

1
2
3
4 **1 Draft genome assembly of the Aral barbell *Luciobarbus brachycephalus* using**
5
6 **2 PacBio sequencing**
7
8

9
10
11
12 4 Longwu Geng^{1#}, Ming Zou^{2#}, Haifeng Jiang³, Minghui Meng⁴, Wei Xu^{1*}
13

14 5 ¹Heilongjiang River Fisheries Research Institute, Chinese Academy of Fishery Sciences, Key
15
16 6 Open Laboratory of Cold Water Fish Germplasm Resources and Breeding of Heilongjiang
17
18
19 7 Province, Harbin, China
20
21

22 8 ²Institute of Zoology, Chinese Academy of Sciences, Beijing, China
23

24 9 ³College of Animal Science and Technology, Northwest A&F University, Xiong Road 22nd,
25
26
27 10 Yangling, China
28
29

30 11 ⁴Diggs (Wuhan) Biotechnology Co., Ltd.
31

32 12 #These authors have contributed equally to this work
33
34

35 13 *Corresponding author: E-mail: xuwei@hrfri.ac.cn.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58

1
2
3
4 **Abstract**
5

6
7 15 The endangered Aral barbell *Luciobarbus brachycephalus* is endemic to the water
8
9 16 systems of the Caspian Sea and Aral Sea. Given the scarcity of genetic data for
10
11 17 the species, we present a draft assembly based on PacBio long read sequencing
12
13
14 18 technology. Approximate 299.4 Gb of long reads representing 166X of the
15
16
17 19 estimated genome size were generated, and the final assembly was composed of
18
19 20 653 contigs totaling approximately 1,698.3 Mb, with a contig N50 length of 4.5
20
21
22 21 Mb. A total of 807.6 Mb represented approximately 47.6% of the assembly and
23
24 22 were identified as repeats. Fifty-four thousand and six hundred possible protein
25
26
27 23 genes were predicted, among which 50,727, representing approximately 92.9%,
28
29
30 24 could be annotated by at least one database. Evolutionary analysis showed that *L.*
31
32 25 *brachycephalus* and *Labeo rohita* diverged by approximately 42.6 Mya, and the
33
34
35 26 obvious expansion of gene families residing in the *L. brachycephalus* genome
36
37
38 27 may be attributed to the specific whole genome duplication of the species. The
39
40 28 first genome assembly of *L. brachycephalus* can not only provide a foundation
41
42
43 29 for genetic conservation and molecular breeding of this species but also contribute
44
45
46 30 to comparative analyses of genome biology and evolution within Cyprinidae.

47
48 31 **Key words:** *Luciobarbus brachycephalus*, PacBio sequencing, *de novo*
49
50 32 assembly, genome annotation, phylogeny.
51
52
53
54
55
56
57
58
59
60

1
2
3
4 34 **Significance**

5
6 35 Aral barbell *Luciobarbus brachycephalus* is an endangered fish species endemic
7
8
9 36 to the water systems of the Caspian Sea and Aral Sea. At present, genetic and
10
11 37 genomic information for the species and the genus *Brachycephalus* is limited. In
12
13
14 38 this study, we obtained a draft genome assembly of *L. brachycephalus*, which will
15
16
17 39 contribute to research on the genomics, evolution, and genetic breeding of this
18
19
20 40 species and the largest and most diverse fish family, Cyprinidae.
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

41 **Introduction**

42 Aral barbell, *Luciobarbus brachycephalus*, is a valuable fish species native to the
43 Aral basin, Chu drainage and southern and western Caspian Sea (Coad 1998). It
44 has been speculated that the population of *L. brachycephalus* has declined by at
45 least 30% in the past 30 years and continues to decline due to the shrinking
46 (increased salinity) of the Aral Sea and damming of its tributaries. Aral barbel is
47 currently listed on the IUCN Red List as a vulnerable (VN) species (Esmaeili, et
48 al. 2017). However, a lack of information on genetic and genomic information for
49 this endangered species hinders efforts to protect, restore and effectively manage
50 the Aral barbell population.

51 The Aral barbell belongs to the Barbinae subfamily, family Cyprinidae
52 (Cypriniformes), in which polyploidization evolved independently on multiple
53 occasions (Xu, et al. 2019). Cyprinids contain many autoallopolyploid fish
54 species, resulting in complicated evolutionary histories and phylogenetic
55 relationships. The chromosome number of *L. brachycephalus* is 100 (Geng, et al.
56 2013), twice that in the common ancestor of Barbinae. It appears that *L.*
57 *brachycephalus* has experienced a polyploidization event beyond teleost-specific
58 whole genome duplication (WGD). Polyploidization plays a significant role in
59 speciation and adaptive evolution since it produces redundant genes, which can
60 serve as an important genetic material basis for complex phenotypes (Xu, et al.
61 2019). Unlike plant lineages where polyploidization occurs frequently (Mayrose
62 and Lysak 2020; Wu, et al. 2021), vertebrate polyploidization is relatively rare,

1
2
3
4 63 and is mainly concentrated in some amphibians and fish (Chen, et al. 2020). The
5
6 64 polyploid history of *L. brachycephalus* provides an excellent model to explore the
7
8
9 65 evolutionary patterns and histories of polyploidization in vertebrates.
10

