



OPEN

DATA DESCRIPTOR

# Improved chromosomal-level genome assembly and re-annotation of leopard coral grouper

Wentao Han, Shaoxuan Wu, Hui Ding, Mingyi Wang, Mengya Wang, Zhenmin Bao, Bo Wang & Jingjie Hu

*Plectropomus leopardus*, as known as leopard coral grouper, is a valuable marine fish that has gradually been bred artificially. To promote future conservation, molecular breeding, and comparative studies, we generated an improved high-quality chromosomal-level genome assembly of leopard coral grouper using Nanopore long-reads, Illumina short reads, and the Hi-C sequencing data. The draft genome is 849.74 Mb with 45 contigs and N50 of 35.59 Mb. Finally, a total of 846.49 Mb corresponding to 99.6% of the contig sequences was anchored to 24 pseudo-chromosomes using Hi-C technology. A final set of 25,965 genes is annotated after manual curation of the predicted gene models, and BUSCO analysis yielded a completeness score of 99.5%. This study significantly improves the utility of the grouper genome and provided a reference for the study of molecular breeding, genomics and biology in this species.

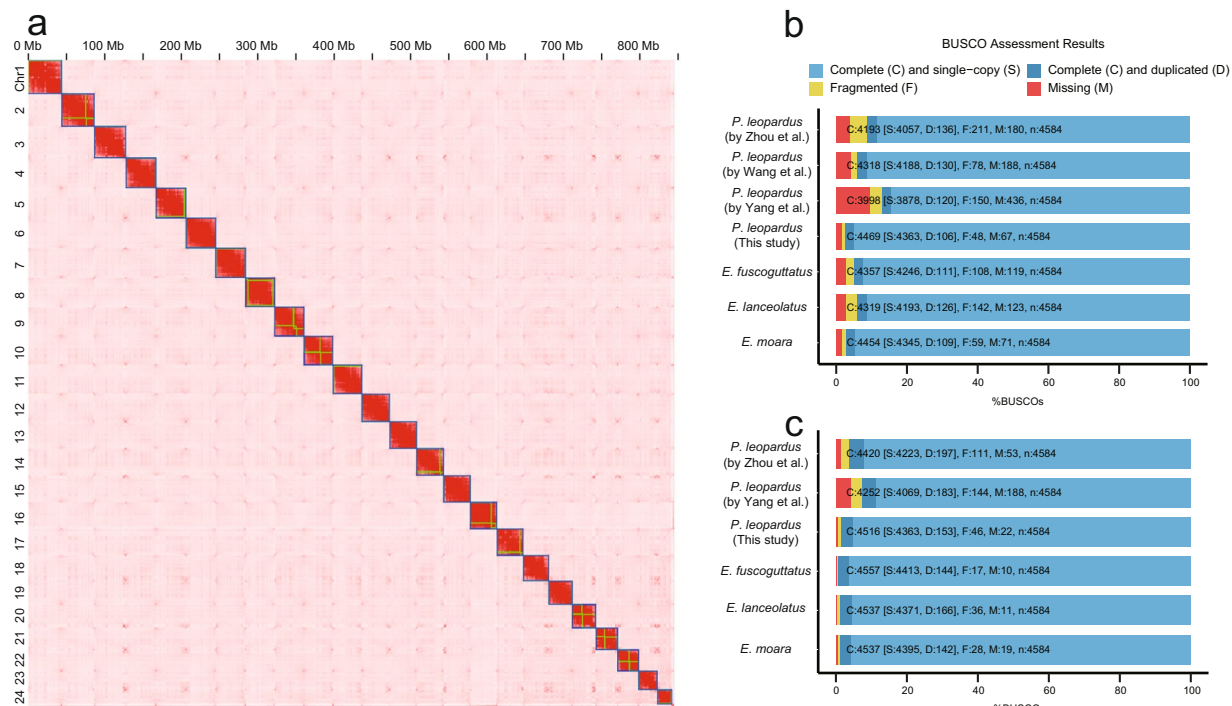
## Background & Summary

Groupers (Family Epinephelidae, Subfamily Epinephelinae) are prominent marine fishes, mostly distributed in tropical and temperate marine areas, comprising 167 species that belong to 15 genera<sup>1</sup>. Due to their high protein, low fat, tender meat quality, and good taste, groupers are high-quality economic fish species in Asia<sup>2,3</sup>. Given the huge commercial interests at stake, groupers are highly susceptible to human-induced impacts, including overfishing, making them considered threatened by the International Union for Conservation of Nature (IUCN)<sup>4</sup>. Therefore, how to scientifically develop and protect their resources has become the top priority<sup>5</sup>.

