tel: +86-053282898717; E-mail: madaoyuan1@163.com; Jun Li, Institute of Oceanology, Chinese Academy of Sciences, 7 Nanhai Road, Qingdao, 266071, China; Tel: +86-053282898718; E-mail: junli@qdio.ac.cn.

†Contributed equally to this work.

29  **Abstract**

30  **Background**

31  The rock bream (*Oplegnathus fasciatus*), a member of the Oplegnathidae family of

32  the Perciformes, is a commerically important rocky reef fish native to East Asia. *O.*

33  *fasciatus* has become an important fishery resource for offshore cage aquaculture and

34  fish stocking of marine ranching in China, Japan and Korea. Recently, growth of

35  sexual dimorphism with neo-sex chromosome and widespread biotic diseases in *O.*

36  *fasciatus* has received increasing concern. However, the adequate genome resources

37  to make insight into sex-determining mechanism and to establish genetically basing

38  resistant breeding systems for *O. fasciatus* have been lacking. Here, we performed

39  whole genome of female fish for *O. fasciatus* using long-read sequencing and Hi-C

40  data to generate chromosome-length scaffolds with highly contiguous genome

41  assembly.

42  **Findings**

43  We assembled the *O. fasciatus* with a total of 245.0 Gb of raw reads, which were

44  generated using both of PacBio Sequel and Illumina Hiseq 2000 platforms. The final

45  draft genome assembly was approximately 778.7 Mb, which reached a remarkable

46  high level of continuity with contig N50 of 2.1 Mb. The genome size was consistent

47  with the estimated genome size (808.9 Mb) based on *k*-mer analysis. The identified

48  repeat sequences account for 32.2% of the whole genome and 24 003 protein-coding

49  genes with an average of 10.1 exons per gene were annotated using *de novo* method

50  and with RNA-seq data and homologies to other teleosts. We combined Hi-C data

51  with draft genome assembly to generate chromosome-length scaffolds. Twenty-four

52  scaffolds corresponding to the twenty-four chromosomes were assembled to a final

53  size of 768.8 Mb with contig N50 2.1 Mb and scaffold N50 of 33.5 Mb using 1372

54  contigs. According to the phylogenetic analysis using protein-coding genes, the *O.*

55  *fasciatus* was close related to *Larimichthys crocea* and *O. fasciatus* diverged from

56  their ancestor was at about 70.3-87.3 million years ago.

57  **Conclusions**

58  We generated high-quality draft genome and chromosomes assembly for *O. fasciatus*

59 using long reads generated using PacBio sequencing technologies, which is the first

60 reference genome for Oplegnathidae species. The genome assembly will provide

61 insight into sex-determining mechanism and serve as a resource for accelerating the

62 genome-assisted improvement of resistant breeding systems.

63 *Keywords*: *Oplegnathus fasciatus*; genome assembly; Hi-C assembly; sex-determining

64 mechanism

## Data description

### Introduction of *O. fasciatus*

67 The family Oplegnathidae belongs to the order Perciformes, including only one genus

68 *Oplegnathus* comprised of two species, *O. fasciatus* and *O. punctatus* that are of

69 commercial values. The rock bream, *O. fasciatus* (Temminck & Schlegel, 1844), is

70 one of the two species in the *Oplegnathus*, which is commonly found at the depth of

71 one to ten meters in association with rocky reefs[1, 2], being distributed in a wide range

72 of shallow waters around Korea, Japan, China and Hawaii[1, 3, 4] (Fig. 1). *O. fasciatus*

73 has become an important fishery resource for offshore cage aquaculture and fish

74 stocking of marine ranching in China, Japan and Korea[5]. It was reported that the male

75 of *Oplegnathus* has a neo-sex chromosome, possibly a sex chromosome Y, and the sex

76 chromosome system for *Oplegnathus* was considered to be $X_1 X_1 X_2 X_2 / X_1 X_2 Y$

77 based on the karyotype analyses[6, 7]. Furthermore, the growth sexual dimorphism was

78 detected in the *O. fasciatus* and the male fish showed a faster growth advantage than

79 the female, may be due to the sex chromosome system of *Oplegnathus*[8]. *O. fasciatus*

80 is vulnerable to viruses (eg. Iridovirus) and genetic degradation caused by inbreeding

81 has led to higher susceptibility to diseases[9, 10]. It is vital to develop genomic resources

82 for making insight into sex-determining mechanism and accelerating the

83 genome-assisted improvement of resistant breeding systems.

84 So far, the genome sequence and the chromosomes assembly of *O. fasciatus*

85 have not been reported. Here we performed a high-quality reference genome assembly

86 for *O. fasciatus* constructed using long reads by the PacBio DNA sequencing platform,

87 and using a genome assembly strategy by taking advantage of genome assemblyer

88 Canu[11]. The genome assembly of *O. fasciatus* is the first reference genome

89 constructed for the family Oplegnathidae. The completeness and continuity of the

90 genome will provide high quality genomic resources for studies on sex-determining

91 mechanism and for accelerating the genome-assisted improvement of resistant

92 breeding systems.