11
12 66 *L. brachycephalus* has high potential for aquaculture exploitation in saline
13
14 67 alkali water due to its strong saline-alkali tolerance (salinity < 10 g/L, alkalinity
15
16 68 < 30 mmol/L) (Geng, et al. 2016). Therefore, *L. brachycephalus* was introduced
17
18
19 69 into China from Uzbekistan as the name large-scale barbell in 2003 (Jiang, et al.
20
21
22 70 2019). Since then, a great deal of studies on artificial breeding, biological
23
24 71 performance and larval rearing have been conducted (Li, et al. 2019; Longwu, et
25
26
27 72 al. 2010). At present, *L. brachycephalus* is becoming more economically
28
29
30 73 important owing to its taste, fast growth and high commercial value, and it has
31
32
33 74 been cultured in more than 20 provinces in China with an annual production of up
34
35 75 to 20,000 - 40,000 tons. However, the genetic resources of *L. brachycephalus* are
36
37
38 76 relatively limited, and the genetic diversity is generally low due to founder effects
39
40
41 77 and the inbreeding of small populations. Therefore, it is urgent to develop
42
43 78 genomic resources to accelerate the genome-assisted breeding of *L.*
44
45 79 *brachycephalus*. However, to date, only the mitochondrial genome and
46
47
48 80 microsatellite DNA have been reported for the Aral barbell (Jiang, et al. 2019;
49
50
51 81 Longwu, et al. 2012). Here, we report a draft genome assembly of *L.*
52
53 82 *brachycephalus*, which will provide valuable genomic resources for studies on
54
55
56 83 conservation and breeding as well as polyploid origin, speciation and adaptation
57
58
59 84 in polyploid Cyprinidae.
60

85 **Results and Discussion**

86 **Genome Sequencing and Assembly**

87 Based on the genomic data generated in the present study, the whole genome of
88 *L. brachycephalus* was assembled. More than 675,436,128 paired-end reads with
89 a length of 150 bp totaling more than 101.3 Gb were generated for the survey
90 analysis (supplementary table S1). The estimated genome size for the species was
91 approximately 1,806.6 Mb (supplementary table S2). The features of the species
92 were similar to those of polyploid-related species such as *C. carpio* and
93 *Schizothorax o'connori* (Xiao, et al. 2020; Xu, et al. 2014). Thus, the genome of
94 *L. brachycephalus* may be tetraploid. The *de novo* assembly following rounds of
95 polishing generated an assembly composed of 653 contigs totaling 1698.3 Mb,
96 with a contig N50 length of 4.5 Mb (supplementary table S3). BUSCO assessment
97 showed that approximately 96.0% of the BUSCO genes were recovered
98 completely by assembly, of which 37.8% were single copy and 58.2% were
99 duplicated (supplementary table S4). The assembly may partially recover another
100 1.6% of the BUSCO genes (supplementary table S4). In this study, approximately
101 40 M paired-end reads totaling nearly 7 Gb were generated each for the brain,
102 heart, kidney, liver, and ovary (supplementary table S1). The mapping rate for the
103 RNA-Seq data to the new assembly reached more than 80% for each of the tissues.
104 Thus, the assembly of the *L. brachycephalus* genome may be of high quality.

105 **Genome Annotation**

106 A number of repeat elements and protein-coding genes residing in the new

1
2
3
4 107 assembly were identified in the study. All the identified repetitive elements
5
6 108 occupied 47.55% of the genome (supplementary tables S5 and S6), which was
7
8
9 109 higher than that of *C. carpio* but lower than that of *S. o'connori* (Xiao, et al. 2020;
10
11 110 Xu, et al. 2014). Briefly, RepeatModeler + RepeatMasker identified that
12
13
14 111 approximately 33.8% of the genome was repeats. Then, on the basis of consensus
15
16
17 112 sequences deposited in Repbase, RepeatMasker identified approximately 28.5%
18
19
20 113 of the genome as repeats (supplementary table S6). RepeatProteinMask and
21
22 114 LTR_retriever identified approximately 11.4% and 6.1% of the genomes as
23
24
25 115 possible TE proteins and LTRs, respectively (supplementary table S6). Moreover,
26
27 116 more than 1.3 M SSRs were identified from the assembly of the *L. brachycephalus*
28
29
30 117 genome.