The leopard coral grouper (*Plectropomus leopardus*) has a beautiful skin color and is a valuable marine fish that commands a higher price<sup>6–8</sup>. Wild populations are suffering sharp declines due to overfishing and the destruction of spawning aggregations<sup>9</sup>. In recent years, the increasing market demands have promoted the development of artificial breeding in leopard coral grouper<sup>10–12</sup>. A high-quality reference genome resource has become increasingly important to facilitate the genomic breeding program, biological phenomena investigation and germplasm conservation<sup>13,14</sup>. Although the leopard coral grouper genome has been released<sup>6,8,15</sup>, the completeness of genome assembly and annotations still need to be further improved. For examples, the reported chromosomal-scale assembly of the sequence contigs only anchored 87.7% of the whole genome sequence using Hi-C technology<sup>6</sup>. Additionally, a wide range of gene structure annotation errors existed in the previous versions<sup>15</sup>, or the annotation information is not released and accessible to the public<sup>8</sup>.

In the present study, we generated an improved high-quality chromosome-level genome assembly of leopard coral grouper using Nanopore long-reads, Illumina short reads, and the Hi-C sequencing data. Approximately 849.74 Mb genome was assembled, consisted of 45 contigs with the contig N50 length of 35.59 Mb. A total of 846.49 Mb (99.6%) of the assembled sequences were anchored to 24 pseudo-chromosomes with low missing bases, only about 2,354 gaps. Based on this improved genome assembly, we have significantly improved upon previous gene annotations combining *de novo* prediction, homology-based searches and transcriptome-assisted methods. BUSCO alignment showed that our final assembly contained 4,469 (97.5%) complete BUSCOs. Taken together, this high-quality reference genome provides a valuable basis for the conservation and utilization of germplasm resources, and the further genetic breeding program in leopard coral grouper.

MOE Key Laboratory of Marine Genetics and Breeding, College of Marine Life Sciences/Key Laboratory of Tropical Aquatic Germplasm of Hainan Province, Sanya Oceanographic Institution, Ocean University of China, Qingdao/Sanya, China. e-mail: [wb@ouc.edu.cn](mailto:wb@ouc.edu.cn); [hujingjie@ouc.edu.cn](mailto:hujingjie@ouc.edu.cn)



**Fig. 1** Statistics on genome assembly and Comparison of four version annotations of the leopard coral grouper, *Plectropomus leopardus*. **(a)** Hi-C interaction heat map for *Plectropomus leopardus*. **(b)** BUSCO evaluation on the genome assembly completeness. **(c)** BUSCO evaluation on the predicted gene models.

	<i>P. leopardus</i>				<i>E. fuscoguttatus</i>	<i>E. lanceolatus</i>	<i>E. moara</i>
	This study	Zhou <i>et al.</i>	Wang <i>et al.</i>	Yang <i>et al.</i>			
Sequenced genome size (Mb)	849.74	881.55	913.38	787.06	1,047.01	1,087.42	1,030.48
Contig N50 (Mb)	35.59	0.86	1.41	1.14	13.80	0.12	2.22
Scaffold N50 (Mb)	38.02	34.15	40.04	33.85	44.42	46.23	43.43
Gap size (Ns per 100 kbp)	2.77	1,793.38	79.43	68.31	1.96	3,609.92	2,988.63
Complete BUSCOs (%)	97.5	91.5	94.2	87.2	95.0	94.2	97.2
Fragmented (%)	1.0	3.9	1.7	3.3	2.4	3.1	1.3
Missing (%)	1.5	4.6	4.1	9.5	2.6	2.7	1.5
Duplicate copy (%)	2.3	3.0	2.8	2.6	2.4	2.7	2.4

