

1 Chromosome-level genome assembly of a butterflyfish, *Chelmon rostratus*

2 Xiaoyun Huang^{1,2,3,*}, Yue Song^{1,2,3,*}, Suyu Zhang^{1,2,3,*}, Yunga A^{1,2,3,*}, Mengqi
3 Zhang^{1,2,3}, Yue Chang^{1,2,3}, He Zhang^{1,2,3}, Chang Li^{1,2,3}, Yong Zhao^{1,2,3}, Meiru Liu^{1,2,3},
4 Inge Seim^{4,5}, Guangyi Fan^{1,2,3}, Xin Liu^{1,2,3,6,#}, Shanshan Liu^{1,2,3,#}

5

⁶ ¹ BGI-Qingdao, BGI-Shenzhen, Qingdao, 266555, China;

⁷ ² BGI-Shenzhen, Beishan Industrial Zone, Yantian District, Shenzhen 518083, China;

⁸ ³ China National Gene Bank, BGI-Shenzhen, Jinsha Road, Shenzhen 518120, China;

⁴ Integrative Biology Laboratory, College of Life Sciences, Nanjing Normal University, Nanjing 210023, China;

⁵ Comparative and Endocrine Biology Laboratory, Translational Research Institute-Institute of Health and Biomedical Innovation, School of Biomedical Sciences, Queensland University of Technology, Brisbane 4102, Queensland, Australia;

⁶ James D. Watson Institute of Genome Science, 310008 Hangzhou, China.

15

13

16

¹⁸ * These authors contributed equally to this work

[#] These authors jointly supervised this work.

Corresponding to liushanshan@genomics.cn (S.L.) and liuxin@genomics.cn (X.L.)

21

22

23 Abstract

24 *Chelmon rostratus* (Teleostei, Perciformes, Chaetodontidae) is a copperband
25 butterflyfish. As an ornamental fish, the genome information for this species might
26 help understanding the genome evolution of Chaetodontidae and adaptation/evolution
27 of coral reef fish

28 In this study, using the stLFR co-Barcode reads data, we assembled a genome of
29 638.70 Mb in size with contig and scaffold N50 sizes of 294.41 kb and 2.61 Mb,
30 respectively. 94.40% of scaffold sequences were assigned to 24 chromosomes using
31 Hi-C data and BUSCO analysis showed that 97.3% (2,579) of core genes were found
32 in our assembly. Up to 21.47 % of the genome was found to be repetitive sequences

1 and 21,375 protein-coding genes were annotated. Among these annotated
2 protein-coding genes, 20,163 (94.33%) proteins were assigned with possible
3 functions.

4 As the first genome for Chaetodontidae family, the information of these data helpfully
5 to improve the essential to the further understanding and exploration of marine
6 ecological environment symbiosis with coral and the genomic innovations and
7 molecular mechanisms contributing to its unique morphology and physiological
8 features.

9 **Keywords:** Chaetodontidae; *Chelmon rostratus*; stLFR; Genome assembly; Gene
10 annotation; Phylogenetic tree

11 **Background & Summary**

12 *Chelmon rostratus* (Teleostei, Perciformes, Chaetodontidae) lives in the western
13 pacific, from southern Japan and Taiwan throughout the Coral Triangle to the
14 Solomon Islands and the northern coast of Australia¹. In the natural environment, it
15 lives in depths from 1- 25 meters underwater, inhabiting coastal and inner reefs and
16 often in turbid water². It usually has an appealing appearance of four yellowish or
17 yellowish-orange vertical bands with black border on the silvery white body, and a
18 false eyespot. Thus, it is one of the protagonists in tropical marine aquarium fish³.

19 Despite as a common coral reef fish, there were limited previous researches in this
20 species. Previous researches were limited to its behavior features such as the prey
21 capture kinematics⁴. More recently, its complete mitochondrial genome was reported
22⁵, revealing its evolutionary position. However, there was no particular study
23 providing the whole genomic information.

24 Here, we reported the first whole genome sequencing data for this species, resulting in
25 a chromosome level genome assembly and annotation, followed by the evolutionary
26 analysis. The data and analysis provided here, can benefit future basic studies and
27 conservation efforts of this species.

28 **Methods**

1 **Sample collection, library construction and sequencing**

2 We captured an adult *Chelmon rostratus* (Figure 1) in the sea area of Qingdao,
3 China, and muscle tissue was stored in liquid nitrogen. Then, high quality DNA was
4 extracted using a modified DNA extraction for vertebrate tissues protocol ⁶ from the
5 muscle tissue. The exacted DNA was fragmented and MGIEasy stLFR Library
6 preparation kit (PN : 1000005622) was used to construct single tube Long Fragment
7 Read (stLFR) library, following the stLFR protocol ⁷. For Hi-C library sequencing,
8 about 1g living muscle tissue was used to DNA extraction and library contraction,
9 according to Wang's method ⁷. Sequencing was conducted on a BGISEQ-500
10 sequencer, generating 286.82Gb raw data (including 134.18Gb raw Hi-C data)
11 (Supplementary Table 1). Data filtering was then carried out using SOAPnuke
12 software (version 1.5) ⁸ with the default parameters, thus low-quality reads (more
13 than 40% bases with quality score lower than 8), PCR duplications, adaptors and
14 reads with high proportion (higher than 10%) of ambiguous bases (Ns) were filtered.
15 After data filtering, 192.71 Gb 'clean data' (including 133.17Gb 'clean Hi-C data')
16 was obtained for the further assembly (Supplementary Table 1).