93

94 **Genomic DNA extraction, genome size estimation and Hi-C library construction**

95 High-quality genomic DNA for Illumina platform (Illumina Inc., San Diego, CA,

96 USA) and PacBio Sequel sequencing (Pacific Biosciences of California, Menlo Park,

97 CA, USA) was extracted from fresh muscle tissue and blood sample of a single

98 female *O. fasciatus*. The fish was collected from the near-shore area of Qingdao city

99 (Yellow Sea), Shandong province. A whole-genome using Illumina DNA sequencing

100 technology was applied to estimate *O. fasciatus* genome size. A short-insert library

101 (300~350 bp) was constructed and generated a total of ~90.7 Gb of raw reads using

102 the standard protocol provided by Illumina Hieq X Ten platform (Illumina Inc., San

103 Diego, CA, USA). After removal of low-quality and redundant reads, we obtained

104 about ~80.8 Gb of clean data for *de novo* assembly to estimate the genome size (S

105 Table 1, Fig. 2). All the cleaned reads were subjected to 17-mer frequency distribution

106 analysis[12]. As the total number of *k*-mers was about $8.09 \times 10^{10}$ and the peak of *k*-mers

107 at a depth of 100, the genome size of *O. fasciatus* was calculated to be 808.9 Mb

108 using the following formula: genome size = *k*-mer number / peak depth (Fig. 2).

109 Meanwhile, the estimated heterozygosity of 0.29% and a repeat content of 38.46%

110 were detected for *O. fasciatus* in this work. A pilot genome assembly was

111 approximately 808.9 Mb with a contig N50 7.2 kb and scaffold N50 84.1kb using the

112 Illumina data and the assembly program Platanus package[13] (S Table 2). The GC

113 content was 41% (S Fig. 1). This genome assembly was of low-quality partly due to

114 its high genomics repeat content.

115 The genomic DNA for Hi-C library was extracted from the whole-blood cell of

116 *O. fasciatus* as described[14]. The cells were fixed with formaldehyde and lysed, and the

117 cross-linked DNA digested with MboI. Sticky ends were biotin-labeled and proximity

118 ligated to form chimeric junctions that were enriched for and then physically sheared

119 to a size of 300–500 bp[14]. Chimeric fragments representing the original cross-linked

120 long-distance physical interactions were then processed into paired-end sequencing

121 libraries, and 629 million 150-bp paired-end Illumina reads (91.5 Gb) with Q20 and

122 Q30 of ~94.0% were produced (S Table 1, S Table 3). As a result, the paired data, data

123 with mate mapped to a different contig (or scaffold) and data with mapped to a

124 different contig (or scaffold) (map $Q5 \geqslant 5$) were 593.7 Mb (94.4%), 240.5 Mb

125 (40.5%) and 205.1 Mb (34.6%), respectively (S Table 3).

126

## Genome assembly using PacBio long reads

128 Two 20 kb genomic DNA libraries were constructed and sequenced using PacBio

129 Sequel platform, generating 62.9 Gb raw DNA reads. We obtain 4.8 million subreads

130 (totally 62.8 Gb) with a read N50 length of ~22 kb after removing adaptor (S Table

131 1).

132 The Canu v1.4 was firstly used to assemble the genome with the

133 Corrected-Error-Rate parameter set at 0.040[11]. As a result, a total length of 875.9 Mb

134 genome assembly was achieved for *O. fasciatus*, which was consistent with the

135 estimated genome size in 17-mer analysis based on the Illumina data (S Table 2). We

136 applied Redundans v0.13c[15] to remove the sequence redundancy and obtain genome

137 assembly size of 778.0 Mb. We then used the Arrow of Smrtlink 5.0 with the

138 minCoverage parameter set at 15 to implement the error correction based on the

139 PacBio long reads data (Table 1). The resulting genome assembly was further

140 polished using NGS data, which were used in the genome survey analysis above. The

141 final draft genome assembly was 778.7 Mb, which reached a remarkable high level of

142 continuity with contig N50 length of 2.1 Mb (Table 1). The contig N50 of *O. fasciatus*

143 was much higher than those of previous fish genome assemblies constructed using

144 NGS DNA sequencing technologies, and is comparable with those of recently

145 reported model fish species (S Table 4)

146

## Genome quality evaluation

148 To assess the completeness of the assembled *O. fasciatus* genome, we subjected the

149 sequences to BUSCO version 3 evaluation (BUSCO, actinopterygii_odb9) [16]. Overall,

150 96.6% and 1.5 % of the 4 584 expected actinopterygii genes were identified in the

151 assembled genome as complete and partial BUSCO profiles, respectively.