**Table 1.** Comparison of genome assembly metrics in groupers.

## Methods

**De novo genome assembly.** First, we estimated the genome size and heterozygosity of leopard coral grouper using GenomeScope v2.0<sup>16</sup> by *k*-mer analysis with clean Illumina short data. Program ontbc (<https://github.com/FlyPythons/ontbc>) was used to filter the Nanopore raw reads with parameters “-min\_score 7 -min\_length 1000”. Then, the filtered Nanopore reads self-corrected the base errors by the long-read assembler NextDenovo v2.3 (<https://github.com/Nextomics/NextDenovo>). Finally, clean long reads were assembled using NextDenovo v2.3 (<https://github.com/Nextomics/NextDenovo>) with the parameters: read\_cutoff = 5k' and 'seed\_cutoff = 40k'. We used purge\_dups v1.2.5<sup>17</sup> to remove the haplotypic duplication after mapping the Nanopore reads with minimap2 v2.1<sup>18</sup>. The assembly sequence was then polished using NextPolish v1.3.1<sup>19</sup> with default parameters based on Nanopore long reads. To ensure high accuracy of the genome assembly, Illumina paired-end clean reads were aligned to the assembly using BWA v0.7.15<sup>20</sup>, and the results were used to conduct another round of polishing by Pilon v1.23<sup>21</sup> with the parameters: --fix SNPs, indels. The contig-level assembly covered 849.74 Mb of the genome consisted of 45 contigs with a contig N50 value of 35.59 Mb.

**Hi-C analysis and chromosome assembly.** To obtain the chromosome-level genome, we further anchored all 45 contigs of the draft assembly onto 24 chromosomes using a 3D-DNA pipeline (version 201008)<sup>22</sup> based on the published high-quality HiC reads<sup>15</sup>. The HiC reads were aligned to the polished genome using Juicer v1.5.7 software<sup>23</sup> with default parameters. Mis-joins, order and orientation were corrected by the 3D-DNA pipeline<sup>22</sup> with the following parameters: -r 2. After the first round of 3D-DNA, we manually adjusted the

	<i>P. leopardus</i>				<i>E. fuscoguttatus</i>	<i>E. lanceolatus</i>	<i>E. moara</i>
	This study	Zhou <i>et al.</i>	Wang <i>et al.</i>	Yang <i>et al.</i>			
Number of protein-coding genes	25,965	25,763	24,700	22,317	23,813	24,067	23,588
Average gene length (bp)	15,512	15,894	16,882	20,758	22,277	21,997	21,583
Average exon length (bp)	174	171	183	276	175	174	174
Average exon number per gene	9.2	8.4	8.7	11.2	10.5	10.3	10.4
Average intron length (bp)	1,840	1,688	1,879	1,890	2,148	2,146	2,094
Percentage of repeat sequence (%)	37.35	33.91	38.02	36.18	41.28	40.17	38.85
LTR (%)	1.69	1.35	2.68	2.12	5.18	3.68	3.45
LINE (%)	3.21	2.87	3.45	3.24	4.84	4.67	4.16
SINE (%)	0.40	0.39	2.17	0.42	0.48	0.50	0.51
DNA transposons (%)	13.58	11.35	12.80	12.79	16.60	16.82	15.87

**Table 2.** Comparison of the genome-wide statistics for annotations of groupers.

ncRNA type		Copy	Proportion in Genome (%)
miRNA		746	0.075
tRNA		1,224	0.011
rRNA	18 S	152	0.023
	28 S	117	0.033
	5.8 S	22	0.001
	5 S	148	0.002
	Subtotal	439	0.059
sRNA	CD-box	135	0.002
	HACA-box	80	0.001
	Splicing	380	0.006
	Subtotal	596	0.009