17 **Genome features revealed by *k*-mer analysis**

18 In order to understand better about its genome features, we applied *k*-mer frequency
19 distribution analysis to estimate its genome size and genome complexity
20 (heterozygosity or/and repetition) ⁹. We randomly selected 59.54 Gb (~90X) clean
21 reads and carried out a 17-mer analysis using KMERFREQ_AR (version 2.0.4) ¹⁰.
22 The genome size was estimated to be ~711.39 Mb (Supplementary Table 2).
23 Observing the distribution of 17-mer frequency (Supplementary Figure 1), we
24 anticipated the genome to be a diploid species with slight heterozygosity proportion.
25 In order to get the percentage of heterozygosity, GCE software ⁹ was used to
26 calculation with 59.54 Gb clean reads and resulted that the percentage of
27 heterozygosity was 0.72%.

1 **Genome assembly and annotation**

2 We assembled the genome using Supernova (version 2.1)¹¹ with default parameters
3 for 55 Gb clean stLFR data. After that, we used Gapcloser software¹⁰ to fill the gaps
4 within the assembly with default parameters. Finally, we obtained an assembly of
5 638.70 Mb in size containing 5,490 scaffolds. The N50 values of contigs and scaffolds
6 were 294.41 Kb and 2.61 Mb (Table 1), respectively, revealing a good contiguity of
7 the genome assembly. The longest scaffold and contig were 9.02 Mb and 1.96 Mb,
8 respectively. The assembled length was account for ~90 % of estimated genome size.
9 To generate a chromosomal-level genome assembly, 133.17 Gb high-quality Hi-C
10 data was used to further assembly. We first used HiC-Pro software
11 (version2.8.0_devel) with default parameters to get ~26 Gb valid sequencing data,
12 accounting for 19.31% of total Hi-C clean reads. Then, Juicer (version 1.5, an
13 opensource tool for analyzing Hi-C datasets)¹², and the 3D de novo assembly pipeline
14¹³ were used to connect the scaffolds to chromosomes with length of 603Mb (account
15 for 94.40% of total genome) and scaffold sequences were assigned to 24
16 chromosomes, with the length from 14.13 Mb to 33.19 Mb (Table 2, Figure 2A and
17 Figure 3).
18 We randomly selected 4 Gb clean data to assemble mitochondrial genome sequence
19 using MitoZ software¹⁴, resulting 16.52 kb of assembly with a cyclic structure
20 (Figure 2B).

21 **Repetitive sequence and gene annotation**

22 We annotated the two major types of repetitive sequences (tandem repeats, TRFs and
23 transposable elements, TEs) in the assembled genome. According to the method of
24 previously research¹⁵⁻¹⁷, TRFs were identified using Tandem Repeats Finder (version
25 4.04)¹⁸. Transposable elements (TEs) were identified by a combination of
26 homology-based and *de novo* approaches. Briefly, for homology-based annotation,
27 known repeats in the database (Repbase16.02)¹⁹ were aligned against the genome
28 assembly using RepeatMasker and RepeatProteinMask (version 3.2.9)²⁰ at both the

1 DNA and protein levels. For *de novo* annotation, RepeatModeler (version 1.1.0.4)²¹
2 was employed to build a *de novo* non-redundant repeat library and then this repeat
3 library was searched against the genome using RepeatMasker²⁰. In this way, up to
4 21.47 % of the assembled sequences were found to be repeat sequences (Figure 3,
5 Supplementary Table 3).

6 Protein-coding gene were then predicted by a combination of two ways: (1) the *ab*
7 *initio* gene prediction and (2) the homology-based annotation²²⁻²⁴. For *ab initio* gene
8 prediction approaches, Augustus²⁵ and GlimmerHMM²⁶ were used with *Danio rerio*
9 as the species of HMM model to predict gene models; For homology-based
10 annotation, four homolog species including *Pundamilia nyererei*, *Maylandia zebra*,
11 *Astatotilapia calliptera* and *Perca flavescens* were aligned against the genome
12 assembly using BLAT software (version 0.36)²⁷ and GeneWise software (version
13 2.4.1)²⁸. 21,375 protein-coding genes were obtained by combining the different
14 evidences using Glean software (version 1.0)²⁹. In the final gene models, the average
15 length was 16,183.81 bp, with an average of 10 exons. The average length of coding
16 sequences, exons and introns were 1,789.82bp, 179.10 bp and 1599.48 bp,
17 respectively, similar to that of the other released fish genomes, such as *Astatotilapia*
18 *calliptera*, *Maylandia zebra*, *Perca flavescens* and *Pundamilia nyererei*³⁰⁻³²
19 (Supplementary Table 4, Supplementary Figure 2). Gene annotation of mitochondria
20 was performed using MitoZ software¹⁴, and 13 protein-coding genes as well as 22
21 tRNA genes were annotated.

22 Functions of the annotated protein-coding genes were inferred by searching homologs
23 in the databases (KEGG, COG, NR, Swissprot and Interpro)³³⁻³⁷. In this way, 18,005
24 (84.23%), 7,343 (34.35%), 20,141 (94.23%), 19,114 (89.42%) and 19,313 (90.35%)
25 of protein-coding genes had their homologous alignment in the above databases,
26 respectively. The remaining 1,212 (5.67%) protein-coding genes with unknown
27 function might be the specific feature of the *Chelmon rostratus* genome
28 (Supplementary Table 5).