152 Approximately 85 genes could be considered missing in our assembly (S Table 5).

153 Among the expected complete actinopterygii genes, both of 4 259 and 171 were

154 identified as single copy and duplicated BUSCOs, respectively (S Table 5). We then

155 used the Minimap2 to estimate the completeness and homogeneity of genome

156 assembly based on the CLR (Continuous Long Reads) subreads. A high quality of

157 completeness and homogeneity was checked for genome assembly, and the mapping

158 rate, coverage rate and average sequencing depth were reached to 90.2%, 99.9% and

159 80.6, respectively (S Table 6).

160 To further evaluate the accuracy of *O. fasciatus* genome assembly, we aligned

161 the NGS-based short reads from whole-genome sequencing data against the reference

162 genome using BWA[17]. We then used the GATK to implement the SNP calling and

163 filter work, the result showed 99.8% and 0.2 % of the $1.6 \times 10^6$ expected SNP reads

164 were identified in the assembled genome as heterozygosis and homology SNPs,

165 respectively. SNP calling on the final assembly also yield a heterozygosity rate of

166 0.20%, supporting the estimate from *k*-mer analysis (0.29%) (S Table 7).

167

168 **Repeat sequence within th**e *O. fasciatus* g**enome assembly**

169 To identify tandem repeats, we utilized Tandem Repeat Finder to annotate repetitive

170 elements in the *O. fasciatus* genome. RepeatModeler (version 1.04) and

171 LTR_FINDER[18] were used to construct a *de novo* repeat library with default

172 parameters. Subsequently, we used RepeatMasker[19] (version 3.2.9) to map our

173 assembled sequences on the Repbase TE (version 14.04) [20] and the *de novo* repeat

174 library to identify known and novel transposable elements (TEs). In addition, the

175 TE-related proteins were annotated by using RepeatProteinMask software (version

176 3.2.2) [19].

177 The total identified repeat sequences accounted for 23.6% of the *O. fasciatus*

178     genome based on the *de novo* repeat library (Table 2). Approximately 23.41% of the

179     *O. fasciatus* genome was identified as interspersed repeats (most often TEs). Among

180     them, DNA transposable elements were the most abundant type of repeat sequences,

181     which occupied 11.5% of the whole genome. The long interspersed nuclear elements

182     (LINE) and long terminal repeat (LTR) took up 7.3% and 4.0% of the whole genome

183     (Table 2, S Fig. 2).

184     **RNA preparation and sequencing**

185     We sequenced cDNA libraries prepared from the eggs of *O. fasciatus* used for genome

186     annotation using Illumina sequencing technologies. High quality of RNA were

187     detected based on the estimation of the absorbance at 260nm / 280nm (OD = 2.0) and

188     the RIN (value = 9.2) by Nanodrop ND-1000 spectrophotometer (LabTech, USA) and

189     2100 Bioanalyzer (Agilent Technologies, USA), respectively. We used the Clontech

190     SMARTer cDNA synthesis kit to complete the process of reverse transcription. The

191     paired-end library was prepared following the manual of the Paired-End Sample

192     Preparation Kit (Illumina Inc., San Diego, CA, USA). Finally, the library with an

193     insert length of 300 bp was sequenced by Illumina HiSeq X Ten in 150PE mode

194     (Illumina Inc., San Diego, CA, USA). As a result, we obtained ~42.2 Gb high-quality

195     transcriptome data from RNA-seq (S Table 1, S Table 8)

196     **Gene annotation**

197     Gene annotation of the *O. fasciatus* genome was performed using *de novo*,

198     homology-based and transcriptome sequencing-based prediction. We employed

199     Augustus (version 2.5.5)[21] and GenScan (version 1.0)[22] softwares to predict

200     protein-coding genes of *O. fasciatus* genome assembly. Protein sequences of closely

201     related fish species including *Larimichthys crocea*，*Lates calcarifer*，*Gasterosteus*

202     *aculeatus*，*Paralichthys olivaceus*，*Cynoglossus semilaevis* and *Gadus morhua* were

203     downloaded from Ensembl[23] and aligned against to *O. fasciatus* genome using

204     TBLASTN software[24]. Subsequently, Genewise2.2.0 software[25] was employed to

205     predict the potential gene structures on all alignments.

206        We also mapped these NGS transcriptome short reads onto our genome assembly

207     using TopHat1.2 software[26], and then we employed Cufflinks[27] to predict the gene

208 structures (S Table 9). All gene models were then integrated using MAKER to obtain

209 a consensus gene set[28]. The final total gene set was composed of 24 003 genes, with

210 an average of 10.1 exons per gene in *O. fasciatus* genome (Table 1). The gene number,

211 gene length distribution, CDS length distribution, exon length distribution and intron

212 length distribution were all comparable with those in other teleost fish species (S

213 Table 9, S Fig. 3).

214     In order to further obtain functional annotation of the protein-coding genes in *O.*

215 *fasciatus* genome, we employed local BLASTX and BLASTN programs to align upon

216 the non-redundant protein (NR), non-redundant nucleotide (NT) and Swissprot

217 database with an e-value $\leqslant$ 1e-5[29]. We also used Blast2GO software to search the

218 Gene ontology (GO), Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway

219 database[30, 31, 32]. Finally, 97.3% (23 364 genes) of the 24 003 genes were annotated by

220 at least one database (S Table10). Four types of non-coding RNAs (microRNAs,

221 transfer RNAs, ribosomal RNAs, and small nuclear RNAs) were also annotated using

222 tRNAscan-SE and the Rfam database in this study[33, 34] (S Table11).