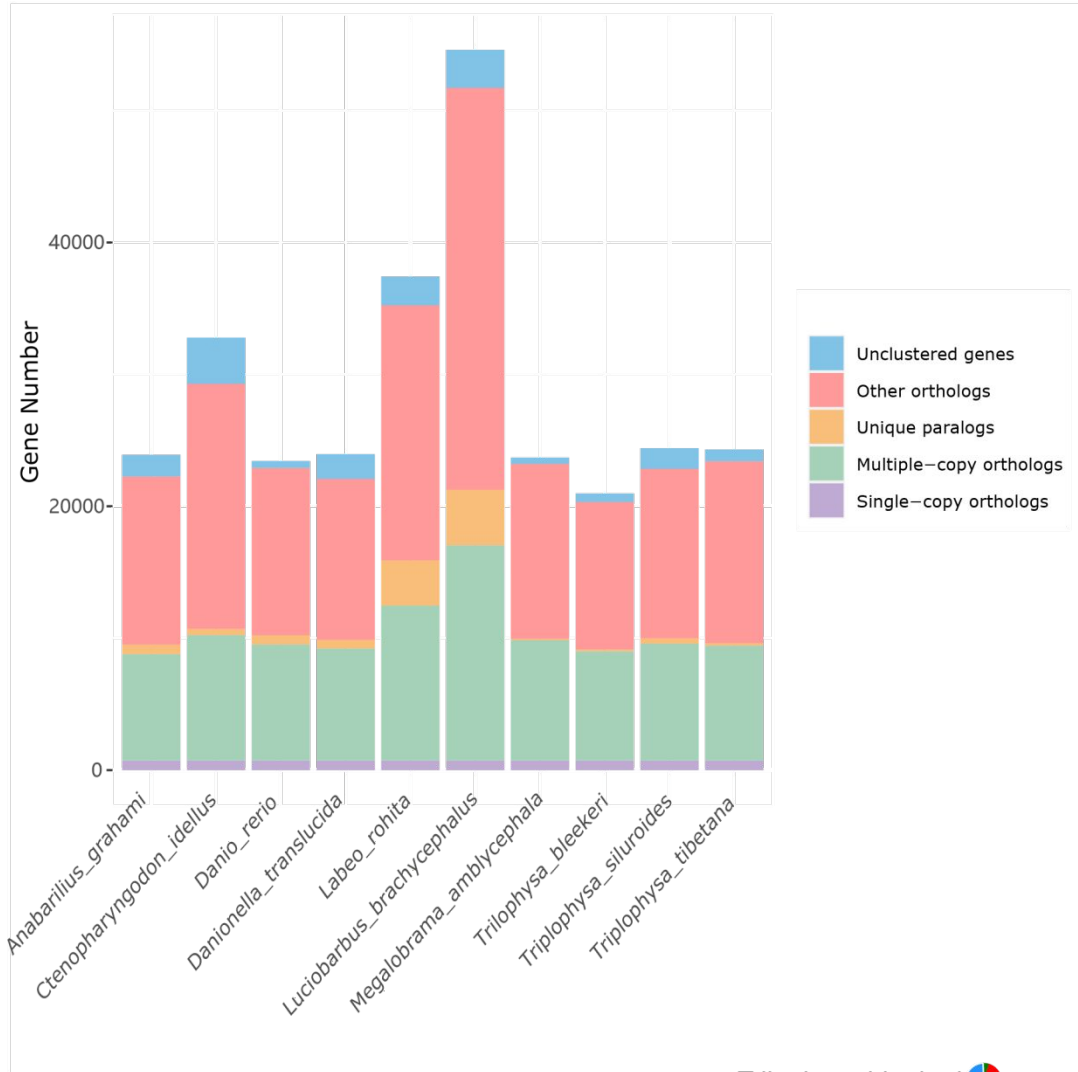
31
32
33 118 The process of gene-model annotation identified a total of 54,600 protein-
34
35 119 coding genes residing in the *L. brachycephalus* genome with a mean transcript
36
37
38 120 length of 15,737.16 bp. The gene number of *L. brachycephalus* is much greater
39
40
41 121 than that of diploid relatives such as *D. rerio* and *C. idellus* (Howe, et al. 2013;
42
43 122 Wang, et al. 2015) but is similar to that of tetraploids *C. carpio* and *S. o'connori*
44
45 123 (Xiao, et al. 2020; Xu, et al. 2014). A total of 50,727 genes accounting for
46
47
48 124 approximately 92.9% of the total genes were annotated against at least one
49
50
51 125 database (supplementary table S7). Among these databases, NR annotated 50,477
52
53 126 genes, accounting for 92.5% of the total genes, followed by SwissProt, which
54
55
56 127 annotated 81.4% (supplementary table S7).