29 **Gene family and phylogenetic analysis**

1 To identify and analyze the gene families, we selected other eleven species with
2 whole genome sequences available (*Lepisosteus oculatus*, *Hippocampus comes*,
3 *Larimichthys crocea*, *Gasterosteus aculeatus*, *Takifugu rubripes*, *Oreochromis*
4 *niloticus*, *Astyanax mexicanus*, *Danio rerio*, *Cynoglossus semilaevis*, *Oryzias latipes*
5 and *Homo sapiens* as outgroup)^{24,38-47}. The protein coding genes of the total twelve
6 species were clustered into 18,502 gene families using TreeFam⁴⁸⁻⁵⁰. Among these
7 gene families, 2,301 were single-copy gene families (one copy in each of these
8 species) (Figure 4B). The 21,375 protein-coding genes of *Chelmon rostratus* were
9 classified into 13,797 gene families, given an average of 1.55 genes per gene family.
10 Compared to the other species, it was similar to the gene family numbers of
11 *Lepisosteus oculatus* (13,967) and *Gasterosteus aculeatus* (13,331), but was quite
12 different from those of *Larimichthys crocea* (14,724), *Homo sapiens* (14,578) and
13 *Takifugu rubripes* (12,554) (Supplementary Table 6). Among the clustered gene
14 families from *Chelmon rostratus*, 8,190 gene families were common to at least one of
15 the other species, the remaining 52 gene families were unique. Between the four
16 species (*Chelmon rostratus*, *Danio rerio*, *Takifugu rubripes* and *Larimichthys crocea*),
17 the number of common shared and unique gene was shown in Figure 4C. To
18 understand the function of these gene families, we further performed GO enrichment
19 with these gene families from *Chelmon rostratus*, compared with the other 11 species.
20 The result reflected that unique gene families from *Chelmon rostratus* were enriched
21 in muscle contraction functions (Supplementary Table 7). Phylogenetic analysis using
22 the concatenated sequence alignment of the 2,301 single-copy genes shared by the
23 twelve species was performed. The PhyML software (version 3.0)⁵¹, based on the
24 method of maximum likelihood, was used to construct the phylogenetic tree. The split
25 time between *Chelmon rostratus* and *Larimichthys crocea* was estimated to be ~92
26 million years ago (Figure 4A). Based on the similarity of the protein sequences, 483
27 syntenic blocks were identified by using McScanX software (version 0.8)⁵²
28 (Supplementary Table 8). The time of the duplication and divergence event in these
29 species was calculated based on the distribution of synonymous mutation rate for the

1 gene pairs in the paralogous syntenic blocks, indicating that whole-genome
2 duplication (WGD) event was not detected in *Chelmon rostratus* genome (Figure 3).
3 The expansion and contraction of the gene family analysis may reveal the
4 evolutionary dynamics of gene families thus provide the clues for understanding the
5 diversity of different species. It is often inferred from the number of genes in the gene
6 family and the phylogenetic tree. In our study, we used the CAFÉ (version 2.1)⁵³
7 software to analyze the expansion and contraction of clustered gene families (Figure
8 4D). As a result, a total of 18,498 gene families from the most recent common
9 ancestor (MRCA) have been identified. Compared to the recent common ancestor
10 between *Chelmon rostratus* and *Larimichthys crocea*, 793 gene families were
11 expanded and the majority of the expanded gene families were found to be involved
12 in synapse organization. (Supplementary Table 9). On the other hand, there was 2,962
13 gene family contracted involved in immune system process (Supplementary Table
14 10).

15 **Data Records**

16 Raw reads from BGISEQ-500 sequencing are deposited in the CNGB Nucleotide
17 Sequence Archive (CNSA) with accession number CNP0000597
18 (<https://db.cngb.org/cnsa>). Data Citation 1: CNGB Nucleotide Sequence
19 Archive CNP0000597).

20 **Technical Validation**

21 To evaluate the genome assembly, we aligned sequencing data which we filtered
22 previously using SOAPaligner (version 2.2)¹⁰ and found that 90.76% could be
23 mapped back to the assembled genome. we also calculate its GC depth to rule out
24 possible biases during sequencing or possible contaminations. We identified the
25 average GC contents of this genome to be ~42.52% and we found a continuous GC
26 depth distribution (Supplementary Figure 3), indicating no obvious assembly errors
27 resulted from GC content or contamination. The genome completeness assessment
28 was estimated with Benchmarking Universal Single-copy Orthologs (BUSCO,

1 version 3.0.1) ⁵⁴. BUSCO analysis showed that 97.3% (2,518) of core genes were
2 found in our assembly with 2,491 (96.3%) were single copy gene and 27 (1.0%) were
3 duplicated (Supplementary Table 11), indicating a good coverage of the genome.
4 To validate the quality of predicted gene sets, we also assessed the completeness using
5 BUSCO (version 3) with the fish core gene database (actinopterygii_odb9) ⁵⁵. We
6 found that about 90.20% of core gene were annotated in our gene set with 4,000 were
7 single copy gene and 132 were duplicated (Supplementary Table 12).

8 **Acknowledgements**

9 This work is supported by the special funding of “Blue granary” scientific and
10 technological innovation of China (2018YFD0900301-05). The work also received
11 the technical support from China National Gene Bank.

12 **Author contributions**

13 Y.S., G.F., S.L. and X.L. conceived the work. M.Z. collected sample and X.H.
14 sequenced the libraries. Y.S., X.H., S.Z., Y.A. collected the public data and performed
15 the analyses. Y.C., M.L., C.L. and Y.Z. helped in the analysis. Y.S., X.H., G.F., I.S.,
16 S.L. and X.L. wrote and revised the manuscript.