223 **Hi-C assembly and chromosome interactions**

224 Hi-C was a sequencing-based approach for determining chromosome interactions by

225 calculating the contact frequency between pairs of loci, which depended strongly on

226 the one-dimensional distance, in base pairs, between a pair of loci[35, 36]. We employed

227 BWA and Lachesis softwares to align paired-end reads to the draft genome assembly

228 and filtered all base sequences other than 500bp from each restriction site[37].

229 According to the conduct of clustering, ordering, and orienting to the assembly

230 contigs (1 692), those were grouped into 24 chromosome clusters and scaffolded

231 using Lachesis software with tuned parameters[38] (Table 3, Fig. 3). Finally, we

232 constructed the chromosome interactions map using Juicer software and employed the

233 JucieBox to complete the visual correction of interactions map. We obtained polished

234 1 756 contigs by interrupting misassembly from the 1 692 contigs. Twenty-four

235 scaffolds corresponding to the 24 chromosomes of *O. fasciatus* based on the

236 karyotype analyses were assembled[6, 7] (Table 3, Fig. 3). A final size of 768.8 Mb

237 accounting for the 98.7% draft genome was assembled, which remarkable high level

238   of continuity with contig N50 of 2.1 Mb and scaffold N50 of 33.5 Mb using 1372

239   contigs. The anchor rate of contigs (> 100 kb) to chromosomes was reached up to the

240   99.7% based on the Hi-C assembly (Table 4). The contig N50 and scaffold N50 of *O.*

241   *fasciatus* were much higher than those of previous fish genome assemblies

242   constructed using NGS DNA sequencing technologies based on the genome assembly

243   using PacBio long reads and Hi-C assembly (S Table 4).

**Gene family identification and phylogenetic tree construction**

245   We employed the BLASTP program[39] with an e-value threshold of 1e-5 to identify

246   gene family based on the transcripts alignments of each gene from *O. fasciatus* and

247   other fish species, which included *Larimichthys crocea*, *Gadus morhua*, *Paralichthys*

248   *olivaceus*, *Cynoglossus semilaevis*, *Notothenia coriiceps*, *Boleophthalmus*

249   *pectinirostris*, *Branchiostoma floridae*, *Gasterosteus aculeatus*, *Callorhinchus milii*,

250   *Danio rerio*, *Salmo salar* and *Oryzias latipes*. 23273 gene families were identified by

251   clustering of homologous gene sequences based on H-scores calculated from

252   Bit-score in Hcluster_sg software (S Fig. 4). Subsequently, we selected 812

253   single-copy orthogroups from the above-mentioned species to construct the

254   phylogenetic relationship between *O. fasciatus* and the other fish species. We used the

255   Clustal W program[40] to extract and align coding sequences of single-copy gene from

256   the 765 orthogroups with length filter, respectively (S Fig. 5). All the alignments were

257   concatenated as a single data set for each species. Nondegenerated sites extracted

258   from the data set were then joined into new sequence of each species to construct a

259   phylogenetic tree based on the maximum-likelihood method implemented in the

260   PhyML package[41] (with the -m PROTGAMMAAUTO model). We used the

261   MCMCtree program to estimate divergence times among species based on the

262   approximate likelihood method[42] and a molecular clock data from the divergence time

263   between medaka from the TimeTree database[43]. According to the phylogenetic

264   analysis *O. fasciatus* were clustered together with *Larimichthys crocea* belonged to

265   the order Perciformes, which was consistent with the fish species taxonomy. The

266   taxonomy of Notothenioidei should be elevated to the order level from the

267   Perciformes and be paralleled with Gasterosteiformes (Fig. 4). The divergence time

268 between *O. fasciatus* and the common ancestor with *Larimichthys crocea* was at about

269 70.3-87.3 Ma.

**Conclusion**

271 We successfully assembled the genome of *O. fasciatus* and reported the first whole

272 genome sequencing, assembly and annotation based on long reads from the

273 third-generation PacBio Sequel sequencing platform. The final draft genome

274 assembly is approximately 778.7 Mb, accounting for 96.3% of the estimated genome

275 size (808.9 Mb) based on *k*-mer analysis. The genome assembly of *O. fasciatus* was

276 also the first high-quality genome of all species in Oplegnathidae family, which

277 reached a remarkable high level of continuity with contig N50 of 2.1 Mb and scaffold

278 N50 of 33.5 Mb. The contig N50 was remarkably longer than those of most fish

279 genome assemblies, and was comparable with those of recently reported model fish

280 species. We also predicated 24 003 protein-coding genes from the generated assembly,

281 and 97.3% (23 364 genes) of all protein-coding genes were annotated. Twenty-four

282 scaffolds corresponding to the twenty-four chromosomes were assembled to a final

283 size of 768.8 Mb using 1372 contigs based on the Hi-C assembly. We found the

284 taxonomy of Notothenioidei should be elevated to the order level and the divergence

285 time between *O. fasciatus* and the common ancestor with *Larimichthys crocea* was at

286 about 70.3-87.3 Ma. The genome assembly, together with gene annotation data

287 generated in this work provided a valuable resource for research on sex-determining

288 mechanism and for accelerating the genome-wide association studies on resistant

289 breeding systems.

