

GigaScience

The draft genome of blunt snout bream (*Megalobrama amblycephala*) reveals the development of intermuscular bone and adaptation to herbivorous diet --Manuscript Draft--

| | | |
|--|--|----------------|
| Manuscript Number: | GIGA-D-16-00088R3 | |
| Full Title: | The draft genome of blunt snout bream (<i>Megalobrama amblycephala</i>) reveals the development of intermuscular bone and adaptation to herbivorous diet | |
| Article Type: | Research | |
| Funding Information: | Fundament Research Funds for the Central Universities (2662015PY019) | Not applicable |
| | Modern Agriculture Industry Technology System Construction Projects of China (CARS-46-05) | Not applicable |
| | International Scientific and Technology Cooperation Program of Wuhan City (2015030809020365) | Not applicable |
| Abstract: | <p>Background: The blunt snout bream, <i>Megalobrama amblycephala</i>, is the economically most important cyprinid fish species. As an herbivore, it can be grown by eco-friendly and resource-conserving aquaculture. However, the large number of intermuscular bones in the trunk musculature is adverse to fish meat processing and consumption.</p> <p>Results: As a first towards optimizing this aquatic livestock, we present a 1.116-Gb draft genome of <i>M. amblycephala</i>, with 779.54 Mb anchored on 24 linkage groups. Integrating spatiotemporal transcriptome analyses we show that intermuscular bone is formed in the more basal teleosts by intramembranous ossification, and may be involved in muscle contractibility and coordinating cellular events. Comparative analysis revealed that olfactory receptor genes, especially of the beta type, underwent an extensive expansion in herbivorous cyprinids, whereas the gene for the umami receptor T1R1 was specifically lost in <i>M. amblycephala</i>. The composition of gut microflora, which contributes to the herbivorous adaptation of <i>M. amblycephala</i>, was found to be similar to that of other herbivores.</p> <p>Conclusions: As a valuable resource for improvement of <i>M. amblycephala</i> livestock, the draft genome sequence offers new insights into the development of intermuscular bone and herbivorous adaptation.</p> | |
| Corresponding Author: | Weimin Wang CHINA | |
| Corresponding Author Secondary Information: | | |
| Corresponding Author's Institution: | | |
| Corresponding Author's Secondary Institution: | | |
| First Author: | Han Liu | |
| First Author Secondary Information: | | |
| Order of Authors: | Han Liu | |
| | Chunhai Chen | |
| | Zexia Gao | |
| | Jiumeng Min | |
| | Yongming Gu | |
| | Jianbo Jian | |
| | Xiewu Jiang | |

| | |
|--|--|
| | Huimin Cai |
| | Ingo Ebersberger |
| | Meng Xu |
| | Xinhui Zhang |
| | Jianwei Chen |
| | Wei Luo |
| | Boxiang Chen |
| | Junhui Chen |
| | Hong Liu |
| | Jiang Li |
| | Ruifang Lai |
| | Mingzhou Bai |
| | Jin Wei |
| | Shaokui Yi |
| | Huanling Wang |
| | Xiaojuan Cao |
| | Xiaoyun Zhou |
| | Yuhua Zhao |
| | Kaijian Wei |
| | Ruibin Yang |
| | Bingnan Liu |
| | Shancen Zhao |
| | Xiaodong Fang |
| | Manfred Scharl |
| | Xueqiao Qian |
| | Weimin Wang |
| Order of Authors Secondary Information: | |
| Response to Reviewers: | <p>2 May 2017 Dr. Hans Zauner Journal: GigaScience</p> <p>Dear Dr. Zauner, Manuscript No.: GIGA-D-16-00088R2 Title: "The draft genome of blunt snout bream (<i>Megalobrama amblycephala</i>) reveals the development of intermuscular bone and adaptation to herbivorous diet" Author(s): Han Liu, Chunhai Chen, Zexia Gao, Jiumeng Min, Yongming Gu, Jianbo Jian, Xiewu Jiang, Huimin Cai, Ingo Ebersberger, Meng Xu, Xinhui Zhang, Jianwei Chen, Wei Luo, Boxiang Chen, Junhui Chen, Hong Liu, Jiang Li, Ruifang Lai, Mingzhou Bai, Jin Wei, Shaokui Yi, Huanling Wang, Xiaojuan Cao, Xiaoyun Zhou, Yuhua Zhao, Kaijian Wei, Ruibin Yang, Bingnan Liu, Shancen Zhao, Xiaodong Fang, Manfred Scharl, Xueqiao Qian, Weimin Wang</p> <p>We have carefully read the referee's comments which you forwarded to us of 27 April 2017. We would like to express our sincere thanks for the positive comments. We have addressed all the suggestions. The common name of the species has been included in the title (Line 1 to 2). The Fishbase number and image of an adult blunt snout bream have also been added in the revised manuscript (Line 73 to 74). The amendments are highlighted in red in the revised manuscript. Responses to the reviewer's comments</p> |

are detailed below in this letter. We hope that with the amendments made in response to you and the reviewer's comments, the manuscript is now acceptable for publication in GigaScience.

I look forward to hearing from you soon.

Yours sincerely,
Weimin Wang (PhD) (Correspondence author)
College of Fisheries
Huazhong Agricultural University
Wuhan 430070, P. R. China
E-mail address: wangwm@mail.hzau.edu.cn
Tel: +86-27-8728 4292; Fax: +86-27-8728 4292

Response to Reviewer

Reviewer Report

Reviewer #2:

1. Have a look at the sentence at line 104-107: "to assess the genome assembly quality"

Author response: This sentence has been modified as "To assess the quality of genome assembly, the short-insert size paired-end libraries reads and published ESTs [14] (Additional file 1: Tables S3 and S4) were mapped onto the genome. The results indicated that the assembled error is low." (Line 105 to 107)

2. Line 151 "single -copy genes"?

Author response: This expression has been changed as "single-copy orthologous genes". (Line 152)

3. 153: "outgroup" not "out group"

Author response: This has been corrected.

4. Rephrase line 552.

Author response: The SRP accession number has been added in the revised manuscript. This sentence has been rephrased as "Raw whole genome sequencing and RAD-Seq data have been deposited at NCBI in the SRA under accession number SRP090157 (BioProject Number: PRJNA343584)". (Line 553 to 554)

5. 3 Additional file 2 Data note 1. Typo: "Expansion" instead of "Expasion"

Author response: This has been corrected.

6. Figure 2B not mentioned in the text.

Author response: We have now mentioned it in the text and expressed as "We found 9349 orthologous gene families shared among five fish species. 246 are specific in the M. Amblycephala (Figure 3B)". (Line 155 to 156)

Additional Information:

Question

Response

Are you submitting this manuscript to a special series or article collection?

No

Experimental design and statistics

Yes

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.

Have you included all the information requested in your manuscript?

| | |
|---|------------|
| <p>Resources</p> <p>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.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p> | <p>Yes</p> |
| <p>Availability of data and materials</p> <p>All datasets and code on which the 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.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p> | <p>Yes</p> |

1 **The draft genome of blunt snout bream (*Megalobrama amblycephala*) reveals the**
2 **development of intermuscular bone and adaptation to herbivorous diet**

3 Han Liu^{1†}, Chunhai Chen^{2†}, Zexia Gao^{1†}, Jiumeng Min^{2†}, Yongming Gu^{3†}, Jianbo Jian^{2†}, Xiewu
4 Jiang³, Huimin Cai², Ingo Ebersberger⁴, Meng Xu², Xinhui Zhang¹, Jianwei Chen², Wei Luo¹,
5 Boxiang Chen^{1,3}, Junhui Chen², Hong Liu¹, Jiang Li², Ruifang Lai¹, Mingzhou Bai², Jin Wei¹,
6 Shaokui Yi¹, Huanling Wang¹, Xiaojuan Cao¹, Xiaoyun Zhou¹, Yuhua Zhao¹, Kaijian Wei¹,
7 Ruibin Yang¹, Bingnan Liu³, Shancen Zhao², Xiaodong Fang², Manfred Schartl^{5,*}, Xueqiao
8 Qian^{3,*}, Weimin Wang^{1,*}