17 **Competing interests**

18 The authors declare no competing interests.

19 **References**

- 20 1 Kuitert, R. H. & Tonozuka, T. *Pictorial guide to Indonesian reef fishes*. (Zoonetics, 2001).
- 21 2 Lieske, E. & Myers, R. *Collins pocket guide: coral reef fishes*. Collins, London (2001).
- 22 3 EA, L. & CK, T. Fine structure of the gastric epithelium of the coral fish, Chelmon rostratus
23 Cuvier. *Okajimas folia anatomica Japonica* **51**, 285-309 (1975).
- 24 4 Ferry-Graham, L. A., Wainwright, P. C. & Bellwood, D. R. Prey capture in long-jawed
25 butterflyfishes (Chaetodontidae): the functional basis of novel feeding habits. *Journal of
26 Experimental Marine Biology and Ecology* **256**, 167-184 (2001).
- 27 5 Wang, L.-J., You, F. & Wu, Z.-H. Complete mitochondrial genome of copperband butterflyfish
28 Chelmon rostratus (Teleostei, Perciformes, Chaetodontidae). *Mitochondrial DNA Part A* **27**,
29 2141-2142 (2016).
- 30 6 Panova, M. *et al.* DNA extraction protocols for whole-genome sequencing in marine

1 organisms in *Marine Genomics*. 13-44 (Springer, 2016).

2 7 Wang, O. *et al.* Efficient and unique cobarcoding of second-generation sequencing reads from
3 long DNA molecules enabling cost-effective and accurate sequencing, haplotyping, and de
4 novo assembly. *Genome research* **29**, 798-808 (2019).

5 8 Chen, Y. *et al.* SOAPnuke: a MapReduce acceleration-supported software for integrated
6 quality control and preprocessing of high-throughput sequencing data. *GigaScience* **7**, gix120
7 (2017).

8 9 Liu, B. *et al.* Estimation of genomic characteristics by analyzing k-mer frequency in de novo
9 genome projects. *arXiv preprint arXiv:1308.2012* (2013).

10 10 Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read de novo
11 assembler. *GigaScience* **1**, 18 (2012).

12 11 Weisenfeld, N. I., Kumar, V., Shah, P., Church, D. M. & Jaffe, D. B. Direct determination of
13 diploid genome sequences. *Genome research* **27**, 757-767 (2017).

14 12 Durand, N. C. *et al.* Juicer provides a one-click system for analyzing loop-resolution Hi-C
15 experiments. *Cell systems* **3**, 95-98 (2016).

16 13 Bonev, B. *et al.* Multiscale 3D genome rewiring during mouse neural development. *Cell* **171**,
17 557-572. e524 (2017).

18 14 Meng, G., Li, Y., Yang, C. & Liu, S. MitoZ: a toolkit for animal mitochondrial genome
19 assembly, annotation and visualization. *Nucleic acids research* **47**, e63-e63 (2019).

20 15 Li, C. *et al.* Draft genome of the Peruvian scallop *Argopecten purpuratus*. *GigaScience* **7**,
21 giy031 (2018).

22 16 Liu, H.-P. *et al.* The sequence and de novo assembly of *Oxygymnocypris stewartii* genome.
23 *Scientific data* **6**, 190009 (2019).

24 17 Song, L. *et al.* Draft genome of the Chinese mitten crab, *Eriocheir sinensis*. *GigaScience* **5**, 5
25 (2016).

26 18 Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids
27 research* **27**, 573-580 (1999).

28 19 Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in
29 eukaryotic genomes. *Mobile DNA* **6**, 11 (2015).

30 20 Smit, A., Hubley, R. & Green, P. RepeatMasker, version 4.0.9. Available from
31 <http://www.repeatmasker.org/> (2015).

32 21 Smit, A. & Hubley, R. RepeatModeler Open-1.0. Available from <http://www.repeatmasker.org>
33 (2008).

34 22 Shao, C. *et al.* Chromosome-level genome assembly of the spotted sea bass, *Lateolabrax
35 maculatus*. *GigaScience* **7**, doi:10.1093/gigascience/giy114 (2018).

36 23 Valenzano, D. R. *et al.* The African turquoise killifish genome provides insights into evolution
37 and genetic architecture of lifespan. *Cell* **163**, 1539-1554 (2015).

38 24 Ao, J. *et al.* Genome sequencing of the perciform fish *Larimichthys crocea* provides insights
39 into molecular and genetic mechanisms of stress adaptation. *PLoS genetics* **11**, e1005118
40 (2015).

41 25 Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic acids
42 research* **34**, W435-W439 (2006).

43 26 Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM: two open source
44 ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878-2879 (2004).

1 27 Kent, W. J. BLAT—the BLAST-like alignment tool. *Genome research* **12**, 656-664 (2002).

2 28 Birney, E., Clamp, M. & Durbin, R. GeneWise and genomewise. *Genome research* **14**, 988-995 (2004).

3 29 Elsik, C. G. *et al.* Creating a honey bee consensus gene set. *Genome biology* **8**, R13 (2007).

4 30 Brawand, D. *et al.* The genomic substrate for adaptive radiation in African cichlid fish. *Nature* **513**, 375 (2014).

5 31 Ozerov, M. Y. *et al.* Highly Continuous Genome Assembly of Eurasian Perch (*Perca fluviatilis*) Using Linked-Read Sequencing. *G3: Genes, Genomes, Genetics* **8**, 3737-3743 (2018).

6 32 Malinsky, M. *et al.* Whole-genome sequences of Malawi cichlids reveal multiple radiations interconnected by gene flow. *Nature ecology & evolution* **2**, 1940 (2018).

7 33 Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids research* **28**, 27-30 (2000).

8 34 Tatusov, R. L. *et al.* The COG database: an updated version includes eukaryotes. *BMC bioinformatics* **4**, 41 (2003).

9 35 Kanz, C. *et al.* The EMBL nucleotide sequence database. *Nucleic acids research* **33**, D29-D33 (2005).

10 36 Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic acids research* **28**, 45-48 (2000).

11 37 Mitchell, A. L. *et al.* InterPro in 2019: improving coverage, classification and access to protein sequence annotations. *Nucleic Acids Research* (2018).