290

**Ethics Statement**

292 This research was approved by the Animal Care and Use committee of Chinese

293 Academic Science. All participates consent the study under the 'Ethics, consent and

294 permissions' heading. All participants consent to publish the work under the 'Consent

295 to publish' heading.

296

**Availability of supporting data**

The numbers on the left margin run vertically: 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

298 Supporting data and materials are available in the GigaScience GigaDB database,

299 with the raw sequences deposited in the SRA under the accession number

300 SRP158313.

301

302 **Competing interests**

303 The authors declare that they have no competing interests.

304

313

314 **Author Contributions**

315 YSX conceived the project. ZZX, DYM collected the samples and extracted the

316 genomic DNA. YSX, JL and JL performed the genome assembly and data analysis.

317 YSX, ZZX, JL, DYM and JL wrote the paper.

318

319

320 **Reference**

321 1.  Schembri, P.J. *et al*. Occurrence of barred kinfejaw, *Oplegnathuf fasciatus* (Actinopterygii:

322 Perciformes: Oplegnathidae), in Malta (Central Mediterranean) with a discussion on possible

323 modes of entry. *Acta Ichthyol Piscat* 40,101-104 (2010).

324 2.  Mundy, B.C. Checklist of the fishes of the Hawaiian Archipelago. *Bishop Mus Bull Zool* 6,

325 1-704 (2005).

326 3.  An, H.S. & Hong, S.W. Genetic diversity of rock bream *Oplegnathus fasciatus* in Southern

327    Korea. *Genes Genom* 30, 451-459 (2008).

328    4.  Xiao, Y.S. *et al.* Pronounced population genetic differentiation in the rock bream

329        *Oplegnathus fasciatus* inferred from mitochondrial DNA sequences. *Mitochondrial DNA A*

330        27, 2045-2052 (2016).

331    5.  Park, H.S. *et al*. Population Genetic Structure of Rock Bream (*Oplegnathus fasciatus*

332        Temminck & Schlegel, 1884) Revealed by mtDNA COI Sequence in Korea and China.

333        Ocean Sci J 53, 261-274 (2018).

334    6.  Xu, D.D. *et al*. Chromosomal mapping of microsatellite repeats in the rock bream fish

335        *Oplegnathus fasciatus*, with emphasis of their distribution in the neo-Y chromosome. *Mol*

336        *Cytogenet* 6, 12 (2013).

337    7.  Xue, R. *et al.* Karyotype and Ag-Nors In Male And Female Of *Oplegnathus Punctatus*.

338        *Oceanol Limnol Sin* 47, 626-632 (2016).

339    8.  Xiao, Z.Z. Study on population genetics and culture biology of *Oplegnathus fasciatus*.

340        Doctor thesis p 162-176 (2015).

341    9.  Zhang, B.C. *et al*. Rock bream (*Oplegnathus fasciatus*) viperin is a virus-responsive protein

342        that modulates innate immunity and promotes resistance against megalocytivirus infection.

343        *Dev Comp Immunol* 45, 35-42 (2014).

344    10. L, H. *et al*. Characterization of an Iridovirus Detected in Rock Bream (Oplegnathus

345        fasciatus ;Temminck and Schlegel). *Chin J Virol* 27, 158-164 (2011).

346    11. Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting

347        and repeat separation. *Genome Res* 27, 722 (2017).

348    12. Marcais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of

349        occurrences of k-mers. *Bioinformatics* 27, 764–70 (2011).

350    13. Kajitani, R. *et al.* Efficient de novo assembly of highly heterozygous genomes from

351        whole-genome shotgun short reads. *Genome Res* 24, 1384-1395 (2014).

352    14. Sandborn, A.L. *et al*. Chromatin extrusion explains key features of loop and domain

353        formation in wild-type and engineered genomes. *Proc Natl Acad Sci USA* 112, E6456 (2015).

354    15. Pryszcz, L. P. & Gabaldón, T. Redundans: an assembly pipeline for highly heterozygous

355    genomes. *Nucleic Acids Res* 44, e113-e113 (2016).

356  16. Simão, F. A. *et al*. BUSCO: assessing genome assembly and annotation completeness with

357    single-copy orthologs. *Bioinformatics* 31, 3210 (2015).

358  17. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform

359    *Bioinformatics* 25, 1754-1760 (2009).

360  18. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*

361    27, 573–80 (1999).