9
10 *Equally contributing corresponding authors: wangwm@mail.hzau.edu.cn; xueqiaoqian@263.net;
11 phch1@biozentrum.uni-wuerzburg.de

12 †Equal contributors

13 ¹College of Fisheries, Key Lab of Freshwater Animal Breeding, Ministry of Agriculture, Key Lab
14 of Agricultural Animal Genetics, Breeding and Reproduction of Ministry of Education, Huazhong
15 Agricultural University, Wuhan 430070, China

16 ²Beijing Genomics Institute (BGI)-Shenzhen, Shenzhen 518083, China

17 ³Guangdong Haid Group Co., Ltd., Guangzhou 511400, China

18 ⁴Dept. for Applied Bioinformatics, Institute for Cell Biology and Neuroscience, Goethe University,
19 Frankfurt D-60438, Germany

20 ⁵Physiological Chemistry, University of Würzburg, Biozentrum, Am Hubland, and Comprehensive
21 Cancer Center Mainfranken, University Clinic Würzburg, Würzburg 97070, Germany and Texas
22 A&M Institute for Advanced Study and Department of Biology, Texas A&M University, College
23 Station, TX 77843, USA

29 **Abstract**

30 **Background:** The blunt snout bream, *Megalobrama amblycephala*, is the economically most
31 important cyprinid fish species. As an herbivore, it can be grown by eco-friendly and
32 resource-conserving aquaculture. However, the large number of intermuscular bones in the trunk
33 musculature is adverse to fish meat processing and consumption.

34 **Results:** As a first towards optimizing this aquatic livestock, we present a 1.116-Gb draft genome
35 of *M. amblycephala*, with 779.54 Mb anchored on 24 linkage groups. Integrating spatiotemporal
36 transcriptome analyses we show that intermuscular bone is formed in the more basal teleosts by
37 intramembranous ossification, and may be involved in muscle contractibility and coordinating
38 cellular events. Comparative analysis revealed that olfactory receptor genes, especially of the beta
39 type, underwent an extensive expansion in herbivorous cyprinids, whereas the gene for the umami
40 receptor *TIR1* was specifically lost in *M. amblycephala*. The composition of gut microflora, which
41 contributes to the herbivorous adaptation of *M. amblycephala*, was found to be similar to that of
42 other herbivores.

43 **Conclusions:** As a valuable resource for improvement of *M. amblycephala* livestock, the draft
44 genome sequence offers new insights into the development of intermuscular bone and herbivorous
45 adaptation.

46
47 **Keywords:** *Megalobrama amblycephala*, whole genome, herbivorous diet, intermuscular bone,
48 transcriptome, gut microflora

58 **Background**

1
2 59 Fishery and aquaculture play an important role in global alimentation. Over the past decades food
3
4 60 fish supply has been increasing with an annual rate of 3.6 percent, about 2 times faster than the
5
6 61 human population [1]. This growth of fish production is meanwhile solely accomplished by an
7
8 62 extension of aquaculture, as over the past thirty years the total mass of captured fish has remained
9
10 63 almost constant [1]. As a consequence of this emphasis on fish breeding, the genomes of various
11
12 64 economically important fish species, e.g. Atlantic cod (*Gadus morhua*) [2], rainbow trout
13
14 65 (*Oncorhynchus mykiss*) [3], European sea bass (*Dicentrarchus labrax*) [4], yellow croaker
15
16 66 (*Larimichthys crocea*) [5], half-smooth tongue sole (*Cynoglossus semilaevis*) [6], tilapia
17
18 67 (*Oreochromis niloticus*) [7] and channel catfish (*Ictalurus punctatus*) [8] have been sequenced.
19
20 68 Yet, the majority of these species are carnivorous, requiring large inputs of protein from wild
21
22 69 caught fish or other precious feed. Reports on draft genomes of herbivorous and omnivorous
23
24 70 species, in particular cyprinid fish are scarce. It is well known that cyprinids are currently the
25
26 71 economically most important group of teleosts for sustainable aquaculture. They grow to large
27
28 72 population sizes in the wild and already now account for the majority of freshwater aquaculture
29
30 73 production worldwide [1]. Among these, the herbivorous blunt snout bream, *Megalobrama*
31
32 74 *amblycephala* (Fishbase Sp. ID: 285) (Figure 1), a particularly eco-friendly and
33
34 75 resource-conserving species, is predominant in aquaculture and has been greatly developed in
35
36 76 China (Additional file 1: Figure S1) [1]. However, most cyprinids, including *M. amblycephala*,
37
38 77 have a large number of intermuscular bones (IBs) in the trunk musculature, which have an adverse
39
40 78 effect on fish meat processing and consumption. IBs—a unique form of bone occurring only in the
41
42 79 more basal teleosts—are completely embedded within the myosepta and are not connected to the
43
44 80 vertebral column or any other bones [9, 10]. Our previous study on IB development of *M.*
45
46 81 *amblycephala* revealed that some miRNA-mRNA interaction pairs may be involved in regulating
47
48 82 bone development and differentiation [11]. However, the molecular genetic basis and the
49
50 83 evolution of this unique structure are still unclear. Unfortunately, the recent sequencing of two
51
52 84 cyprinid genomes common carp (*Cyprinus carpio*) [12] and grass carp (*Ctenopharyngodon idellus*)
53
54 85 [13], which provided valuable information for their genetic breeding, contributed little to the
55
56 86 understanding of IB formation.
57
58
59
60
61
62
63
64
65

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

87 In an initial genome survey of *M. amblycephala*, we identified 25,697 single-nucleotide
88 polymorphism (SNP) [14], 347 conserved miRNAs [15], and 1,136 miRNA-mRNA interaction
89 pairs [11]. However, lack of a whole genome sequence resource limited a thorough investigation
90 of *M. amblycephala*. Here we report the first high-quality draft genome sequence of *M.*
91 *amblycephala*. Integrating this novel genome resource with tissue- and developmental
92 stage-specific gene expression information, as well as with meta-genome data to investigate the
93 composition of the gut microbiome (Workflow shown in additional file 1: Figure S2) provides
94 relevant insights into the function and evolution of two key features characterizing this species:
95 The formation of IB and the adaptation to herbivory. By that our study lays the foundation for
96 genetically optimizing *M. amblycephala* to further increase its relevance for securing human food
97 supply.

98 99 **Data description**

100 **Genome Assembly and Annotation**

101 The *M. amblycephala* genome was sequenced and assembled by a whole-genome shotgun strategy
102 using genomic DNA from a double-haploid fish (Additional file 1: Table S1). We assembled a
103 1.116 Gb reference genome sequence from 142.55 Gb (approximately 130-fold coverage) of clean
104 data [16] (Additional file 1: Tables S1 and S2, Figure S3). The contig and scaffold N50 lengths
105 reached 49 Kb and 839 Kb, respectively (Table 1). The largest scaffold spans 8,951 Kb and the
106 4,034 largest scaffolds cover 90% of the assembly. To assess the quality of genome assembly, the
107 short-insert size paired-end libraries reads and published ESTs [14] (Additional file 1: Tables S3
108 and S4) were mapped onto the genome. The results indicated that the assembled error is low. To
109 further estimate the completeness of the assembly and gene prediction, the benchmarking
110 universal single-copy orthologs (BUSCO) (BUSCO , RRID:SCR_015008) [17] analysis was used
111 and the results showed that the assembly contains 81.4% complete and 9.1% partial vertebrate
112 BUSCO orthologues (Additional file 1: Table S5).