58 128 **Comparative Genome Analysis**

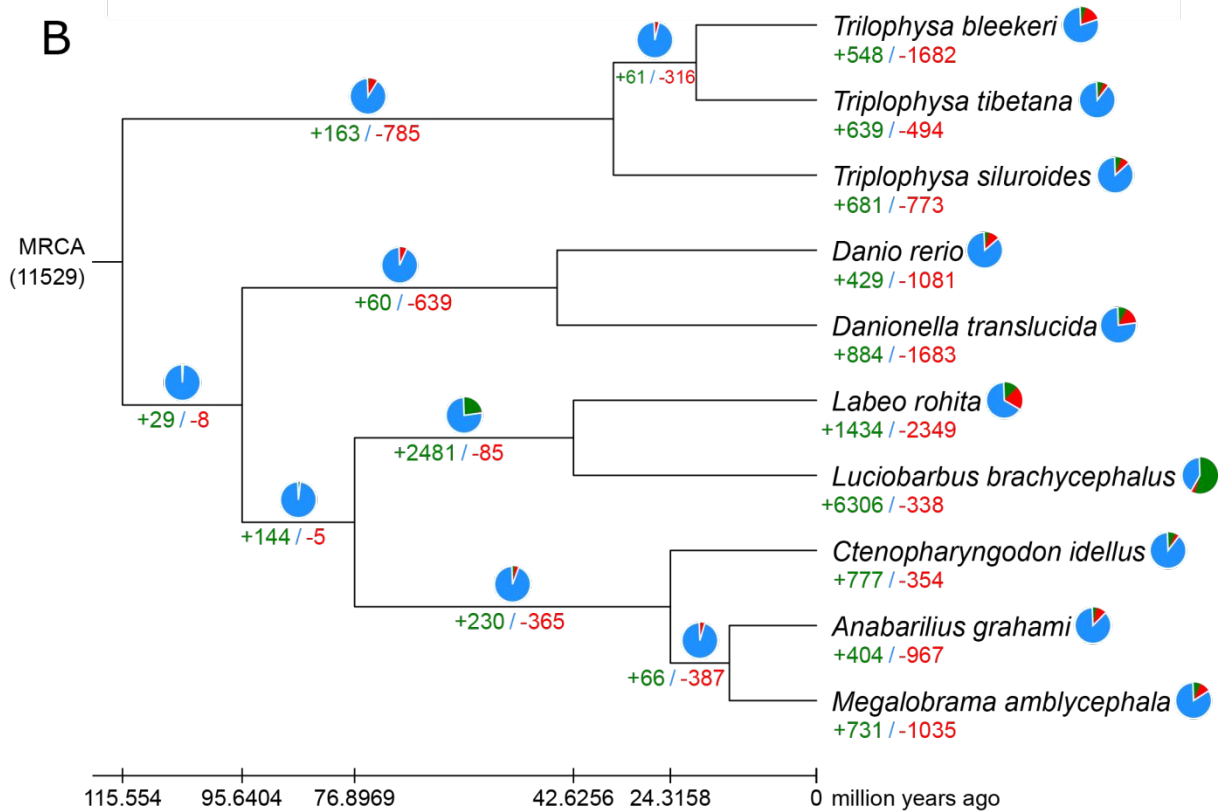
59
60

1
2
3
4 129 The evolutionary dynamics related to *L. brachycephalus* were deduced by
5
6 130 comparison to the relative species. Seven hundred thirty-two one-to-one orthologs
7
8
9 131 were identified among these species (fig. 1), and the concatenated supermatrix of
10
11 132 the alignments of these genes suggested that *L. brachycephalus* and *L. rohita*
12
13
14 133 formed sister groups with high confidence (fig. 1). Divergence time estimation
15
16 134 suggested that *L. brachycephalus* and *L. rohita* diverged by approximately 42.6
17
18
19 135 Mya. Approximately 6,306 and 338 gene families were probably subject to
20
21
22 136 expansions and contractions in the species, respectively (fig. 1). The obvious
23
24
25 137 expansion of gene families within the *L. brachycephalus* genome may coincide
26
27 138 with its whole genome duplication (Voltaire, et al. 2017).
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

A



B



139

1
2
3
4 140 **FIG. 1.**—Comparison between *L. brachycephalus* and relative species. (A)
5
6 141 Comparison of copy numbers in gene clusters reside in the genomes of *L. brachycephalus* and
7
8 142 the other 9 relative fish species. Single-copy orthologs denote the family have and only have
9
10 143 one gene for each species, and multi-copy orthologs denote the family clustered more than one
11
12 144 gene for each species. Other orthologs denote the family can have any number of genes for each
13
14 145 species except the single-copy and multi-copy orthologs. Unique paralogs denote species-
15
16 146 specific gene families, and Unclustered genes denote species-specific genes cannot cluster with
17
18 147 any other genes. (B) Phylogenetic relationships between *L. brachycephalus* and relatives
19
20 148 recovered based on the maximum likelihood method. All nodes are fully supported by the
21
22 149 bootstrap resampling test. The numbers and sectors marked with green, red, and blue denote
23
24 150 gene families subject to expansion, contraction, and stability for each species, respectively.

25 26 151 **Conclusions**

27
28 152 In the present study, we sequenced and *de novo* assembled the whole genome of
29
30 153 tetraploid *L. brachycephalus* based on PacBio sequencing technology. The
31
32 154 estimated genome size of the species is approximately 1,806.6 Mb, and the
33
34 155 assembly recovered more than 94.0% of the estimated size. Approximately 47.55%
35
36 156 of the assembly may be composed of repeat elements, and a total of 54,600 protein
37
38 157 genes were identified. The first genome sequences for the species should be a
39
40 158 valuable resource for studies on the genomics, evolution, and ecology of polyploid
41
42 159 fish in the family Cyprinidae.

43 44 160 **Materials and Methods**

45 46 161 **Sample Collection and Sequencing**

47
48 162 An individual sample of *L. brachycephalus* was obtained from a breeding farm at
49
50 163 Hulan Experimental Station, Heilongjiang River Fisheries Research Institute
51
52 164 (Harbin, P. R. China), for genome sequencing. After the injection of tricaine

1
2
3
4 165 methanesulfonate (MS-222), it was dissected, and a number of tissues, including
5
6 166 white muscle, brain, heart, kidney, liver, and ovary, were sampled and
7
8
9 167 immediately frozen in liquid nitrogen. For genome sequencing, genomic DNA
10
11 168 from white muscle was extracted using a standard phenol/chloroform method.
12
13
14 169 The integrity of the extracted DNA was assessed by 0.75% agarose gel
15
16
17 170 electrophoresis, and the concentration was quantified by a Qubit 4 Fluorometer
18
19
20 171 (Thermo Fisher Scientific, Inc., USA). Ten micrograms of DNA was then used to
21
22 172 construct the library for PacBio SMRT Sequencing using the SMRTbell Express
23
24
25 173 Template Prep Kit (PacBio, Menlo Park, CA, USA), and the library was
26
27 174 sequenced using the PacBio Sequel II System with CLR mode. To estimate the
28
29
30 175 profiles, including genome size, heterozygosity, and repeat contents of the
31
32
33 176 genome, a paired-end library with a 400 bp insert size was constructed using
34
35 177 genomic DNA with an Illumina TruSeq DNA Nano Preparation Kit (Illumina,
36
37
38 178 San Diego, CA, USA) and sequenced using the Illumina HiSeq Xten platform
39
40
41 179 (Illumina). To assist the annotation of gene models, total RNA from all the
42
43 180 sampled tissues, that is, brain, heart, kidney, liver, and ovary, was extracted using
44
45
46 181 TRIzol reagent (Invitrogen, Carlsbad, CA) following the manufacturer's protocol,
47
48 182 and libraries with 400 bp insert sizes were prepared and sequenced based on the
49
50
51 183 Illumina HiSeq Xten platform. All experiments involving the handling and
52
53 184 treatment of fish in this study were approved by the Animal Care and Use
54
55
56 185 Committee of the HRFRI of the Chinese Academy of Fishery Sciences.