12 38 Braasch, I. *et al.* The spotted gar genome illuminates vertebrate evolution and facilitates human-teleost comparisons. *Nature Genetics* **48**, 427-437 (2016).

13 39 Conte, M. A., Gammerdinger, W. J., Bartie, K. L., Penman, D. J. & Kocher, T. D. A high quality assembly of the Nile Tilapia (*Oreochromis niloticus*) genome reveals the structure of two sex determination regions. *Bmc Genomics* **18**, 341 (2017).

14 40 Hinaux, H. *et al.* De novo sequencing of *Astyanax mexicanus* surface fish and Pachón cavefish transcriptomes reveals enrichment of mutations in cavefish putative eye genes. *PLoS one* **8**, e53553 (2013).

15 41 Howe, K. *et al.* The zebrafish reference genome sequence and its relationship to the human genome. *Nature* **496**, 498 (2013).

16 42 Chen, S. *et al.* Whole-genome sequence of a flatfish provides insights into ZW sex chromosome evolution and adaptation to a benthic lifestyle. *Nature genetics* **46**, 253 (2014).

17 43 Kasahara, M. *et al.* Kasahara, M. *et al.* The medaka draft genome and insights into vertebrate genome evolution. *Nature* **447**, 714-719. *Nature* **447**, 714-719 (2007).

18 44 Wong, K. H. Y., Michal, L.-S. & Pui-Yan, K. De novo human genome assemblies reveal spectrum of alternative haplotypes in diverse populations. *Nature Communications* **9**, 3040- (2018).

19 45 Lin, Q. *et al.* The seahorse genome and the evolution of its specialized morphology. *Nature* **540**, 395-399 (2016).

20 46 Jones, F. C. *et al.* The genomic basis of adaptive evolution in threespine sticklebacks. *Nature* **484**, 55-61 (2012).

21 47 Aparicio & S. Whole-Genome Shotgun Assembly and Analysis of the Genome of Fugu rubripes. *Science* **297**, 1301-1310 (2002).

22 48 Ruan, J. *et al.* TreeFam: 2008 update. *Nucleic acids research* **36**, D735-D740 (2007).

1 49 Li, H. *et al.* TreeFam: a curated database of phylogenetic trees of animal gene families.
2 *Nucleic acids research* **34**, D572-D580 (2006).

3 50 Schreiber, F., Patricio, M., Muffato, M., Pignatelli, M. & Bateman, A. TreeFam v9: a new
4 website, more species and orthology-on-the-fly. *Nucleic acids research* **42**, D922-D925
5 (2013).

6 51 Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood phylogenies:
7 assessing the performance of PhyML 3.0. *Systematic biology* **59**, 307-321 (2010)

8 52 Wang, Y. *et al.* MCScanX: a toolkit for detection and evolutionary analysis of gene synteny
9 and collinearity. *Nucleic acids research* **40**, e49-e49 (2012).

10 53 De Bie, T., Cristianini, N., Demuth, J. P. & Hahn, M. W. CAFE: a computational tool for the
11 study of gene family evolution. *Bioinformatics* **22**, 1269-1271 (2006).

12 54 Waterhouse, R. M. *et al.* BUSCO applications from quality assessments to gene prediction and
13 phylogenomics. *Molecular Biology & Evolution* **35** (2017).

14 55 Yang, X. *et al.* Chromosome-level genome assembly of *Triphlophysa tibetana*, a fish adapted
15 to the harsh high-altitude environment of the Tibetan Plateau. *Molecular ecology resources*
16 (2019).

17 Data Citations

18 1. Ao, J. *et al.* Genome sequencing of the perciform fish *Larimichthys crocea* provides insights into
19 molecular and genetic mechanisms of stress adaptation. *PLoS genetics* **11**, e1005118 (2015).

20 2. Brawand, D. *et al.* The genomic substrate for adaptive radiation in African cichlid fish. *Nature* **513**,
21 375 (2014).

22 3. Ozerov, M. Y. *et al.* Highly Continuous Genome Assembly of Eurasian Perch (*Perca fluviatilis*)
23 Using Linked-Read Sequencing. *G3: Genes, Genomes, Genetics* **8**, 3737-3743 (2018).

24 4. Malinsky, M. *et al.* Whole-genome sequences of Malawi cichlids reveal multiple radiations
25 interconnected by gene flow. *Nature ecology & evolution* **2**, 1940 (2018).

26 5. Braasch, I. *et al.* The spotted gar genome illuminates vertebrate evolution and facilitates
27 human-teleost comparisons. *Nature Genetics* **48**, 427-437 (2016).

28 6. Conte, M. A., Gammerdinger, W. J., Bartie, K. L., Penman, D. J. & Kocher, T. D. A high quality
29 assembly of the Nile Tilapia (*Oreochromis niloticus*) genome reveals the structure of two sex
30 determination regions. *Bmc Genomics* **18**, 341 (2017).

31 7. Hinaux, H. *et al.* De novo sequencing of *Astyanax mexicanus* surface fish and Pachón cavefish
32 transcriptomes reveals enrichment of mutations in cavefish putative eye genes. *PloS one* **8**, e53553
33 (2013).

34 8. Howe, K. *et al.* The zebrafish reference genome sequence and its relationship to the human
35 genome. *Nature* **496**, 498 (2013).

36 9. Chen, S. *et al.* Whole-genome sequence of a flatfish provides insights into ZW sex chromosome
37 evolution and adaptation to a benthic lifestyle. *Nature genetics* **46**, 253 (2014).

38 10. Kasahara, M. *et al.* Kasahara, M. *et al.* The medaka draft genome and insights into vertebrate
39 genome evolution. *Nature* **447**, 714-719. *Nature* **447**, 714-719 (2007).