113 The *M. amblycephala* genome has an average GC content of 37.3%, similar to cyprinid *C.*
114 *carpio* and *Danio rerio* (Additional file 1: Figure S4). Using a comprehensive annotation strategy
115 combining RNA-seq derived transcript evidence, *de-novo* gene prediction and sequence similarity

116 to proteins from five further fish species, we annotated a total of 23,696 protein-coding genes
117 (Additional file 1: Table S6). Of the predicted genes, 99.44% (23,563 genes) are annotated by
118 functional database. In addition, we identified 1,796 non-coding RNAs including 474 miRNAs,
119 220 rRNA, 530 tRNAs, and 572 snRNAs. Transposable elements (TEs) comprise approximately
120 34.18% (381.3 Mb) of the *M. amblycephala* genome (Additional file 1: Table S7). DNA
121 transposons (23.80%) and long terminal repeat retrotransposons (LTRs) (9.89%) are the most
122 abundant TEs in *M. amblycephala*. The proportion of LTRs in *M. amblycephala* is highest in
123 comparison with other teleosts: *G. morhua* (4.88%) [2], *L. crocea* (2.2%) [5], *C. semilaevis*
124 (0.08%) [6], *C. carpio* (2.28%) [12], *C. idellus* (2.58%) [13] and stickleback (*Gasterosteus*
125 *aculeatus*) (1.9%) [18] (Additional file 1: Tables S7 and S8, Figure S5). The distribution of
126 divergence between the TEs in *M. amblycephala* peaks at 7% (Additional file 1: Figure S6),
127 indicating a more recent activity of these TEs when compared with *O. mykiss* (13%) [3] and *C.*
128 *semilaevis* (9%) [6].

130 **Anchoring Scaffolds and Shared Synteny Analysis**

131 Sequencing data from 198 F1 specimens, including the parents, were used as the mapping
132 population to anchor the scaffolds on to 24 pseudo-chromosomes of the *M. amblycephala* genome.
133 Following RAD-Seq and sequencing protocol, 1883.5 Mb of 125-bp reads (on average 30.6 Mb
134 and 9.3 Mb of read data for each parent and progeny, respectively) were generated on the HiSeq
135 2500 next-generation sequencing platform. Based on the SOAP bioinformatic pipeline
136 (SOAPdenovo2 , RRID:SCR_014986), we generated 5,317 SNP markers for constructing a
137 high-resolution genetic map. The map spans 1,701 cM with a mean marker distance of 0.33 cM
138 and facilitated an anchoring of 1,434 scaffolds comprising 70% (779.54 Mb) of the *M.*
139 *amblycephala* genome assembly to form 24 linkage groups (LG) (Additional file 1: Table S9). Of
140 the anchored scaffolds, 598 could additionally be oriented (678.27 Mb, 87.01% of the total
141 anchored sequences) (Figure 2A). A subsequent comparison of the gene order between *M.*
142 *amblycephala* and its close relative *C.idellus* revealed 607 large shared syntenic blocks
143 encompassing 11,259 genes, and 190 chromosomal rearrangements. The values change to 1,062
144 regions, 13,152 genes and 279 rearrangements when considering *D. rerio*. The unexpected higher

145 number of genes in syntenic regions shared with the more distantly related *D. rerio* is most likely
146 an effect of the more complete genome assembly of this species compared to *C. idellus*. The
147 rearrangement events are distributed across all *M. amblycephala* linkage groups without evidence
148 for a local clustering (Figure 2B). The most prominent event is a chromosomal fusion in *M.*
149 *amblycephala* LG02 that joined two *D. rerio* chromosomes, Dre10 and Dre22. The same fusion is
150 observed in *C. idellus* but not in *C. carpio* suggesting that it probably occurred in a last common
151 ancestor of *M. amblycephala* and *C. idellus*, approximately 13.1 million years ago (Additional file
152 1: Figure S7).

153

154 **Results**

155 **Evolutionary Analysis**

156 A phylogenetic analysis of 316 single-copy orthologous genes in the genomes of 10 other fish
157 species, and coelacanth (*Latimeria chalumnae*) and elephant shark (*Callorhynchus milii*), as
158 outgroup served as a basis for investigating the evolutionary trajectory of *M. amblycephala*
159 (Figure 3A, Additional file 1: Figure S8). We found 9349 orthologous gene families shared among
160 five fish species. 246 are specific in the *M. Amblycephala* (Figure 3B). To illuminate the
161 evolutionary process resulting in the adaptation to a grass diet, we analyzed the functional
162 categories of expanded genes in the *M. amblycephala* and *C. idellus* lineage (Additional file 1:
163 Figure S9, Additional file 2: Data Note1), two typical herbivores mainly feeding on aquatic and
164 terrestrial grasses. Among the significantly over-represented KEGG pathways (KEGG ,
165 RRID:SCR_012773) (Fisher's exact test, $P < 0.01$), we find olfactory transduction (ko04740),
166 immune-related pathways (ko04090, ko04672, ko04612 and ko04621), lipid metabolic related
167 process (ko00590, ko03320, ko00591, ko00565, ko00592 and ko04975), as well as xenobiotics
168 biodegradation and metabolism (ko00625 and ko00363) (Figure S10). Indeed, when tracing
169 positively selected genes (PSG) in *M. amblycephala* and *C. idellus* (Additional file 3: Date Note2),
170 we identified 10 candidates involved in starch and sucrose metabolism (ko00500), in citrate cycle
171 (ko00020) and in other types of O-glycan biosynthesis (ko00514). Moreover, 10 genes encoding
172 enzymes involved in lipid metabolism appear positively selected in both fish species (Additional
173 file 1: Table S10).

174 **Development of Intermuscular Bones**

1
2
3 175 To explain the genetic basis of IB, their formation and their function in cyprinids, we first
4
5 176 analyzed the functional annotation of genes that expanded in this lineage (Figure 3C). Many of
6
7 177 these genes are involved in cell adhesion (GO: 0007155, $P=5.26E-32$, 357 genes), myosin
8
9 178 complex (GO:0016459, $P=2.74E-08$, 100 genes) and cell-matrix adhesion (GO:0007160,
10
11 179 $P=1.59E-21$, 69 genes) (Figure 3C). As a second line of evidence, we performed transcriptome
12
13 180 analyses of early developmental stages (stage1: whole larvae without IBs) and juvenile *M.*
14
15 181 *amblycephala* (stage2: trunk muscle with partial IBs; stage3: trunk muscle with completed IBs)
16
17 182 (Figure 4A). Compared with stage1, 388 and 651 differentially expressed genes (DEGs) are
18
19 183 up-regulated in stage2 and stage3, respectively. And 249 of them are significantly up-regulated
20
21 184 both in stage2 and stage3. KEGG analyses indicate many of these genes involved in tight junction
22
23 185 (ko04530), regulation of actin cytoskeleton (ko04810), cardiac muscle contraction (ko04260) and
24
25 186 vascular smooth muscle contraction (ko04270) (Additional file 1: Figure S11). Specifically, 26
26
27 187 genes encoding proteins related to muscle contraction, including titin, troponin, myosin, actinin,
28
29 188 calmodulin and other Ca^{2+} transporting ATPases (Figure 4A) point to a strong remodeling of the
30
31 189 musculature compartment. To confirm that the observed differences in gene expression are indeed
32
33 190 linked to IB formation and function and are not simply due to the fact that different developmental
34
35 191 stages were compared, we performed differential expression analysis of muscle tissues, IB, and
36
37 192 connective tissues from the same six months old individual of *M. amblycephala* (Figure 4B,
38
39 193 Additional file 1: Figure S12). 1,290 DEGs and 5,231 DEGs are significantly up-regulated in IB
40
41 194 compared with connective tissues and muscle, respectively. 24 of these DEGs encode extracellular
42
43 195 matrix (ECM) proteins (collagens and integrin-binding protein), Rho GTPase family (*RhoA*, *Rho*
44
45 196 *GAP*, *Rac*, *Ras*), motor proteins (myosin, dynein, actin), and calcium channel regulation proteins
46
47 197 (Additional file 1: Figure S13 and Table S11). In addition, GO annotations of 963 IB-specific
48
49 198 genes indicative of abundance in protein binding (GO:0005515), calcium ion binding
50
51 199 (GO:0005509), GTP binding (GO:0005525) and iron ion binding (GO:0005506) were found
52
53 200 (Figure 4C).