**Table 3.** The statistics of functional annotation in the leopard coral grouper.

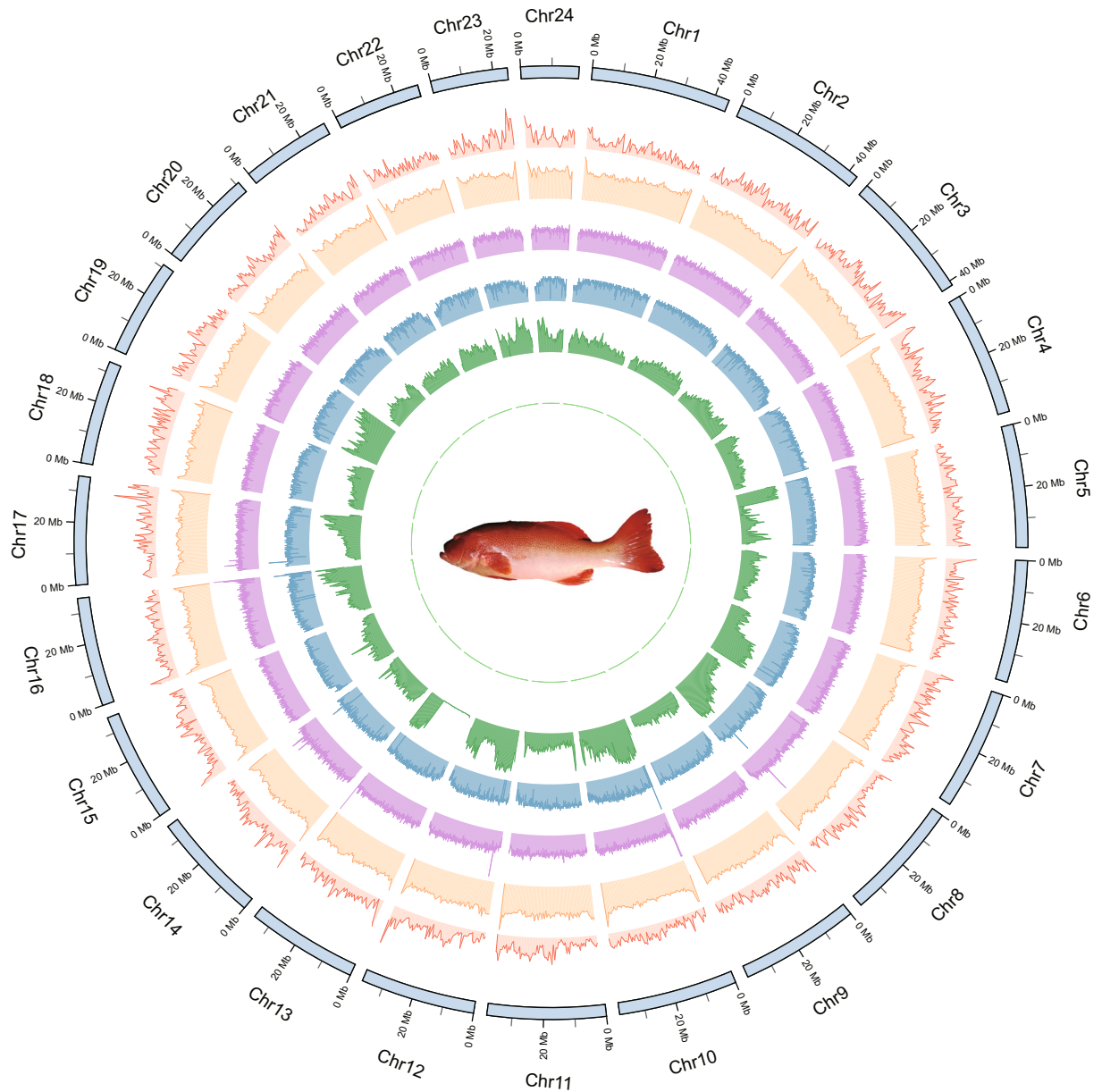
Type	Number of overall predicted genes	Percentage of overall predicted genes (%)
Total	25,965	—
SwissProt	21,331	82.2
KEGG	15,813	61.0
NR	23,027	88.7
GO	15,965	61.5
Pfam	20,201	77.8
Annotated	25,927	99.9
Unannotated	38	0.1

**Table 4.** The statistics of functional annotation in the leopard coral grouper.

assembly with Juicebox<sup>23</sup> and rerun the 3D-DNA. The Hi-C scaffolding resulted in 24 chromosome-length scaffolds (Fig. 1a).

**Repeat annotation.** *De novo* and structure-based searches were used to identify repetitive sequences with both RepeatModeler v2<sup>24</sup> (<http://www.repeatmasker.org/RepeatModeler/>) and RepeatMasker v4.0.9<sup>25</sup> (<http://www.repeatmasker.org>). Candidate LTR-RTs repetitive sequence library was identified using LTR\_finder<sup>26</sup> with parameters '-D 15000 -d 1000 -L 7000 -l 100 -p 20 -C -M 0.9' and LTRharvest v1.5.8<sup>27</sup> with parameters '-minlenltr 100 -maxlenltr 7000 -mintsd 4 -maxtsd 6 -motif TGCA -motifmis 1 -similar 85 -vic 10 -seed 20 -seqids yes'. The identified LTR-RT candidates were filtered with LTR\_retriever v2.5<sup>28</sup> program with default parameters. RepeatScout v1.0.5<sup>29</sup> LTR\_retriever v2.5<sup>28</sup> and RepeatModeler v2<sup>24</sup> were used to build *de novo* repeat libraries. The combined repeat library was used as the final library to identify repetitive sequences using RepeatMasker v4.0.9<sup>25</sup> with parameters '-q -no\_is -norna -nolow -div 40'.