40 11. Wong, K. H. Y., Michal, L.-S. & Pui-Yan, K. De novo human genome assemblies reveal spectrum
41 of alternative haplotypes in diverse populations. *Nature Communications* **9**, 3040- (2018).

42 12. Lin, Q. *et al.* The seahorse genome and the evolution of its specialized morphology. *Nature* **540**,
43 395-399 (2016).

- 1 13. Jones, F. C. *et al.* The genomic basis of adaptive evolution in threespine sticklebacks. *Nature* **484**,
- 2 55-61 (2012).
- 3 14. Aparicio & S. Whole-Genome Shotgun Assembly and Analysis of the Genome of Fugu rubripes.
- 4 *Science* **297**, 1301-1310 (2002).

5

6

1 **Tables**

2

3 **Table 1 Summary of the scaffold-level assembly.** The N50 values of contigs and
4 scaffolds were 294.41 Kb and 2.61 Mb, respectively, indicating good contiguity of the
5 scaffold-level assembly for further Hi-C data assembly.

Species name	<i>Chelmon rostratus</i>
Estimated genome size (bp)	711,393,276
Assembly size (bp)	638,700,992
Scaffold N50 (bp)	2,610,810
Longest scaffold (bp)	9,024,273
Contig N50 (bp)	294,414
Longest contig (bp)	1,961,159
Number of genes	21,375
GC content (%)	42.52
BUSCO (%)	98.40
Mapping ratio (%)	90.76
Coverage (%)	99.28

6

7

1 **Table 2 Summary of the chromosome-level assembly.** About ~26 Gb valid Hi-C
2 sequencing data was used to connect scaffold-level assembly into chromosomes with
3 length of 603Mb (account for 94.40% of scaffold-level assembly size) and
4 constructed 24 chromosomes, with lengths ranging from 14.13 (chr24) Mb to 33.19
5 Mb (chr1).

Chromosome	Length (bp)	GC (%)
chr1	33,188,845	41.94
chr2	29,501,685	42.26
chr3	29,388,801	42.61
chr4	29,119,845	42.33
chr5	29,062,128	42.38
chr6	28,079,113	42.55
chr7	28,043,518	42.42
chr8	26,564,978	42.56
chr9	26,447,630	42.51
chr10	26,405,724	42.19
chr11	26,164,899	42.37
chr12	25,929,763	42.31
chr13	25,766,544	42.17
chr14	25,529,937	42.41
chr15	24,867,447	42.81
chr16	24,414,281	43.00
chr17	24,306,480	42.65
chr18	23,537,532	42.75
chr19	23,051,343	42.21
chr20	22,222,359	42.56
chr21	21,649,575	42.91
chr22	20,076,678	43.05
chr23	15,560,595	43.28
chr24	14,125,005	43.64
Average	25,125,196	42.58

6

7

1 **Figure Legends**

2

3 **Figure 1 Photograph of *Chelmon rostratus*.**

4 **Figure 2 Heat map of interactive intensity between chromosome sequences (A)**
5 **and Physical map of mitochondrial assembly (B).**

6 **Figure 3 Distribution of basic genomic elements of *Chelmon rostratus* genome.**

7 (A) Chromosome karyotype. Different colored represented different chromosome we
8 assembled. (B) Gene density. The histogram indicated number of genes per 1 Mb
9 ranges from a minimum of 0.17 to a maximum of 1, illustrated by blue bar. (C)
10 Repeat sequence density. The histogram indicated average DNA TE ratio per 1 Mb
11 ranges from 0.16 to 1, illustrated by orange bar. Synteny blocks of each chromosome
12 was illustrated by color lines, indicating that whole-genome duplication (WGD) event
13 was not detected in *Chelmon rostratus* genome. Circos (Krzywinski et al. 2009)
14 (<http://circos.ca>) was used for constructing this diagram.

15 **Figure 4 Comparative analysis of the *Chelmon rostratus* genome.** (A)
16 Phylogenetic analysis among *Lepisosteus oculatus*, *Hippocampus comes*,
17 *Larimichthys crocea*, *Gasterosteus aculeatus*, *Takifugu rubripes*, *Oreochromis*
18 *niloticus*, *Astyanax mexicanus*, *Danio rerio*, *Cynoglossus semilaevis*; *Oryzias latipes*
19 and *Homo sapiens* as outgroup, by using the single-copy gene families. The species
20 differentiation time between *Chelmon rostratus* and *Larimichthys crocea* was ~92
21 million years ago. (B) The protein coding genes of the total twelve species were
22 clustered into 18,502 gene families. Among these gene families, 2,301 were
23 single-copy gene families (one copy in each of these species). (C) Venn diagram
24 showing overlaps of gene families between *Chelmon rostratus*, *Danio rerio*, *Takifugu*
25 *rubripes* and *Larimichthys crocea*. A total of 322 gene families were unique to
26 *Chelmon rostratus* and 10,711 were commonly shared by the other species genome.
27 (D) Compared to the recent common ancestor between *Chelmon rostratus* and
28 *Larimichthys crocea*, 793 gene families were expanded and 2,962 gene family
29 contracted in *Chelmon rostratus* genome.