54
55
56
57
58 201 During development of *M. amblycephala*, the first IB appears in muscles of caudal vertebrae
59
60 202 as early as 28 days post fertilization (dpf) when body length is 12.95 mm (Additional file 1: Figure

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

203 S14). The system then develops and ossifies predominantly from posterior to anterior (Additional
204 file 1: Figure S15). IBs are present throughout the body within two months (Additional file 1:
205 Figure S16) and develop into multiple morphological types in adults (Additional file 1: Figure
206 S17). The bone is formed directly without an intermediate cartilaginous stage (Additional file 1:
207 Figures S18 and S19). We also found a large number of mature osteoblasts distributed at the edge
208 of the bone matrix while some osteocytes were apparent in the center of the mineralized bone
209 matrix (Additional file 1: Figures S20 and S21). These primary bone-forming cells predominantly
210 regulate bone formation and function throughout life. Notably, among the genes up-regulated in
211 IB, 35 bone formation regulatory genes were identified (shown in colored boxes in Figure 4D). In
212 particular, genes involved in bone morphogenetic protein (BMP) signaling including *Bmp3*,
213 *Smad8*, *Smad9*, and *Id2*, in fibroblast growth factor (FGF) signaling including *Fgf2*, *Fgfr1a*,
214 *Fgfbp2*, *Col6a3*, and *Col4a5*, and in Ca²⁺ channels including *Cacna1c*, *CaM*, *Creb5* and *Nfatc*
215 were highly expressed (>2-fold change) in IB (Additional file 1: Figure S22).

216

217 **Adaptation to Herbivorous Diet**

218 Next to the presence of IB, the herbivory of *M. amblycephala* is the second key feature in
219 connection to the use of this species as aquatic livestock. Olfaction, the sense of smell, is crucial
220 for animals to find food. The perception of smell is mediated by a large gene family of olfactory
221 receptor (OR) genes. In the *M. amblycephala* genome, we identified 179 functional olfactory
222 receptor (OR) genes (Figure 5A), and based on the classification of Niimura [19], 158, 117 and
223 153 receptors for water-borne odorants were identified in *M. amblycephala*, *C. idellus* and *D.*
224 *rerio*, respectively (Additional file 1: Table S12). Overall, these receptor repertoires are
225 substantially larger than those of other and carnivorous teleosts (*G. morhua*, *C. semilaevis*, *O.*
226 *latipes*, *X. maculatus*) (Additional file 1: Figures S23 and S24, Table S12). In addition, we found
227 a massive expansion of beta-type OR genes in the genomes of the herbivorous *M. amblycephala*
228 and *C. idellus*, while very few exist in other teleosts (Figure 5B, Additional file 1: Table S12).

229 Taste is also an important factor in the development of dietary habits. Most animals can
230 perceive five basic tastes, namely sourness, sweetness, bitterness, saltiness and umami [20]. *TIR1*,

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

231 the receptor gene necessary for sensing umami, has been lost in herbivorous *M. amblycephala* but
232 is duplicated in carnivorous *G. morhua*, *C. semilaevis* and omnivorous *O. latipes* and *X. maculatus*
233 (Figures 5C and 5D, Additional file 1: Figures S25-26 and Table S13). In contrast, *T1R2*, the
234 receptor gene for sensing sweet, has been duplicated in herbivorous *M. amblycephala* and *C.*
235 *idellus*, omnivorous *C. carpio* and *D. rerio*, while it has been lost in carnivorous *G. morhua* and *C.*
236 *semilaevis* (Additional file 1: Figure S27 and Table S13). Also the *T2R* gene family, most likely
237 important in the course switching to a diet that contains a larger fraction of bitterness containing
238 food, has been expanded in *M. amblycephala*, *C. idellus*, *C. carpio* and *D. rerio* (Additional file 1:
239 Figure S28).

240 To obtain further insights into the genetic adaptation to herbivorous diet, we focused on
241 further genes that might be associated with digestion. Genes that encode proteases (including
242 pepsin, trypsin, cathepsin and chymotrypsin) and amylases (including alpha-amylase and
243 glucoamylase) were identified in the genomes of *M. amblycephala*, carnivorous *C. semilaevis*, *G.*
244 *morhua* and omnivorous *D. rerio*, *O. latipes* and *X. maculatus*, indicating that herbivorous *M.*
245 *amblycephala* has a protease repertoire that is not substantially different from those of carnivorous
246 and omnivorous fishes (Additional file 1: Table S14). We did not identify any genes encoding
247 potentially cellulose-degrading enzymes including endoglucanase, exoglucanase and
248 beta-glucosidase in the genome of *M. amblycephala*, suggesting that utilization of the herbivorous
249 diet may largely depend on the gut microbiome. To elucidate this further, we determined the
250 composition of the gut microbial communities of juvenile, domestic, wild adult *M. amblycephala*
251 and wild adult *C. idellus* using bacterial 16S rRNA sequencing. A total of 549,020 filtered high
252 quality sequence reads from 12 samples were clustered at a similarity level of 97%. The resulting
253 8,558 operational taxonomic units (OTUs) are dominated at phylum level by Proteobacteria,
254 Fusobacteria, Bacteroidetes, Firmicutes and Actinobacteria (Additional file 1: Table S15, Figure
255 5E). Increasing the resolution to the genus level, the composition and relative abundance of the
256 gut microbiota of wild adult *M. amblycephala* and *C. idellus* are still very similar (Additional file
257 1: Table S16) and we could identify more than 7% cellulose-degrading bacteria (Additional file 1:
258 Table S17).

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

260 Discussion

261 The evolutionary trajectory analysis of *M. amblycephala* and other teleosts revealed that *M.*
262 *amblycephala* has the closest relationship to *C. idellus*. Both the species are herbivorous fish but
263 which endogenous and exogenous factors affected their feeding habits and how they adapted to
264 their herbivorous diet is not known. Our results from the expanded genes and PSG in the lineage
265 of the two herbivores uncovered a number of genes that are involved in glucose, lipid and
266 xenobiotics metabolism, which would enhance the ability of an herbivore to detoxify the
267 secondary compounds present in grasses that are adverse or even toxic to the organism.
268 Furthermore, the high-fiber but low-energy grass diet requires a highly effective intermediate
269 metabolism that accelerates carbohydrate and lipid catabolism and conversion into energy to
270 maintain physiological functions.

271 Olfaction and taste are also crucial for animals to find food and to distinguish whether
272 potential food is edible or harmful [21, 22]. The ORs of teleosts are predominantly expressed in
273 the main olfactory epithelium of the nasal cavity [21, 23] and can discriminate, like those of other
274 vertebrates, different kinds of odor molecules. Previous studies have demonstrated that the beta
275 type OR genes are present in both aquatic and terrestrial vertebrates, indicating that the
276 corresponding receptors detect both water-soluble and airborne odorants [19, 21]. In the present
277 study, the search for genes encoding OR showed that herbivorous *M. amblycephala* and *C. idellus*
278 have a large number of beta-type OR, while other omnivorous and carnivorous fish only have one
279 or two. This might be attributed to their particular herbivorous diet consisting not only of aquatic
280 grasses but also the duckweed and terrestrial grasses, which they ingest from the water surface.

281 It is known that the receptor for umami is formed by the T1R1/T1R3 heterodimer, while
282 T1R2/T1R3 senses sweet taste [24]. We found that the umami gene *T1R1* was lost in herbivorous
283 *M. amblycephala* but duplicated in the carnivorous *G. morhua* and *C. Semilaevis*. The loss of the
284 *T1R1* gene in *M. amblycephala* might exclude the expression of a functional umami taste receptor.
285 Such situations in other organism, e.g. the Chinese panda, have previously been related to feeding
286 specialization [25]. Bitterness sensed by the *T2R* is particularly crucial for animals to protect them
287 from poisonous compounds [22]. Interestingly, the bitter receptor *T2R* genes are expanded in the
288 herbivorous fish but few or no copy was found in carnivorous fish. These results not only indicate

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

289 the genetic adaptation to herbivorous diet of *M. amblycephala*, but also provided a clear and
290 comprehensive picture of adaptive evolutionary mechanisms of sensory systems in other fish
291 species with different trophic specializations.