**Gene prediction and annotation.** To comprehensively annotate genes, protein-coding genes prediction was undertaken using the BRAKER v2.1.5<sup>30</sup> annotation pipeline which integrated different evidence, including *de novo* prediction, homology-based searches and transcriptome-assisted methods. First, for *de novo* gene



**Fig. 2** Global genome landscape of the leopard coral grouper, *Plectropomus leopardus*. From outer to inner circles: Density of genes with 500 kbp windows, ranging from 0 to 70; GC content with 500 kbp windows, ranging from 0.30 to 0.45; depth of coverage of Nanopore reads with 100 kbp windows, ranging from 20 to 150; depth of coverage of Illumina short reads with 100 kbp windows, ranging from 10 to 35; distribution of heterozygous SNPs with 500 kbp windows, ranging from 0 to 3,420; distribution of homozygous SNPs with 500 kbp windows, ranging from 0 to 3,420.

prediction, we downloaded published RNA-seq (SRP201943<sup>31</sup> and SRP329031<sup>32</sup>) and then mapped to the soft masked genome using Hi-SAT2 v. 2.1.0<sup>33</sup>. Then, all mapping results were used to build transcript models using BRAKER v2.1.5<sup>30</sup> and StringTie v2.1.6<sup>34</sup>. BRAKER v2.1.5<sup>30</sup> was run with Semi-HMM-based Nucleic Acid Parser (SNAP, v2013.11.29)<sup>35</sup> and Augustus v3.3.3<sup>36</sup> which pre-trained using released gene models of *P. leopardus*<sup>6,15</sup>. Second, protein-coding sequences of from *P. leopardus*<sup>6,15</sup>, *E. fuscoguttatus*<sup>37</sup>, *E. lanceolatus*<sup>38</sup>, and *E. moara*<sup>39</sup> were aligned to the genome assembly using TBLASTN and GeneWise v2.2.0<sup>40</sup>. Third, Trinity v2.1.1<sup>41</sup> was used to generate the transcripts. The transcriptome data were further assembled using the PASA pipeline v2.5.2<sup>42</sup> with BLAT v35<sup>43</sup> and GMAP (version 20150921)<sup>44</sup> as the aligner. Finally, all evidences were merged to form a consensus gene set using EvidenceModeler v1.1.1<sup>45</sup>. Finally, we identified a total of 25,965 protein-coding genes (Table 2). The non-coding RNA genes including rRNAs, tRNAs, snRNAs and miRNAs were screened using INFERNAL v 1.1.2<sup>46</sup> and tRNAscan-SE v1.4<sup>47</sup>. Four types of noncoding RNAs, including 746 miRNAs, 1,224 tRNAs, 439 rRNAs and 596 sRNAs, were identified from the *P. leopardus* genome (Table 3).

Type	Number	Percentage (%)
All SNP	2,326,997	0.2738
Heterozygous SNP	2,320,097	0.2730
Homologous SNP	6,900	0.0008

**Table 5.** The statistics of the leopard coral grouper (MGB\_pleo\_1.0) SNPs.

In order to explore the function of predicted protein-coding genes in leopard coral grouper, InterPro30, Pfam32, PANTHER 14.1, Superfamily 1.75, Gene3D 4.2.0, SMART 7.1 and TrEMBL32 databases were respectively used to predict protein function based on the conserved protein domains by InterProScan v5.36<sup>48</sup>. We performed functional annotation by aligning the protein sequences to NCBI nr databases and SwissProt using BLASTP. The result showed more than 99.9% (25,927) of protein-coding genes were annotated (Table 4).