58 186 **Genome Assembly and Assessment**

1
2
3
4 187 The Illumina short reads generated for the survey were quality controlled using
5
6 188 fastp (Chen, et al. 2018) and were used to estimate the genome profiles using
7
8
9 189 GenomeScope2 (Ranallo-Benavidez, et al. 2020) based on the 17-mer counts
10
11
12 190 generated by Jellyfish2.0 (Marcais and Kingsford 2011). Based on the estimated
13
14 191 genome size, the genomic long reads were used to *de novo* assemble the genome
15
16
17 192 using FALCON 0.3.0 (Chin, et al. 2016) with the following parameters:
18
19 193 “length_cutoff = 21000, length_cutoff_pr = 20000, pa_HPCdaligner_option = -v
20
21
22 194 -B4 -t16 -e.70 -l1000 -s1000 -T12 -M50, overlap_filtering_setting = --bestn 10 -
23
24
25 195 -n_core 16 --min_cov 3 --max_diff 100 --max_cov 100”. To improve the
26
27 196 assembly, all the long reads were mapped back to the new assembly using blasr
28
29
30 197 v5.3.1 with default settings (Chaisson and Tesler 2012) and were used to polish
31
32
33 198 the genome using gcpp2 (<https://github.com/PacificBiosciences/gcpp>). The
34
35 199 improved assembly was also polished using pilon v1.8 (Walker, et al. 2014) based
36
37
38 200 on the alignments of genomic short reads mapped back to the assembly using
39
40
41 201 BWA-MEM with default settings (Li 2013). The process was repeated 3 times
42
43 202 iteratively to try to adjust the assembly to increase its accuracy. Evolutionarily
44
45 203 informed expectations of the gene content of near-universal single-copy orthologs
46
47
48 204 within Actinopterygii were estimated using Benchmarking Universal Single-
49
50 205 Copy Orthologs (BUSCO 3) software to estimate the completeness of the new
51
52
53 206 assembly (Simão, et al. 2015). To estimate the accuracy of the new assembly, all
54
55
56 207 the cleaned short reads from each of the RNA-Seq libraries were mapped back to
57
58
59 208 the genome using hisat2 with default settings (Kim, et al. 2015), and the mapping
60

1
2
3
4 209 rate was estimated.
5

6 210 **The Annotation of Repetitive Elements**

7
8
9 211 Two strategies were recruited to identify repeats residing in the genome. The first
10
11 212 was *de novo* identification based on structural characters, and the second was
12
13
14 213 based on homology to known sequences. On the basis of their structural
15
16
17 214 characteristics, tandem repeats were detected by tandem repeats finder TRF
18
19 215 v4.07b (Benson 1999), and simple sequence repeats (SSRs) residing in the
20
21
22 216 genome were identified using MISA (Beier, et al. 2017) with default settings.
23
24
25 217 Long terminal repeats (LTRs) were identified using LTR_retriever (Ou and Jiang
26
27 218 2018) on the basis of the results of LTRharvest (Ellinghaus, et al. 2008) and
28
29
30 219 LTR_Finder (Ou and Jiang 2019) with the suggested parameters described in the
31
32
33 220 manual. The *de novo* identification of other repeats was implemented using
34
35 221 RepeatModeler followed by genome-scale detection using RepeatMasker v4.0.6
36
37
38 222 (Tarailo-Graovac and Chen 2009) based on homology to the consensus sequences.
39
40
41 223 The consensus sequences deposited in Repbase were also used for the genome-
42
43 224 scale detection of repeat regions residing in the *L. brachycephalus* genome using
44
45 225 RepeatMasker v4.0.6. Possible TE proteins were detected by RepeatProteinMask.
46
47
48 226 All aforementioned results were combined and merged to generate a
49
50
51 227 nonredundant list of repeat elements residing in the genome.

52 53 228 **Gene Prediction and Functional Annotation**

54
55
56 229 Three strategies are typically used to identify gene models residing in eukaryotic
57
58
59 230 genomes: ab initio, homology-based, and transcriptome-based predictions. The ab
60

1
2
3
4 231 initio prediction of genes was performed using Augustus (Stanke, et al. 2006),
5
6 232 GlimmerHMM (Majoros, et al. 2004) and Geneid (Blanco, et al. 2007). The model
7
8
9 233 parameters for these tools were trained based on the genes predicted during
10
11 234 BUSCO v.3.0.1 (Simão, et al. 2015) estimation of the assembly completeness. For
12
13
14 235 the prediction based on homology, proteomes downloaded from public databases
15
16
17 236 for 11 relatives, including *Hypophthalmichthys molitrix*, *Danio rerio*,
18
19 237 *Megalobrama amblycephala*, *Trilophysa bleekeri*, *Triplophysa siluroides*, and
20
21 238 *Triplophysa tibetana*, were used to scan the genome of *L. brachycephalus* using
22
23
24 239 tblastn (Camacho, et al. 2009) with an E-value $\leq 1e-05$, and GeneWise v.2.4.0
25
26
27 240 (Birney, et al. 2004) was used to align the proteins to the homologous genome
28
29
30 241 sequences for accurate spliced alignments. The generated RNA-seq data were also
31
32
33 242 used to obtain direct evidence of the protein coding sequences following the
34
35 243 PASA pipeline (Haas, et al. 2008). The results from the above gene prediction
36
37
38 244 methods were merged by EVM (Haas, et al. 2008) with different weights. After
39
40 245 that, all the genes with transposable elements were removed by TransposonPSI to
41
42
43 246 obtain the final gene set (<http://transposonpsi.sourceforge.net>). The gene set was
44
45 247 annotated by homology scanning to proteins deposited in SwissProt, the
46
47
48 248 nonredundant dataset (NR), and KEGG using Blastp (Camacho, et al. 2009) with
49
50
51 249 an E-value of $1e-05$. Gene Ontology (GO) and protein domain annotations were
52
53
54 250 implemented using InterProScan (Quevillon, et al. 2005) with default settings.