292 It has been reported that some insects such as *Tenebrio molitor* [26] and *Neotermes*
293 *koshunensis* [27], and the mollusc *Corbicula japonica* [28] have genes encoding endogenous
294 cellulose degradation-related enzymes. However, all so far analyzed herbivorous vertebrates lack
295 these genes and always rely on their gut microbiome to digest food [25, 29]. In herbivorous *M.*
296 *amblycephala* and *C. idellus*, we also did not find any homologues of digestive cellulase genes.
297 Interestingly, our work on the composition of gut microbiota of the two fish species identifies
298 more than 7% cellulose-degrading bacteria, suggesting that the cellulose degradation of
299 herbivorous fish largely depend on their gut microbiome.

300 IB has evolved several times during teleost evolution [9, 30]. The developmental mechanisms
301 and ossification processes forming IB are dramatically distinct from other bones such as ribs,
302 skeleton, vertebrae or spines. These usually develop from cartilaginous bone and are derived from
303 the mesenchymal cell population by endochondral ossification [31, 32]. However, IB form directly
304 by intramembranous ossification and differentiate from osteoblasts within connective tissue,
305 forming segmental, serially homologous ossifications in the myosepta. Although various methods
306 of ossification of IB have been proposed, few experiments have been conducted to confirm the
307 ossification process and little is known about the potential role of IB in teleosts. Based on our
308 findings of expanded genes in cyprinid lineage and evidence from transcriptome of developmental
309 stages of IB formation, a number of genes were found to interact dynamically to mediate efficient
310 cell motility, migration and muscle construction [33-36]. In addition, transcriptome analyses of
311 three tissues indicated that ECM, Rho GTPase, motor and calcium channel regulation protein
312 displayed high expression in IB. It is known that ECM proteins bound to integrins influence cell
313 migration by actomyosin-generated contractile forces [34, 37]. Rho GTPases, acting as molecular
314 switches, are also involved in regulating the actin cytoskeleton and cell migration, which in turn
315 initiates intracellular signaling and contributes to tissue repair and regeneration [38-40]. Thus, our
316 results provide molecular evidence that IB might play significant roles not only in regulating
317 muscle contraction but also in active remodeling at the bone-muscle interface and coordination of

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

318 cellular events.

319 Some major developmental signals including BMP, FGF, WNT, together with
320 calcium/calmodulin signaling [31, 41-43], are essential for regulating the differentiation and
321 function of osteoblasts and osteocytes and for regulating the RANKL signaling pathway for
322 osteoclasts [44]. In agreement with this concept, we found 35 bone formation regulatory genes
323 involved in these signals were highly up-regulated in IB. Among these signaling pathways, in
324 particular, *Bmp*, *Fgf2*, and *Fgfr1* are closely related to intramembranous bone development and
325 affect the expression and activity of other osteogenesis related transcription factors [31, 45]. The
326 calcium-sensitive transcription factor *NFATc1* together with *CREB* induces the expression of
327 osteoclast-specific genes [46]. Taken together, these results suggest that IB indeed undergoes an
328 intramembranous ossification process, is regulated by bone-specific signaling pathways, and
329 underlies a homeostasis of maintenance, repair and remodeling.

330 **Conclusions**

331 Our results provide novel functional insights into the evolution of cyprinids. Importantly, the *M.*
332 *amblycephala* genome data come up with novel insights shedding light on the adaptation to
333 herbivorous nutrition and evolution and formation of IB. Our results on the evolution of gene
334 families, digestive and sensory system, as well as our microbiome meta-analysis and
335 transcriptome data provide powerful evidence and a key database for future investigations to
336 increase the understanding of the specific characteristics of *M. amblycephala* and other fish
337 species.

338 339 **Methods**

340 **Sampling and DNA Extraction**

341 DNA for genome sequencing was derived from a double haploid fish from the *M. amblycephala*
342 genetic breeding center at Huazhong Agricultural University (Wuhan, Hubei, China). Fish blood
343 was collected from adult female fish caudal vein using sterile injectors with pre-added
344 anticoagulant solutions following anesthetized with MS-222 and sterilization with 75% alcohol.
345 Genomic DNA was extracted from the whole blood.

346

347 **Genomic Sequencing and Assembly**

348 Libraries with different insert sized inserts of 170 bp, 500 bp, 800 bp, 2 Kb, 5 Kb, 10 Kb and 20
349 Kb were constructed from the genomic DNA at BGI-Shenzhen. The libraries were sequenced
350 using a HiSeq2000 instrument. In total, 11 libraries, sequenced in 23 lanes were constructed. To
351 obtain high quality data, we applied filtering criteria for the raw reads. As a result, 142.55 Gb of
352 filtered data were used to complete the genome assembly using SOAPdenovo_V2.04
353 (SOAPdenovo2 , RRID:SCR_014986) [16]. Only filtered data were used in the genome assembly.
354 First, the short insert size library data were used to construct a de Bruijn graph. The tips, merged
355 bubbles and connections with low coverage were removed before resolving the small repeats.
356 Second, all high-quality reads were realigned with the contig sequences. The number of shared
357 paired-end relationships between pairs of contigs was calculated and weighted with the rate of
358 consistent and conflicting paired ends before constructing the scaffolds in a stepwise manner from
359 the short-insert size paired ends to the long-insert size paired ends. Third, the gaps between the
360 constructed scaffolds were composed mainly of repeats, which were masked during scaffold
361 construction. These gaps were closed using the paired-end information to retrieve read pairs in
362 which one end mapped to a unique contig and the other was located in the gap region.
363 Subsequently, local assembly was conducted for these collected reads. To assess the genome
364 assembly quality, approximately 42.82 Gb Illumina reads generated from short-insert size libraries
365 were mapped onto the genome. Bwa0.5.9-r16 software (BWA , RRID:SCR_010910) [47] with
366 default parameters was used to assess the mapping ratio and Soap coverage 2.27 was used to
367 calculate the sequencing depth. We also assessed the accuracy of the genome assembly by Trinity
368 (Trinity , RRID:SCR_013048) [48], including number of ESTs and new mRNA reads from early
369 stages of embryos and multiple tissues, by aligning the scaffolds to the assembled transcriptome
370 sequences.

371 After obtaining K-mers from the short-insert-size (<1Kb) reads with just one bp slide,
372 frequencies of each K-mer were calculated. The K-mer frequency fits Poisson distribution when a
373 sufficient amount of data is present. The total genome size was deduced from these data in the
374 following way: Genome size = K-mer num / Peak_depth.

375

376 **Genome Annotation**

377 The genome was searched for repetitive elements using Tandem Repeats Finder (version 4.04)
378 [49]. TEs were identified using homology-based approaches. The Repbase (version 16.10) [50]
379 database of known repeats and a *de novo* repeat library generated by RepeatModeler
380 (RepeatModeler, RRID:SCR_015027) were used. This database was mapped using the software of
381 RepeatMasker (RepeatMasker , RRID:SCR_012954) (version 3.3.0). Four types of non-coding
382 RNAs (microRNAs, transfer RNAs, ribosomal RNAs and small nuclear RNAs) were also
383 annotated using tRNAscan-SE (version 1.23) and the Rfam database45 (Rfam ,
384 RRID:SCR_007891) (Release 9.1) [51].