362  19. Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive elements in

363    genomic sequences. In: Editoral board, Baxevanis Andreas D et al. (eds.), Current Protocols

364    in Bioinformatics, Chapter 4:Unit 4 10 (2009).

365  20. Jurka, J. *et al*. Repbase Update, a database of eukaryotic repetitive elements. Cytogenetic

366    Genome Res 110, 462–67 (2005).

367  21. Stanke, M. *et al*. AUGUSTUS: a web server for gene finding in eukaryotes. *Nucleic Acids*

368    *Res* 32(Web Server issue), W309-12 (2004).

369  22. Cai, Y. *et al*. Computational systems biology methods in molecular biology, chemistry

370    biology, molecular biomedicine, and biopharmacy. *BioMed Res Int* 2014, 746814 (2014).

371  23. Flicek, P. *et al*. Ensembl 2014. *Nucleic Acids Res* 42, D749-D755 (2014).

372  24. Gertz, E. M. *et al*. Composition-based statistics and translated nucleotide searches:

373    Improving the TBLASTN module of BLAST. *MBC Biology* 4, 41 (2006).

374  25. Birney, E. et al. GeneWise and Genomewise. *Genome Res* 14, 988-995 (2004).

375  26. Trapnell, C. *et al*. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* 25,

376    105-1111 (2009).

377  27. Ghosh, S. & Chan, C. K. K. Analysis of RNA-Seq data using TopHat and Cufflinks. *Methods*

378    *Mole Biol* 1374, 339 (2016).

379  28. Campbell, M. S. *et al*. Genome Annotation and Curation Using MAKER and MAKER-P.

380    *Current Protocols in Bioinformatics* 48, 4.11.11 (2014).

381  29. Lobo, I. Basic Local Alignment Search Tool (BLAST). *J Mol Biol* 215, 403-410 (2008).

382  30. Harris, M. A. *et al*. The Gene Ontology (GO) database and informatics resource. *Nucleic*

383    *Acids Res* (2004).

384  31. Ogata, H. *et al*. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Res* 27,

385    29-34 (2000).

386    32. Conesa, A. *et al.* Blast2GO: a universal tool for annotation, visualization and analysis in

387    functional genomics research. *Bioinformatics* 21, 3674 (2005).

388    33. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer RNA

389    genes in genomic sequence. *Nucleic Acids Res* 25, 955-964 (1997).

390    34. Griffithsjones, S. *et al.* Rfam: an RNA family database. *Nucleic Acids Res* 31, 439 (2003).

391    35. Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions reveals

392    folding principles of the human genome. *Science* 326, 289-293 (2009).

393    36. Rao, S.S.P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of

394    chromatin looping. *Cell* 159, 1665-1680 (2014).

395    37. Flot, J.F. *et al.* Contact genomics: scaffolding and phasing (meta) genomes using

396    chromosome be 3D physical signatures. *FEBs Letters* 589, 2966-2974 (2015).

397    38. Burton, J.N. *et al.* Chromosome-scale scaffolding of de novo genome assemblies based on

398    chromatin interactions. *Nat Biotechnol* 31, 1119-1125 (2013).

399    39. Lobo, I. Basic Local Alignment Search Tool (BLAST). J Mol Biol 215, 403-410 (2008).

400    40. Thompson, J. D. *et al.* Multiple sequence alignment using ClustalW and ClustalX. *Curr*

401    *Protoc Bioinformatics* 2.3. 1-2.3. 22 (2002).

402    41. Guindon, S. *et al.* PhyML: Fast and Accurate Phylogeny Reconstruction by Maximum

403    Likelihood. *Infect Genet Evol* 9, 384-385 (2009).

404    42. Yang, Z. & Rannala, B. Bayesian estimation of species divergence times under a molecular

405    clock using multiple fossil calibrations with soft bounds. *Mol Biol Evol* 23, 212-226 (2006).

406    43. Hedges, S. B. *et al.* Tree of life reveals clock-like speciation and diversification. *Mol Biol*

407    *Evol* 32, 835-845 (2015).

408

Table 1 Summary of *Oplegnathus fasciatus* genome assembly and annotation

| Genome assembly | values |
|---|---|
| Contig N50 size (Mb) | 2.1 |
| Contig number | 1 692 |
| Scaffold N50 size (Mb) | 33.5 |
| Scaffold N50 number | 24 |
| Total length (Mb) | 778.7 |
| Genome coverage (X) | 314.6 |
| Contig number ($\geq$ 1 Mb) | 219 |
| Length of contig ($\geq$ 1 Mb) (bp) | 565 184 128 |
| The longest contig (bp) | 8 891 851 |
| The longest scaffold (bp) | 38 619 456 |
| Genome annotation | |
| Protein-coding gene number | 24 003 |
| Mean transcript length (kb) | 16.1 |
| Mean exons per gene | 10.1 |
| Mean exon length (bp) | 217.7 |
| Mean intron length (bp) | 1527.4 |