251 **Comparative analyses with relative species**

252 By comparing the predicted protein sequences of *L. brachycephalus* with those of

1
2
3
4 253 the related species, including *Anabarilius grahami*, *Ctenopharyngodon idellus*,
5
6 254 *Danionella translucida*, *D. rerio*, *Labeo rohita*, *M. amblycephala*, *T. bleekeri*, *T.*
7
8
9 255 *siluroides*, and *T. tibetana*, their phylogenetic relationships, divergence times, and
10
11 256 the expansions and contractions of gene families were deduced. First, OrhoMCL
12
13
14 257 v2.0 was used to cluster gene families among these species with default settings
15
16
17 258 (Li, et al. 2003). One-to-one gene families among these species were then selected
18
19 259 and aligned using MAFFT V7 (Katoh, et al. 2002). After that, all the alignments
20
21
22 260 were concatenated to form a supermatrix and were used to deduce their
23
24
25 261 phylogenetic relationships using RAXML7 (Stamatakis 2015). One hundred rapid
26
27 262 bootstraps were also implemented using RAXML7 to estimate the robustness of
28
29
30 263 the phylogeny (Stamatakis, et al. 2008). Based on the phylogeny and the 4-fold
31
32
33 264 degenerate sites extracted from the concatenated supermatrix, their divergence
34
35 265 times were estimated using MCMCTREE included in the PAML4 software
36
37
38 266 package (Yang 2007). Based on the estimated divergence times and the contents
39
40
41 267 for each gene family clustered using OrhoMCL, the possible gene families subject
42
43 268 to expansions and contractions for each species were deduced using CAFÉ2 (De
44
45 269 Bie, et al. 2006).

47 270 **Acknowledgments**

48
49
50 271 We appreciate Pan Xu, general manager of Diggs (Wuhan) Biotechnology Co.,
51
52
53 272 Ltd., for his helpful suggestions on genome assembly. This study was supported
54
55
56 273 by the National Key R & D Program of China (2019YFD0900405), the
57
58 274 Heilongjiang province natural science foundation of China (LH2019C088), the
59
60

1
2
3
4 275 Central Public-interest Scientific Institution Basal Research Fund. HRFRI [NO.
5
6 276 HSY202013PT, HSY201809M, HSY201803M], the Central Public-interest
7
8
9 277 Scientific Institution Basal Research Fund, CAFS under Grant [2020TD56].
10

11 278 **Author Contributions**

12
13
14 279 L.G., H.J. and W.X. conceived the study and collected the samples for sequencing.
15
16 280 M.M. performed the genomics analysis. L.G., H.J. and M.Z. wrote and revised
17
18
19 281 the manuscript. All authors approved the final submission.
20
21

22 282 **Data Availability**

23
24 283 All raw sequencing data, genome assembly and annotations for this article have
25
26
27 284 been deposited at Genome Sequence Archive (GSA, <https://bigd.big.ac.cn/gsa/>)
28
29
30 285 under accession number PRJCA004709.
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