385 For gene prediction, *de novo* gene prediction, homology-based methods and RNA-seq data
386 were used to perform gene prediction. For the sequence similarity based prediction, *D. rerio*, *G.*
387 *aculeatus*, *O. niloticus*, *O. latipes* and *G. morhua* protein sequences were downloaded from
388 Ensembl (Ensembl , RRID:SCR_002344)(release 73) and were aligned to the *M. amblycephala*
389 genome using TBLASTN (TBLASTN , RRID:SCR_011822). Then homologous genome
390 sequences were aligned against the matching proteins using GeneWise [52] to define gene models.
391 Augustus was employed to predict coding genes using appropriate parameters in *de novo*
392 prediction. For the RNA-seq based prediction, we mapped transcriptome reads to the genome
393 assembly using TopHat (TopHat , RRID:SCR_013035) [53]. Then, we combined TopHat mapping
394 results together and applied Cufflinks (Cufflinks , RRID:SCR_014597)[54] to predict transcript
395 structures. All predicted gene structures were integrated by GLEAN [55]
396 (<http://sourceforge.net/projects/glean-gene/>) to obtain a consensus gene set. Gene functions were
397 assigned to the translated protein-coding genes using Blastp tool (BLASTP , RRID:SCR_001010),
398 based on their highest match to proteins in the SwissProt and TrEMBL [56] databases (UniProt ,
399 RRID:SCR_002380)(Uniprot release 2011-01). Motifs and domains in the protein-coding genes
400 were determined by InterProScan (InterProScan , RRID:SCR_005829)(version 4.7) searches
401 against six different protein databases: ProDom, PRINTS, Pfam, SMART, PANTHER and
402 PROSITE. Gene Ontology (GO , RRID:SCR_002811)[57] IDs for each gene were obtained from
403 the corresponding InterPro entries. All genes were aligned against KEGG (KEGG ,

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

404 RRID:SCR_012773)[58] (Release 58) database, and the pathway in which the gene might be
405 involved was derived from the matched genes in KEGG. tRNA genes were *de novo* predicted by
406 tRNAscan-SE software (tRNAscan-SE , RRID:SCR_010835)[59], with eukaryote parameters on
407 the repeat pre-masked genome. The rRNA fragments were identified by aligning the rRNA
408 sequences using BlastN at E-value 1e-5 (BLASTN, RRID:SCR_001598). The snRNA and miRNA
409 were searched by the method of aligning and searching with INFERNAL (Infernal ,
410 RRID:SCR_011809)(version 0.81) [60] against Rfam database (Rfam ,
411 RRID:SCR_007891)(release 9.1).

412 413 **Genetic Map Construction**

414 To anchor the scaffolds into pseudo-chromosomes, 198 F1 population individuals were used to
415 obtain the genetic map. Each of the individual genomic DNA was digested with the restriction
416 endonuclease EcoR I, following the RAD-Seq protocol [61]. The SNP calling process was carried
417 out using the SOAP bioinformatic pipeline. The RAD-based SNP calling was done by SOAPsnp
418 software (SOAPsnp , RRID:SCR_010602)[62] after each individual's paired-end RAD reads was
419 mapped onto the assembled reference genome with the alignment software SOAP2
420 (SOAPaligner/soap2 , RRID:SCR_005503)[63]. The potential SNP markers were used for the
421 linkage analysis if the following criteria were satisfied: for parents - sequencing depth ≥ 8 and
422 ≤ 100 , base quality ≥ 25 , copy number ≤ 1.5 ; for progeny - sequencing depth ≥ 5 , base quality \geq
423 20, copy number ≤ 1.5 . If the markers were showing significantly distorted segregation (P -value $<$
424 0.01), they were excluded from the map construction. Linkage analysis was performed only for
425 markers present in at least 80% of the genomes, using JoinMap 4.0 software (JOINMAP ,
426 RRID:SCR_009248) with CP population type codes and applying the double pseudo-test cross
427 strategy [64]. The linkage groups were formed at a logarithm of odds threshold of 6.0 and ordered
428 using the regression mapping algorithm.

429 430 **Construction of Gene Families**

431 We identified gene families using TreeFam software (Tree families database ,
432 RRID:SCR_013401)[65] as follows: Blast was used to compare all the protein sequences from 13

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

433 species: *M. amblycephala*, *C. idellus*, *C. semilaevis*, *C. carpio*, *D. rerio*, *Callorhinchus milii*, *G.*
434 *morhua*, *G. aculeatus*, *Latimeria chalumnae*, *Oncorhynchus mykiss*, *O. niloticus*, *O. latipes*, *Fugu*
435 *rubripes*, with the E-value threshold set as 1e-7. In the next step, HSP segments of each protein
436 pair were concatenated by Solar software. H-scores were computed based on Bit-scores and these
437 were taken to evaluate the similarity among genes. Finally, gene families were obtained by
438 clustering of homologous gene sequences using Hcluster_sg (Version 0.5.0). Specific genes of *M.*
439 *amblycephala* were those that did not cluster with other vertebrates that were chosen for gene
440 family construction, and those that did not have homologs in the predicted gene repertoire of the
441 compared genomes. If these genes had functional motifs, they were annotated by GO.

442 443 **Phylogenetic Tree Reconstruction and Divergence Time Estimation**

444 The coding sequences of single-copy gene families conserved among *M. amblycephala*, *C. idellus*,
445 *C. carpio*, *D. rerio*, *C. semilaevis*, *G. morhua*, *G. aculeatus*, *Latimeria chalumnae*, *O. mykiss*, *O.*
446 *niloticus*, *O. latipes*, *C. milii* and *Fugu rubripes* (Ensembl Gene v.77) were extracted and aligned
447 with guidance from amino-acid alignments created by the MUSCLE program (MUSCLE ,
448 RRID:SCR_011812)[66]. The individual sequence alignments were then concatenated to form one
449 supermatrix. PhyML (PhyML , RRID:SCR_014629)[67, 68] was applied to construct the
450 phylogenetic tree under an HKY85+gamma model for nucleotide sequences. ALRT values were
451 taken to assess the branch reliability in PhyML. The same set of codon sequences at position 2
452 was used for phylogenetic tree construction and estimation of the divergence time. The PAML
453 mcmctree program (PAML , RRID:SCR_014932) (PAML version 4.5) [69, 70] was used to
454 determine divergence times with the approximate likelihood calculation method and the correlated
455 molecular clock and REV substitution model.

456 457 **Gene Family Expansion and Contraction Analyses**

458 Protein sequences of *M. amblycephala* and 11 other related species (Ensembl Gene v.77)) were
459 used in BLAST searches to identify homologs. We identified gene families using CAFÉ [71],
460 which employs a random birth and death model to study gene gains and losses in gene families
461 across a user-specified phylogeny. The global parameter λ , which describes both the gene birth (λ)

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

462 and death ($\mu = -\lambda$) rate across all branches in the tree for all gene families, was estimated using
463 maximum likelihood. A conditional P -value was calculated for each gene family, and families
464 with conditional P -values less than the threshold (0.05) were considered as having notable gain or
465 loss. We identified branches responsible for low overall P -values of significant families.

467 **Detection of Positively Selected Genes**

468 We calculated Ka/Ks ratios for all single copy orthologs of *M. amblycephala* and *C. semilaevis*, *D.*
469 *rerio*, *G. morhua*, *O. niloticus* and *C. carpio*. Alignment quality was essential for estimating
470 positive selection. Thus, orthologous genes were first aligned by PRANK [72], which is
471 considerably conservative for inferring positive selection. We used Gblocks [73] to remove
472 ambiguously aligned blocks within PRANK alignments and employed ‘codeml’ in the PAML
473 package with the free-ratio model to estimate Ka, Ks, and Ka/Ks ratios on different branches. The
474 differences in mean Ka/Ks ratios for single-copy genes between *M. amblycephala* and each of the
475 other species were compared using paired Wilcoxon rank sum tests. Genes that showed values of
476 Ka/Ks higher than 1 along the branch leading to *M. amblycephala* were reanalyzed using the
477 codon based branch-site tests implemented in PAML (PAML , RRID:SCR_014932). The
478 branch-site model allowed ω to vary both among sites in the protein and across branches, and was
479 used to detect episodic positive selection.