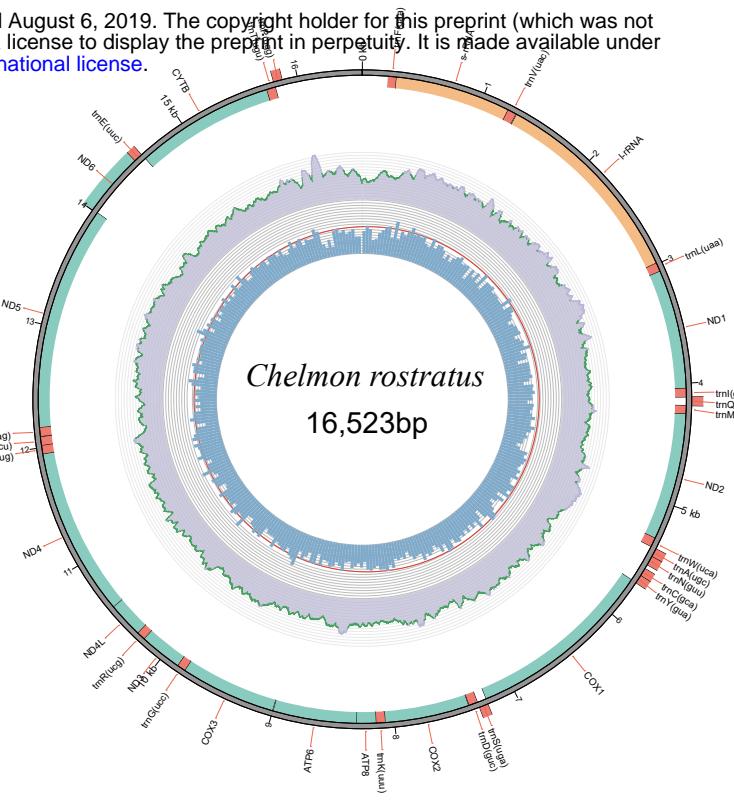
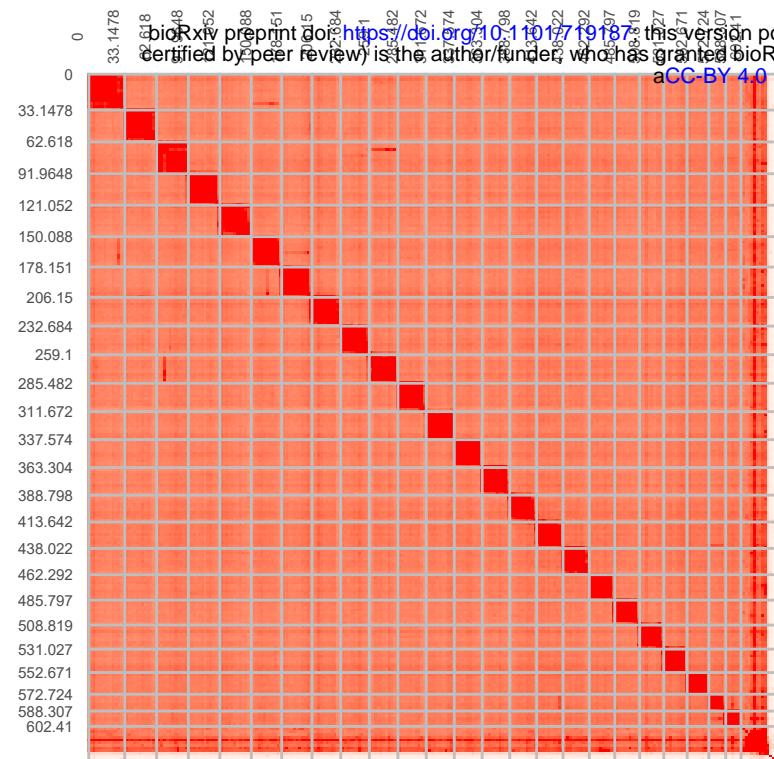