Table 2 The detailed classification of repeat sequences of *Oplegnathus fasciatus*

| Type | Repbase TEs | | TE proteins | | De novo | | Combined TEs | |
|---|---|---|---|---|---|---|---|---|
| | Length (bp) | % in genome | Length (bp) | % in genome | Length (bp) | % in genome | Length (bp) | % in genome |
| DNA | 39 147 527 | 5.03 | 5 390 266 | 0.69 | 93 089 344 | 11.95 | 124 417 402 | 15.98 |
| LINE | 23 983 322 | 3.08 | 16 460 762 | 2.11 | 57 167 551 | 7.34 | 85 761 250 | 11.01 |
| SINE | 875 585 | 0.11 | 0 | 0.00 | 914 559 | 0.12 | 1 747 250 | 0.22 |
| LTR | 10 163 601 | 1.31 | 5 770 483 | 0.74 | 31 126 639 | 4.00 | 42 465 968 | 5.45 |
| Satellite | 2 028 992 | 0.26 | 0 | 0.00 | 2 613 480 | 0.34 | 4 361 048 | 0.56 |
| Simple_repeat | 1 556 026 | 0.20 | 0 | 0.00 | 5 179 965 | 0.67 | 6 386 303 | 0.82 |
| Other | 6 545 | 0.00 | 0 | 0.00 | 0 | 0.00 | 6 545 | 0.00 |
| Unknown | 331 430 | 0.04 | 0 | 0.00 | 20 636 768 | 2.65 | 20 967 052 | 2.69 |
| Total | 73 544 786 | 9.44 | 27 613 880 | 3.55 | 183 954 095 | 23.62 | 250 611 845 | 32.18 |

Table 3 Hi-C libraries for chromosome-scale assembly of *Oplegnathus fasciatus*

| Chromosome | Number of contigs | Length of contigs | Length of chromosome |
|---|---|---|---|
| Chr1 | 36 | 19 852 463 | 19 869 963 |
| Chr2 | 51 | 34 905 999 | 34 930 999 |
| Chr3 | 43 | 33 654 321 | 33 675 321 |
| Chr4 | 74 | 35 290 762 | 35 327 262 |
| Chr5 | 54 | 38 592 956 | 38 619 456 |
| Chr6 | 72 | 38 156 734 | 38 192 234 |
| Chr7 | 60 | 35 029 969 | 35 059 469 |
| Chr8 | 64 | 37 546 719 | 37 578 219 |
| Chr9 | 45 | 31 457 603 | 31 479 603 |
| Chr10 | 52 | 35 302 682 | 35 328 182 |
| Chr11 | 80 | 31 971 344 | 32 010 844 |
| Chr12 | 46 | 30 287 574 | 30 310 074 |
| Chr13 | 52 | 33 665 353 | 33 690 853 |
| Chr14 | 101 | 31 190 130 | 31 240 130 |

| | | | |
|---|---|---|---|
| Chr15 | 48 | 30 038 946 | 30 062 446 |
| Chr16 | 59 | 28 825 591 | 28 854 591 |
| Chr17 | 33 | 28 220 078 | 28 236 078 |
| Chr18 | 50 | 26 754 155 | 26 778 655 |
| Chr19 | 52 | 34 380 882 | 34 406 382 |
| Chr20 | 52 | 25 675 509 | 25 701 009 |
| Chr21 | 64 | 31 397 692 | 31 429 192 |
| Chr22 | 63 | 30 492 179 | 30 523 179 |
| Chr23 | 70 | 33 514 462 | 33 548 962 |
| Chr24 | 51 | 31 930 140 | 31 956 140 |
| Total | 1 372 | 768 134 243 | 768 808 243 |

Table 4 Genome assembly of *Oplegnathus fasciatus* based on chromosome-length scaffolds

| | Draft scaffolds | Chromosome-length scaffolds based on Hi-C |
|---|---|---|
| Length of genome (bp) | 778 731 089 | 768 808 243 |
| Number of contigs | 1 692 | 1 372 |
| Contigs N50 (bp) | 2 149 025 | 2 130 780 |
| Number of scaffold | / | 24 |
| Scaffold N50 (bp) | / | 33 548 962 |
| Number of contigs (≥ 100 kb) | 693 | 708 |
| Total length of contigs (≥ 100 kb) | 735 235 962 | 732 827 446 |
| Mapping rate of contigs (≥ 100 kb) (%) | / | 99.67 |