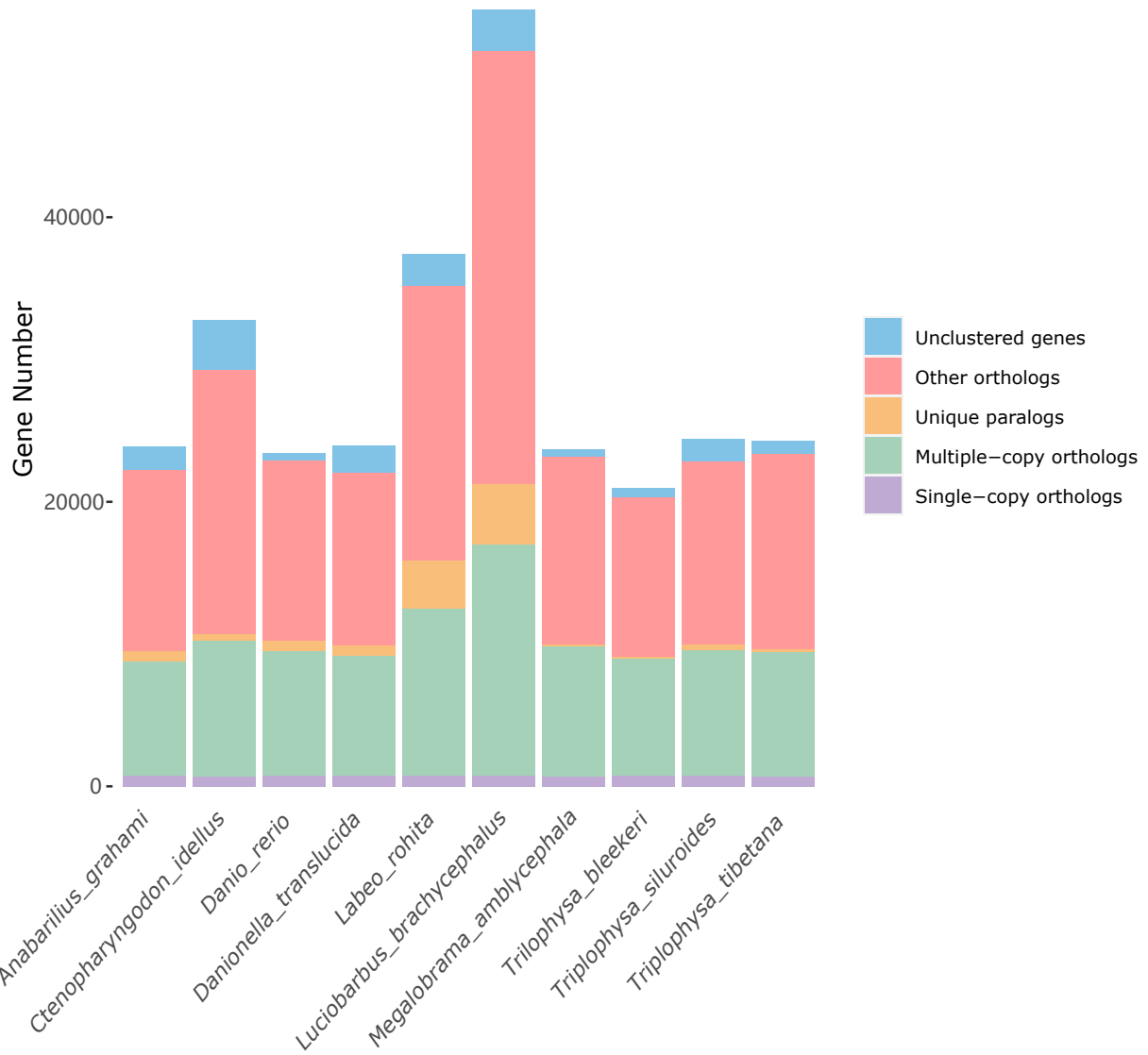
Literature Cited

- 286 **Literature Cited**
- 287 Beier S, Thiel T, Münch T, Scholz U, Mascher M. 2017. MISA-web: a web server for microsatellite
288 prediction. *Bioinformatics* 33: 2583-2585.
- 289 Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids
290 Research* 27: 573-580. doi: 10.1093/nar/27.2.573.
- 291 Birney E, Clamp M, Durbin R. 2004. GeneWise and genomewise. *Genome Research* 14: 988-995.
292 doi: 10.1101/gr.1865504.
- 293 Blanco E, Parra G, Guigó R. 2007. Using geneid to Identify genes. *Current Protocols in
294 Bioinformatics* 18: 4.3.1-4.3.28. doi: 10.1002/0471250953.bi0403s18.
- 295 Camacho C, et al. 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10: 421. doi:
296 10.1186/1471-2105-10-421.
- 297 Chaisson MJ, Tesler G. 2012. Mapping single molecule sequencing reads using basic local
298 alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics*
299 13: 238. doi: 10.1186/1471-2105-13-238.
- 300 Chen S, Zhou Y, Chen Y, Gu J. 2018. fastp: an ultra-fast all-in-one FASTQ preprocessor.
301 *Bioinformatics* 34: i884-i890. doi: 10.1093/bioinformatics/bty560.
- 302 Chin CS, et al. 2016. Phased diploid genome assembly with single-molecule real-time sequencing.
303 *Nat Methods* 13: 1050-1054. doi: 10.1038/nmeth.4035.
- 304 Coad BW. 1998. Systematic biodiversity in the freshwater fishes of Iran. *Canadian Journal of
305 Physiology and Pharmacology* 65: 101-108.
- 306 De Bie T, Cristianini N, Demuth JP, Hahn MW. 2006. CAFE: a computational tool for the study of
307 gene family evolution. *Bioinformatics* 22: 1269-1271. doi: 10.1093/bioinformatics/btl097.
- 308 Ellinghaus D, Kurtz S, Willhoeft U. 2008. LTRharvest, an efficient and flexible software for de novo
309 detection of LTR retrotransposons. *BMC Bioinformatics* 9: 18. doi: 10.1186/1471-2105-9-18.
- 310 Esmaeili HR, Mehraban H, Abbasi K, Keivany Y, Coad BW. 2017. Review and updated checklist of
311 freshwater fishes of Iran: taxonomy, distribution and conservation status. *Iranian Journal of
312 Ichthyology* 4: 1-114.
- 313 Geng L, Tong G, Jiang H, Wei X. 2016. Effect of salinity and alkalinity on *Luciobarbus capito* gill
314 Na⁺/K⁺-ATPase enzyme activity, plasma ion concentration, and osmotic pressure. *BioMed
315 Research International* 2016: 1-7.
- 316 Geng LW, Xu W, Jiang HF, Tong GX. 2013. Karyotype analysis of *Barbus capito* (Güldenstädt, 1773)
317 using curve measurement software. *Journal of Applied Ichthyology* 29: 922-924.
- 318 Haas BJ, et al. 2008. Automated eukaryotic gene structure annotation using EVidenceModeler and
319 the program to assemble spliced alignments. *Genome Biology* 9: R7. doi: 10.1186/gb-2008-
320 9-1-r7.
- 321 Howe K, et al. 2013. The zebrafish reference genome sequence and its relationship to the human
322 genome. *Nature* 496: 498-503. doi: 10.1038/nature12111.
- 323 Jiang H, Geng L, Yang J, Tong G, Xu W. 2019. The complete mitochondrial genome of the aral
324 barbel *Luciobarbus brachycephalus* (Cypriniformes: Cyprinidae: Barbinae). *Mitochondrial
325 DNA Part B* 4: 3685-3686.
- 326 Katoh K, Misawa K, Kuma Ki, Miyata T. 2002. MAFFT: a novel method for rapid multiple sequence
327 alignment based on fast Fourier transform. *Nucleic Acids Research* 30: 3059-3066. doi:
328 10.1093/nar/gkf436.
- 329 Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory requirements.