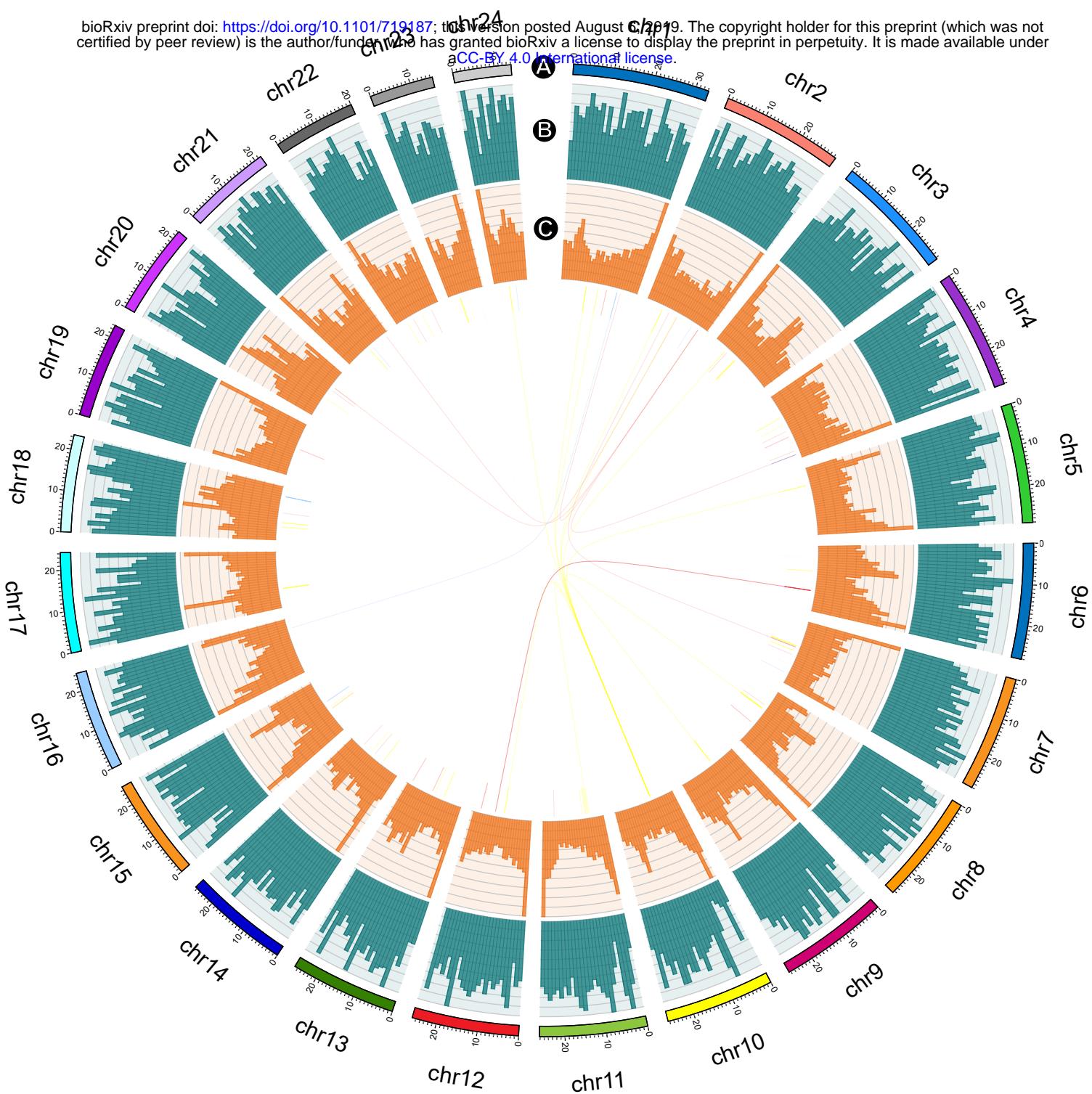
20180609 XM 006

A

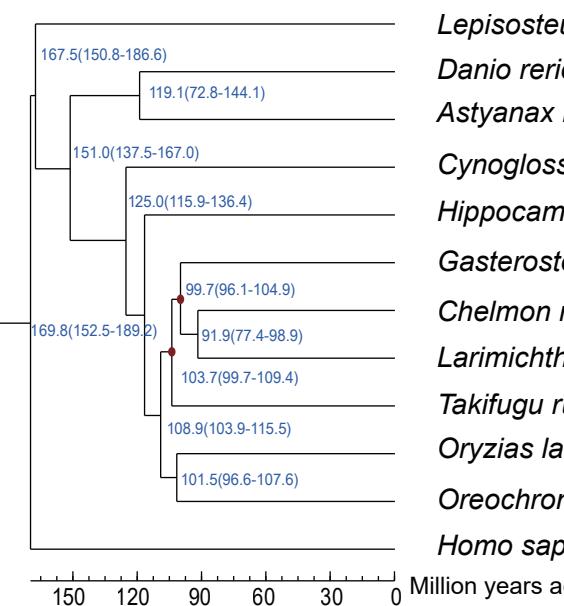
B

this version posted August 6, 2019. The copyright holder for this preprint (which was not granted a license to display the preprint in perpetuity) is made available under aCC-BY 4.0 International license.

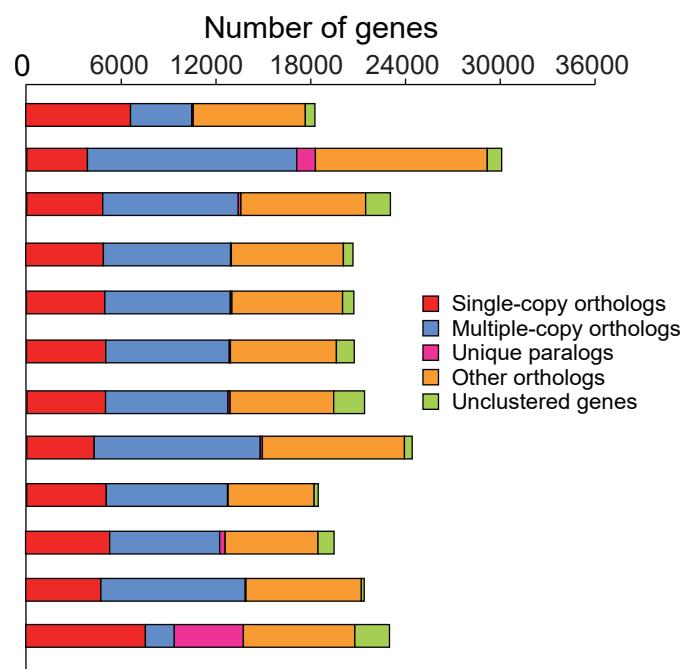




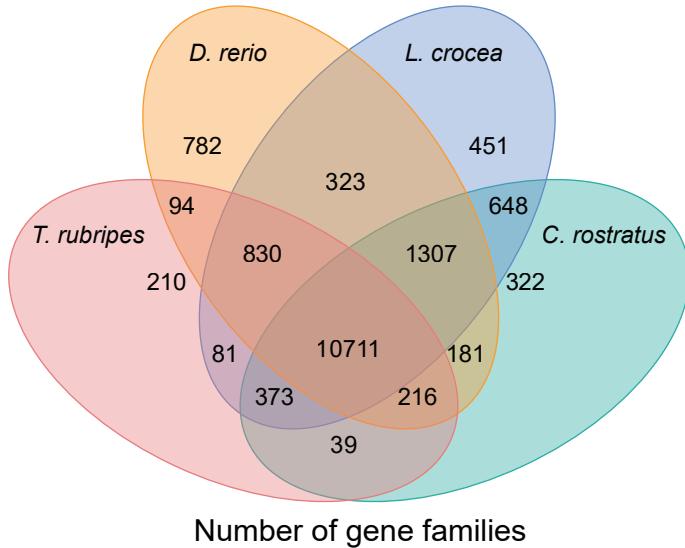
A



B



C



D