481 **Developmental Process of Intermuscular Bone in *M. amblycephala***

482 To better understand the number and morphological types of IBs in adult *M. amblycephala*,
483 specimens with a body length ranging from 15.5 to 20.5 cm were collected and each individual
484 was wrapped in gauze and boiled. The fish body was divided into two sections: anterior (snout to
485 cloaca) and posterior (cloaca to the base of caudal fin), and the length of each section was
486 measured. The IBs were retrieved, counted, arranged in order and photographed with a digital
487 camera. Fertilized *M. amblycephala* eggs were brought from hatching facilities at Freshwater Fish
488 Genetics Breeding Center of Huazhong Agricultural University (Wuhan, Hubei, China) to our
489 laboratory. *M. amblycephala* larvae were maintained in a re-circulating aquaculture system at $23 \pm$
490 1°C with a 14-hr photoperiod. To explore the early development of IBs, larvae at different stages

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

491 from 15 to 40 dpf were collected and fixed in 4% paraformaldehyde and transferred to 70%
492 ethanol for storage. Specimens were stained with alizarin red for bone following the method
493 described by Dawson [74]. The appearance of red color was recorded as the appearance of IB
494 because bone ossification is accompanied by the uptake of alizarin red, resulting in red staining of
495 the mineralized bone matrix. Myosepta, either not yet ossified, or poorly ossified, are not visible
496 with alizarin red staining. For histologic analysis, specimens were paraffin-embedded and
497 sectioned following standard protocols. Sections were stained with hematoxylin and eosin (HE)
498 and Masson trichrome [75] and photographed using a Nikon microscope (Nikon, Tokyo, Japan)
499 with a DP70 digital camera (Olympus, Japan). Scanning electron microscopy (SEM) and
500 transmission electron microscopy (TEM) were also conducted to analyze the ultrastructure of IB.
501 The specimens were fixed with 2.5% (v/v) glutaraldehyde in a solution of 0.1 M sodium
502 cacodylate buffer (pH 7.3) for 2 h at room temperature. The SEM and TEM samples were
503 prepared according to a standard protocol described by Ott [76]. The samples were then visualized
504 with a JSM-6390LV scanning electron microscope (SEM, Japan) and the stained ultrathin sections
505 with a H-7650 transmission electron microscope (Hitachi, Japan).

506 507 **RNA Sequencing Analysis**

508 *M. amblycephala* specimens belonging to three different developmental stages of IBs (stage 1: whole
509 larvae without distribution of IB; stage 2: muscle tissues with partial distribution of IBs; stage 3:
510 muscle tissues with completed distribution of IBs) were identified under microscope and immediately
511 frozen in liquid nitrogen. In addition, dorsal white muscle, IBs and connective tissue surrounding the
512 IBs from six months old fish were also collected. RNA was extracted from total fish samples at
513 different stages and from individual muscle, connective tissue, and intermuscular bone samples of
514 *M. amblycephala* using RNAisoPlus Reagent (TaKaRa, China) according to the manufacturer's
515 protocol. The integrity and purity of the RNA was determined by gel electrophoresis and Agilent
516 2100 BioAnalyzer (Agilent, USA) before preparing the libraries for sequencing. Paired-end RNA
517 sequencing was performed using the Illumina HiSeq 2000 platform. Low quality score reads were
518 filtered and the clean data were aligned to the reference genome using Bowtie [77]. Genes and
519 isoforms expression level were quantified by a software package: RSEM (RNASeq by

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

520 Expectation Maximization) (RSEM , RRID:SCR_013027)[78]. Gene expression levels were
521 calculated by using the RPKM method (Reads per kilobase transcriptome per million mapped
522 reads) [79] and adjusted by a scaling normalization method [80]. We detected DEGs from three
523 stages of IBs with software NOIseq and three different tissues with PossionDis as requested.
524 NOIseq is based on noisy distribution model, performed as described by Tarazona [81]. The
525 parameters were set as: fold change ≥ 2.00 and probability ≥ 0.7 . PossionDis is based on the
526 Poisson distribution, performed as described by Audic [82]. The parameters were set as: fold
527 change ≥ 2.00 and FDR ≤ 0.001 . Annotation of DEGs were mapped to GO categories in the
528 database (<http://www.geneontology.org/>) and the number of genes for every term were calculated
529 to identify GO terms that were significantly enriched in the input list of DEGs. The calculated
530 *P*-value was adjusted by the Bonferroni Correction, taking corrected *P*-value ≤ 0.05 as a threshold.
531 KEGG automatic annotation was used to perform pathway enrichment analysis of DEGs.

533 **Prediction of Olfactory Receptor Genes**

534 Olfactory receptor genes were identified by previously described methods [83], with the exception
535 of a first-round TBLASTN (TBLASTN , RRID:SCR_011822)[84] search, in which 1,417
536 functional olfactory receptor genes from *H. sapiens*, *D. rerio*, *L. chalumnae*, *Lepisosteus oculatus*,
537 *L. vexillifer*, *O. niloticus*, *O. latipes*, *F. rubripes* and *Xenopus tropicalis* were used as queries. We
538 then predicted the structure of sequenced genes using the blast-hit sequence with the software
539 GeneWise [52] extending in both 3' and 5' directions along the genome sequences. The results
540 were further confirmed by NR annotation. Then the coding sequences from the start (ATG) to stop
541 codons were extracted, while sequences that contained interrupting stop codons or frame-shifts
542 were regarded as pseudogenes. To construct phylogenetic trees, the amino-acid sequences encoded
543 by olfactory receptor genes were first aligned using the program MUSCLE nested in MEGA 5.10
544 (MEGA Software , RRID:SCR_000667)[85]. We then constructed the phylogenetic tree using the
545 neighbor-joining method with Poisson correction distances using the program MEGA 5.10. We
546 also identified and compared the genes for five basic tastes (sour, sweet, bitter, umami and salty)
547 using a similar method as in OR gene identification.

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

549 Gut Microbiota Analysis

550 To characterize the microbial diversity of herbivorous *M. amblycephala*, a total of 12 juvenile
551 (LBSB), domestic adult (DBSB), wild adult *M. amblycephala* (BSB) and wild adult *C. idellus*
552 (GC) intestinal fecal samples were collected. Bacterial genomic DNA was extracted from 200 mg
553 gut content of each sample using a QIAamp DNA Stool Mini Kit (Qiagen, Valencia, USA).
554 Quality and integrity of each DNA sample were determined by 1% agarose gel electrophoresis in
555 Tris-acetate-EDTA (TAE) buffer. DNA concentration was quantified using NanoDrop ND-2000
556 spectrophotometer (Thermo Scientific, Waltham, MA, USA). To determine the diversity and
557 composition of the bacterial communities of each sample, a total of 20 µg of genomic DNA were
558 sequenced using the Illumina MiSeq sequencing platform. PCR amplifications were conducted
559 from each sample to produce the V4 hypervariable region (515F and 806 R) of the 16S rRNA
560 gene according to the previously described method [86]. We used the UPARSE pipeline [87] to
561 pick operational taxonomic units (OTUs) at an identity threshold of 97% and picked
562 representative sequences for each OTU and used the RDP classifier to assign taxonomic data to
563 each representative sequence.

565 Additional files

566 Additional file 1: Tables S1 to S17 and Figures S1 to S28.

567 Additional file 2: Data Note1 Expanded genes in the *M. amblycephala* and *C. idellus* lineage.

568 Additional file 3: Data Note2 Positively selected genes in the *M. amblycephala* and *C. idellus*
569 genomes.

570 Abbreviations

571 BMP, bone morphogenetic protein; BUSCO, benchmarking universal single-copy orthologs;
572 DEGs, differentially expressed genes; dpf, days post fertilization; ECM, extracellular matrix; FGF,
573 fibroblast growth factor; HE, hematoxylin and eosin; IB, intermuscular bone; LG, linkage group;
574 LTR, long terminal repeat retrotransposon; PSG, positively selected gene; OR, olfactory receptor;
575 OTU, operational taxonomic unit; SEM, scanning electron microscopy; SNP, single-nucleotide
576 polymorphism; TAE, Tris-acetate-EDTA; TEM, transmission electron microscopy; TE,

577 transposable element

578

579 **Acknowledgements**

580 This work was supported by the Fundament Research Funds for the Central Universities
581 (2662015PY019), the Modern Agriculture Industry Technology System Construction Projects of
582 China titled as—Staple Freshwater Fishes Industry Technology System (No. CARS-46-05),
583 Guangdong Haid Group Co., Ltd and the International Scientific and Technology Cooperation
584 Program of Wuhan City (2015030809020365).