### Data Records

The assembled genome has been deposited at GenBank under the accession GCA\_026936395.1<sup>49</sup>. Moreover, the whole genome sequence data reported in this paper have been deposited in the Genome Warehouse in National Genomics Data Center<sup>50,51</sup>, Beijing Institute of Genomics, Chinese Academy of Sciences/China National Center for Bioinformation, under accession number GWHBPCI00000000 that is publicly accessible at <https://ngdc.cncb.ac.cn/gwh/Assembly/29542/show><sup>52</sup>. In addition, the genome annotation files had been submitted at the figshare<sup>53</sup>. The Nanopore long reads, Illumina genomic sequencing data and Hi-C data were downloaded from CNGBdb<sup>51,54</sup> under the accession CNP0000859<sup>55</sup>. Transcriptomic sequences can be retrieved under the following accession numbers: SRP201943<sup>31</sup> and SRP329031<sup>32</sup>.

### Technical Validation

To evaluate the quality of genome assembly, first, we assessed genome continuity with QUASt v5.0.2<sup>56</sup>. Contig N50 (the length such that half of all sequence is in contigs of this size) has achieved a significant improvement to 35.59 Mb, which is much higher than other versions<sup>6,8,15</sup> or closely related species (*Epinephelus fuscoguttatus*, *Epinephelus lanceolatus*, *Epinephelus moara*) assembled with long-read sequencing from 0.12 to 13.8 Mb. Meanwhile, in the latest version, there are very few gaps in the genome (2.77 per 100 kbp), which is remarkably less than the previous from 68.31 per 100 kbp to 1793.38 per 100 kbp<sup>6,8,15</sup> (Table 1; Fig. 2). Second, Illumina paired-end clean reads and Nanopore long reads were mapped to the final reference genome assembly by using BWA v0.7.15<sup>20</sup> and Minimap2 v2.1<sup>18</sup>, respectively. The mapping rate of Illumina and Nanopore reads reached 99.18% and 99.95%. We only detected 6,900 (0.0008%) conflicting sites in the final assembly, indicating that this is a high level of the complete genome (Fig. 2; Table 5). Finally, we evaluated the completeness of our genome assembly using Benchmarking Universal Single-Copy Orthologs (BUSCO, v3.0)<sup>57</sup> with the actinopterygii\_odb9 database. The actinopterygii\_odb9 database contained 4,584 conserved core genes while our assembled genome contained 4,469 (97.5%) of the expected actinopterygii genes (including 4,393 (95.2%) single and 106 (2.3%) duplicated ones). Obviously, our data had complete gene coverage, and 48 (1.0%) were identified as fragmented, respectively, while 67 (1.5%) were missing in our assembled genome (Fig. 1b). Furthermore, we also used BUSCO to evaluate the completeness of gene annotations<sup>57</sup>, and only 22 (0.5%) genes were missing in the final annotation version (Fig. 1c) Table 5.

### Code availability

The data analyses were performed according to the manuals by the developers of corresponding bioinformatics tools and all software, and codes used in this work are publicly available, with corresponding versions indicated in Methods.

Received: 22 December 2022; Accepted: 6 March 2023;

Published online: 22 March 2023

### References

- Félix-Hackradt, F. C., Hackradt, C. W. & García-Charton, J. A. *Biology and Ecology of Groupers*. (CRC Press, 2022).
- Fabinyi, M. Historical, cultural and social perspectives on luxury seafood consumption in China. *Environ. Conserv.* **39**, 83–92 (2012).
- Sale P. F. *Coral reef fishes: dynamics and diversity in a complex ecosystem*. (Academic Press, 2002).
- Lui, O. J., Woods, R. M., Madin, E. M. P. & Madin, J. S. Predicting IUCN Extinction Risk Categories for the World's Data Deficient Groupers (Teleostei: Epinephelidae). *Conserv. Lett.* **9**, 342–350 (2016).
- Valderrama, S. P. *et al.* Marine protected areas in Cuba. *B. Mar. Sci.* **94**, 423–442 (2018).
- Zhou, Q. *et al.* De novo sequencing and chromosomal-scale genome assembly of leopard coral grouper, *Plectropomus leopardus*. *Mol. Ecol. Resour.* **20**, 1403–1413 (2020).
- Wang, L., Yu, C. P., Guo, L., Lin, H. R. & Meng, Z. N. In silico comparative transcriptome analysis of two color morphs of the common coral trout (*Plectropomus leopardus*). *PLoS One* **10**, e0145868 (2015).
- Yang, Y. *et al.* Whole-genome sequencing of leopard coral grouper (*Plectropomus leopardus*) and exploration of regulation mechanism of skin color and adaptive evolution. *Zool. Res.* **41**, 328 (2020).
- Agustina, S., Panggabean, A. S., Natsir, M., Retroningtyas, H. & Yulianto, I. Yield-per-recruit modeling as biological reference points to provide fisheries management of Leopard Coral Grouper (*Plectropomus leopardus*) in Saleh Bay, West Nusa Tenggara. *IOP Conference Series: Earth and Environmental Science* **278**, 012005 (2019).
- Ottolenghi, F., Silvestri, C., Giordano, P., Lovatelli, A. & New, M. B. *Capture-based aquaculture: the fattening of eels, groupers, tunas and yellowtails*. (FAO, 2004).