- 1
2
3 330 Nat Methods 12: 357-360. doi: 10.1038/nmeth.3317.
- 4 331 Li A, et al. 2019. Morphological characteristics and parameters measurement of *Barbus capito*.
5 332 Chinese Agricultural Science Bulletin. p. 144-152.
- 6 333 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv
7 334 1303.3997.
- 8 335 Li L, Stoeckert CJ, Roos DS. 2003. OrthoMCL: Identification of ortholog groups for eukaryotic
9 336 genomes. Genome Research 13: 2178-2189. doi: 10.1101/gr.1224503.
- 10 337 Longwu G, Cuiyun L, Guangxiang T, Chao L, Wei X. 2012. Development and characterization of
11 338 twenty microsatellite markers for the endangered fish *Luciobarbus capito*. Conservation
12 339 Genetics Resources 4: 865-867.
- 13 340 Longwu G, et al. 2010. Technique of artificial reproduction of *Barbus capito*. Journal of Jilin
14 341 Agricultural University. p. 218-220.
- 15 342 Majoros WH, Pertea M, Salzberg SL. 2004. TigrScan and GlimmerHMM: two open source ab initio
16 343 eukaryotic gene-finders. Bioinformatics 20: 2878-2879. doi: 10.1093/bioinformatics/bth315.
- 17 344 Marcais G, Kingsford C. 2011. A fast, lock-free approach for efficient parallel counting of
18 345 occurrences of k-mers. Bioinformatics 27: 764-770. doi: 10.1093/bioinformatics/btr011.
- 19 346 Ou S, Jiang N. 2019. LTR_FINDER_parallel: parallelization of LTR_FINDER enabling rapid
20 347 identification of long terminal repeat retrotransposons. Mob DNA 10: 48. doi:
21 348 10.1186/s13100-019-0193-0.
- 22 349 Ou S, Jiang N. 2018. LTR_retriever: a Highly accurate and sensitive program for identification of
23 350 long terminal repeat retrotransposons. Plant Physiol 176: 1410-1422. doi:
24 351 10.1104/pp.17.01310.
- 25 352 Quevillon E, et al. 2005. InterProScan: protein domains identifier. Nucleic Acids Research 33:
26 353 W116-W120.
- 27 354 Ranallo-Benavidez TR, Jaron KS, Schatz MC. 2020. GenomeScope 2.0 and Smudgeplot for
28 355 reference-free profiling of polyploid genomes. Nat Commun 11: 1432. doi: 10.1038/s41467-
29 356 020-14998-3.
- 30 357 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO: assessing
31 358 genome assembly and annotation completeness with single-copy orthologs. Bioinformatics
32 359 31: 3210-3212. doi: 10.1093/bioinformatics/btv351.
- 33 360 Stamatakis A. 2015. Using RAxML to infer phylogenies. Current Protocols in Bioinformatics 51:
34 361 6.14.11-16.14.14. doi: 10.1002/0471250953.bi0614s51.
- 35 362 Stamatakis A, Hoover P, Rougemont J. 2008. A rapid bootstrap algorithm for the RAxML Web
36 363 servers. Syst Biol 57: 758-771. doi: 10.1080/10635150802429642.
- 37 364 Stanke M, et al. 2006. AUGUSTUS: ab initio prediction of alternative transcripts. Nucleic Acids
38 365 Research 34: W435-W439. doi: 10.1093/nar/gkl200.
- 39 366 Tarailo-Graovac M, Chen N. 2009. Using repeatMasker to identify repetitive elements in genomic
40 367 sequences. Current Protocols in Bioinformatics 25: 4.10.11-14.10.14. doi:
41 368 10.1002/0471250953.bi0410s25.
- 42 369 Voltaire E, Brunet F, Naville M, Volff J-N, Galiana D. 2017. Expansion by whole genome duplication
43 370 and evolution of the sox gene family in teleost fish. PloS one 12: e0180936-e0180936. doi:
44 371 10.1371/journal.pone.0180936.
- 45 372 Walker BJ, et al. 2014. Pilon: an integrated tool for comprehensive microbial variant detection and
46 373 genome assembly improvement. PLoS One 9: e112963. doi: 10.1371/journal.pone.0112963.
- 47 374 Wang Y, et al. 2015. The draft genome of the grass carp (*Ctenopharyngodon idellus*) provides

- 1
2
3 375 insights into its evolution and vegetarian adaptation. *Nature Genetics* 47: 625-631. doi:
4 376 10.1038/ng.3280.
5
6 377 Xiao S, et al. 2020. Genome of tetraploid fish *Schizothorax o'connori* provides insights into early
7 378 re-diploidization and high-altitude adaptation. *iScience* 23: 101497. doi:
8 379 10.1016/j.isci.2020.101497.
9
10 380 Xu P, et al. 2019. The allotetraploid origin and asymmetrical genome evolution of the common
11 381 carp *Cyprinus carpio*. *Nat Commun* 10: 4625. doi: 10.1038/s41467-019-12644-1.
12 382 Xu P, et al. 2014. Genome sequence and genetic diversity of the common carp, *Cyprinus carpio*.
13 383 *Nature Genetics* 46: 1212-1219. doi: 10.1038/ng.3098.
14
15 384 Yang Z 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and*
16 385 *Evolution* 24: 1586-1591. doi: 10.1093/molbev/msm088.
17
18 386
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

A



B