# Figure Legends



Fig. 1 A representative individual of *O. fasciatus*

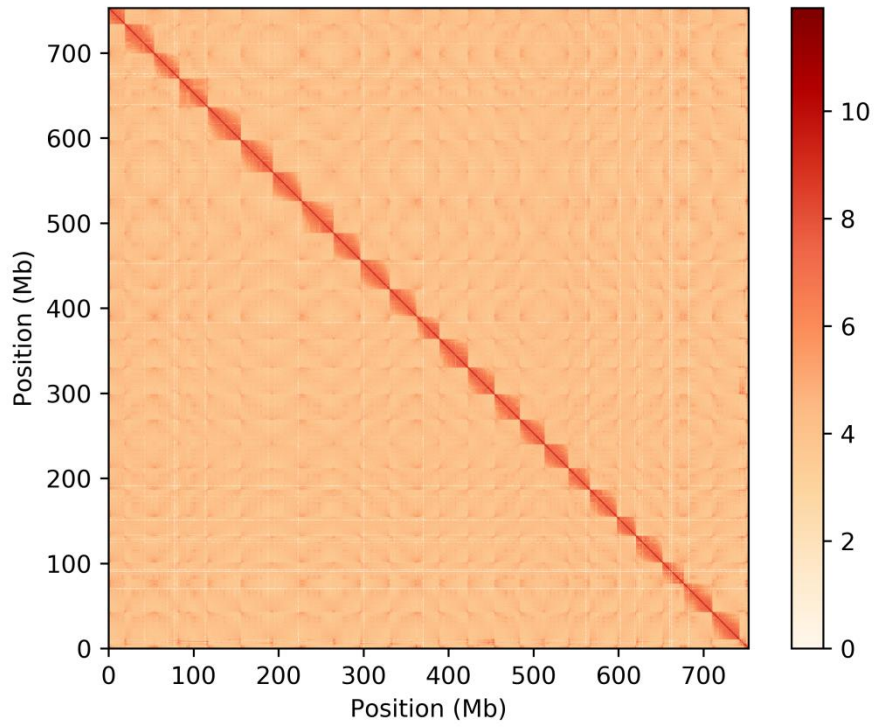Fig. 2 *k*-mer distribution of the *O. fasciatus* genome

Fig. 3 Hi-C interaction heatmap for *O. fasciatus* reference genome, showing
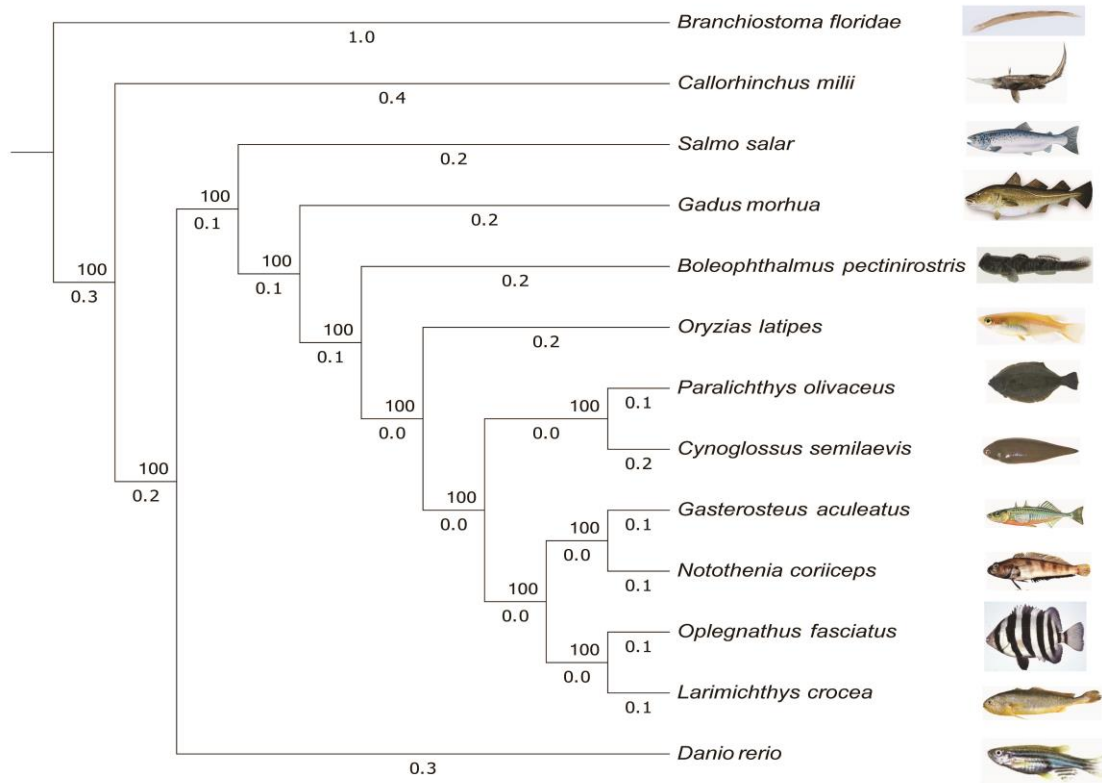
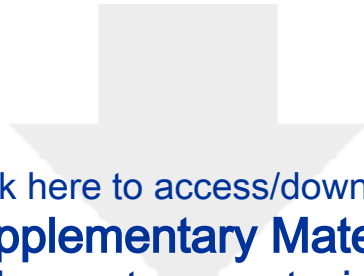interactions between the 24 chromosomes

Fig. 4 The phylogenetic relationships of *O. fasciatus* with other fishes

Click here to access/download
**Supplementary Material**
4-supplementary materials.docx