11. Nguyen, T. T. T., Davy, F. B., Rimmer, M. A. & De Silva, S. S. Use and exchange of genetic resources of emerging species for aquaculture and other purposes. *Rev Aquacult.* **1**, 260–274 (2009).
12. Kongkeo, H., Wayne, C., Murdjani, M., Bunliptanon, P. & Chien, T. Current practices of marine finfish cage culture in China, Indonesia, Thailand and Vietnam. *Aquac. Asia* **15**, 32–40 (2010).
13. Allendorf, F. W., Hohenlohe, P. A. & Luikart, G. Genomics and the future of conservation genetics. *Nat. Rev. Genet.* **11**, 697–709 (2010).
14. Mohanty, B. P. *et al.* Omics technology in fisheries and aquaculture. *Adv. Fish Res.* **7**, 1–30 (2019).
15. Wang, Y. B. *et al.* Chromosome genome assembly of the leopard coral grouper (*Plectropomus leopardus*) with Nanopore and Hi-C sequencing data. *Front. Genet.* **11** (2020).
16. Ranallo-Benavidez, T. R., Jaron, K. S. & Schatz, M. C. GenomeScope 2.0 and Smudgeplot for reference-free profiling of polyploid genomes. *Nature Communications* **11**, 1432 (2020).
17. Guan, D. F. *et al.* Identifying and removing haplotypic duplication in primary genome assemblies. *Bioinformatics.* **36**, 2896–2898 (2020).
18. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
19. Hu, J., Fan, J. P., Sun, Z. Y. & Liu, S. L. NextPolish: a fast and efficient genome polishing tool for long-read assembly. *Bioinformatics* **36**, 2253–2255 (2019).
20. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
21. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS one* **9**, e112963 (2014).
22. Dudchenko, O. *et al.* De novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length scaffolds. *Science* **356**, 92–95 (2017).
23. Durand, N. C. *et al.* Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst.* **3**, 95–98 (2016).
24. Flynn, J. M. *et al.* RepeatModeler2 for automated genomic discovery of transposable element families. *Proc. Natl. Acad. Sci.* **117**, 9451–9457 (2020).
25. Chen, N. S. Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr. Protoc. Bioinformatics* **5**, 4.10.11–14.10.14 (2004).
26. Xu, Z. & Wang, H. LTR\_FINDER: an efficient tool for the prediction of full-length LTR retrotransposons. *Nucleic Acids Res.* **35**, W265–W268 (2007).
27. Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible software for de novo detection of LTR retrotransposons. *BMC Bioinformatics* **9**, 18 (2008).
28. Ou, S. & Jiang, N. LTR\_retriever: A Highly accurate and sensitive program for identification of long terminal repeat retrotransposons. *Plant Physiol.* **176**, 1410–1422 (2017).
29. Price, A. L., Jones, N. C. & Pevzner, P. A. De novo identification of repeat families in large genomes. *Bioinformatics* **21**, i351–i358 (2005).
30. Hoff, K. J., Lomsadze, A., Borodovsky, M. & Stanke, M. in *Gene Prediction: Methods and Protocols* (ed M., Kollmar) 65–95 (Springer New York, 2019).
31. *NCBI Sequence Read Archive* <https://identifiers.org/insdc.sra:SRP201943> (2021).
32. *NCBI Sequence Read Archive* <https://identifiers.org/insdc.sra:SRP329031> (2021).
33. Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory requirements. *Nat. Methods* **12**, 357–360 (2015).
34. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).
35. Korf, I. Gene finding in novel genomes. *BMC Bioinformatics* **5**, 59 (2004).
36. Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res.* **34**, W435–W439, <https://doi.org/10.1093/nar/gkl200> (2006).
37. Yang, Y. *et al.* Whole-genome sequencing of brown-marbled grouper (*Epinephelus fuscoguttatus*) provides insights into adaptive evolution and growth differences. *Mol. Ecol. Resour.* **22**, 711–723 (2022).
38. Zhou, Q. *et al.* A chromosome-level genome assembly of the giant grouper (*Epinephelus lanceolatus*) provides insights into its innate immunity and rapid growth. *Mol. Ecol. Resour.* **19**, 1322–1332 (2019).
39. Zhou, Q., Gao, H. Y., Xu, H., Lin, H. R. & Chen, S. L. A Chromosomal-scale reference genome of the kelp grouper *Epinephelus moara*. *Mar. Biotechnol.* **23**, 12–16 (2021).
40. Birney, E., Clamp, M. & Durbin, R. GeneWise and genomewise. *Genome Res.* **14**, 988–995 (2004).
41. Grabherr, M. G. *et al.* Trinity: reconstructing a full-length transcriptome without a genome from RNA-Seq data. *Nat. Biotechnol.* **29**, 644 (2011).
42. Haas, B. J. *et al.* Improving the *Arabidopsis* genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res.* **31**, 5654–5666 (2003).
43. Kent, W. J. BLAT—the BLAST-like alignment tool. *Genome Res.* **12**, 656–664 (2002).
44. Wu, T. D. & Watanabe, C. K. GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics* **21**, 1859–1875 (2005).
45. Haas, B. J. *et al.* Automated eukaryotic gene structure annotation using EVidenceModeler and the Program to Assemble Spliced Alignments. *Genome Biol.* **9**, R7 (2008).
46. Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics* **29**, 2933–2935 (2013).
47. Lowe, T. M. & Chan, P. P. tRNAscan-SE On-line: integrating search and context for analysis of transfer RNA genes. *Nucleic Acids Res.* **44**, W54–W57 (2016).
48. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification. *Bioinformatics* **30**, 1236–1240, <https://doi.org/10.1093/bioinformatics/btu031> (2014).
49. *NCBI Assembly* [https://identifiers.org/ncbi/insdc.gca:GCA\\_026936395.1](https://identifiers.org/ncbi/insdc.gca:GCA_026936395.1) (2022).
50. Chen, M. L. *et al.* Genome warehouse: a public repository housing genome-scale data. *Genom. Proteom. Bioinforma.* **19**, 584–589 (2021).
51. CNCB-NGDC Members and Partners Database Resources of the National Genomics Data Center, China National Center for Bioinformatics in 2022. *Nucleic Acids Res.* **50**, D27–D38 (2021).
52. *National Genomics Data Center* <https://ngdc.cncb.ac.cn/gwh/Assembly/29542/show> (2022).
53. Han, W. *Plectropomus leopardus* genome. *Figshare* <https://doi.org/10.6084/m9.figshare.21441396.v3> (2022).
54. FAIRsharing.org: CNGBdb; China National GeneBank DataBase; <https://doi.org/10.25504/FAIRsharing.9btRvC>.
55. Zhang X. & Institute of Biodiversity Conservation. leopard coral grouper genome. *CNGBdb* <https://db.cngb.org/search/project/CNP0000859/> (2020).
56. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for genome assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
57. Waterhouse, R. M. *et al.* BUSCO applications from quality Assessments to gene prediction and phylogenomics. *Mol. Biol. Evol.* **35**, 543–548 (2017).

## Acknowledgements

This research was funded by the National Key Research and Development Program of China (2022YFD2400501), the Project of Sanya Yazhouwan Science and Technology City Management Foundation (SKJC-2020-02-009), the Key R&D Project of Hainan Province (ZDYF2021XDNY133), and the China Postdoctoral Science Foundation (2021703030).

## Author contributions

J.H., B.W. and Z.B. conceived and designed the study. J.H. and B.W. coordinated and supervised the whole study. W.H. conducted the genome assembly and analysis. S.W., M.W., H.D. and M.W. participated in discussions and provided suggestions for manuscript improvement. W.H., B.W. and J.H. did most of the writing with input from other authors.

## Competing interests

The authors declare no competing interests.

## Additional information

**Correspondence** and requests for materials should be addressed to B.W. or J.H.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023