585 **Availability of data and materials**

586 Datasets and source images supporting the results of this article including are available in the
587 GigaDB repository associated with this publication [88]. Raw whole genome sequencing and
588 RAD-Seq data have been deposited at NCBI in the SRA under accession number SRP090157
589 (BioProject Number: PRJNA343584).

590 **Authors' contributions**

591 W.W. initiated and conceived the project and provided scientific input. X.Q. organized financial
592 support and designed the project. M.S. discussed the data, wrote and modified the paper. H.L. and
593 C.C. conducted the biological experiments, analyzed the data and wrote the paper with input from
594 other authors. I.E. wrote and modified the manuscript and discussed the data. The RAD-Seq data
595 analyses and the genetic map construction were performed by Z.G., Y.G., J.J. and X.J. Genome
596 assembly and annotation were performed by J.M., H.C., M.X. and J.C. X.Z., W.L., R.L., B.C.,
597 J.W., H.L., S.Y., H.W., X.C., X.Z., Y.Z., K.W., R.Y. and B. L. carried out the samples preparation
598 and data collection. J.L. and J.C. identified the gene families and analyzed the RNA-seq data. M.B.
599 coordinated the project. S.Z. and X.F. modified the manuscript and discussed the data. All authors
600 read the manuscript and provided comments and suggestions for improvements. The authors
601 declare no competing financial interests.

602 **Competing interests**

603 The authors declare that they have no competing interests.

604 **Ethics approval and consent to participate**

605 All experimental procedures involving fish were performed in accordance with the guidelines and
606 regulations of the National Institute of Health Guide for the Care and Use of Laboratory Animals
607 and the Institutional Review Board on Bioethics and Biosafety of BGI (BGI-IRB).

608 **References**

- 609 1. FAO Fisheries and Aquaculture Department. FAO yearbook Fishery and Aquaculture Statistics
610 2014 (Food and Agriculture Organization of the United Nations, Rome, 2016).
- 611 2. Star B, Nederbragt AJ, Jentoft S, Grimholt U, Malmstrøm M, Gregers TF, et al. The genome
612 sequence of Atlantic cod reveals a unique immune system. *Nature*. 2011; 477:207–10.
- 613 3. Berthelot C, Brunet F, Chalopin D, Juanchich A, Bernard M, Noël B, et al. The rainbow trout
614 genome provides novel insights into evolution after whole-genome duplication in vertebrates.
615 *Nat. Commun.* 2014; 5:3657.
- 616 4. Tine M., Kuhl H., Gagnaire PA, Louro B, Desmarais E, Martins RS, et al. European sea bass
617 genome and its variation provide insights into adaptation to euryhalinity and speciation. *Nat.*
618 *Commun.* 2014; 5:5770.
- 619 5. Wu C, Zhang D, Kan M, Lv Z, Zhu A, Su Y, et al. The draft genome of the large yellow
620 croaker reveals well-developed innate immunity. *Nat. Commun.* 2014; 5:5227.
- 621 6. Chen S, Zhang G, Shao C, Huang Q, Liu G, Zhang P, et al. Whole-genome sequence of a
622 flatfish provides insights into ZW sex chromosome evolution and adaptation to a benthic
623 lifestyle. *Nat. Genet.* 2014; 46:253–60.
- 624 7. Brawand D, Wagner CE, Li YI, Malinsky M, Keller I, Fan S, et al. The genomic substrate for
625 adaptive radiation in African cichlid fish. *Nature*. 2014; 513:375–82.
- 626 8. Chen X, Zhong L, Bian C, Xu P, Qiu Y, You X, et al. High-quality genome assembly of
627 channel catfish, *Ictalurus punctatus*. *GigaScience*. 2016; 5:39.
- 628 9. Gemballa S, Britz R. Homology of intermuscular bones in acanthomorph fishes. *Am. Mus.*
629 *Novit.* 1998; 3241:1–25.
- 630 10. Danos N, Ward AB. The homology and origins of intermuscular bones in fishes: Phylogenetic
631 or biomechanical determinants? *Biol. J. Linn. Soc.* 2012; 106:607–22.
- 632 11. Wan SM, Yi SK, Zhong J, Nie CH, Guan NN, Zhang WZ, et al. Dynamic mRNA and miRNA
633 expression analysis in response to intermuscular bone development of blunt snout bream
634 (*Megalobrama amblycephala*). *Sci. Rep.* 2016; 6:31050.
- 635 12. Xu P, Zhang X, Wang X, Li J, Liu G, Kuang Y, et al. Genome sequence and genetic diversity

- 636 of the common carp, *Cyprinus carpio*. *Nat. Genet.* 2014; 46:1212–9.
- 1
2 637 13. Wang Y, Lu Y, Zhang Y, Ning Z, Li Y, Zhao Q, et al. The draft genome of the grass carp
3 638 (*Ctenopharyngodon idellus*) provides insights into its evolution and vegetarian adaptation.
4
5 639 *Nat. Genet.* 2015; 47:625–31.
6
- 7 640 14. Gao Z, Luo W, Liu H, Zeng C, Liu X, Yi S, et al. Transcriptome analysis and SSR/SNP
8 641 markers information of the blunt snout bream (*Megalobrama amblycephala*). *PLoS One.*
9 642 2012; 7:e42637.
10
- 11 643 15. Yi S, Gao ZX, Zhao H, Zeng C, Luo W, Chen B, et al. Identification and characterization of
12 644 microRNAs involved in growth of blunt snout bream (*Megalobrama amblycephala*) by
13 645 Solexa sequencing. *BMC Genomics.* 2013; 14:754.
14
- 15 646 16. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved
16 647 memory-efficient short-read de novo assembler. *Gigascience.* 2012;1:18.
17
- 18 648 17. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing
19 649 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.*
20 650 2015;31:3210–2.
21
- 22 651 18. Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, et al. The genomic basis
23 652 of adaptive evolution in threespine sticklebacks. *Nature.* 2012; 484:55–61.
24
- 25 653 19. Niimura Y. On the origin and evolution of vertebrate olfactory receptor genes: comparative
26 654 genome analysis among 23 chordate species. *Genome Biol. Evol.* 2009; 1:34–44.
27
- 28 655 20. Lindemann B. Receptors and transduction in taste. *Nature.* 2001; 413:219–25.
29
- 30 656 21. Nei M, Niimura Y, Nozawa M. The evolution of animal chemosensory receptor gene
31 657 repertoires: roles of chance and necessity. *Nat. Rev. Genet.* 2008; 9:951–63.
32
- 33 658 22. Chandrashekar J, Mueller KL, Hoon MA, Adler E, Feng L, Guo W, et al. T2Rs function as
34 659 bitter taste receptors. *Cell.* 2000; 100:703–11.
35
- 36 660 23. Niimura Y, Nei M. Evolutionary dynamics of olfactory and other chemosensory receptor
37 661 genes in vertebrates. *J. Hum. Genet.* 2006; 51:505–17.
38
- 39 662 24. Nelson G, Chandrashekar J, Hoon MA, Feng L, Zhao G, Ryba NJ, et al. An amino-acid taste
40 663 receptor. *Nature.* 2002; 416:199–202.
41
- 42 664 25. Li R, Fan W, Tian G, Zhu H, He L, Cai J, et al. The sequence and *de novo* assembly of the
43 665 giant panda genome. *Nature.* 2010; 463:311–7.
44
- 45 666 26. Ferreira AHP, Marana SR, Terra WR, Ferreira C. Purification, molecular cloning, and
46 667 properties of a β -glycosidase isolated from midgut lumen of *Tenebrio molitor* (Coleoptera)
47 668 larvae. *Insect Biochem. Mol. Biol.* 2001; 31:1065–76.
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 669 27. Tokuda G, Saito H, Watanabe H. A digestive β -glucosidase from the salivary glands of the
670 termite, *Neotermes koshunensis* (Shiraki): Distribution, characterization and isolation of its
671 precursor cDNA by 5'- and 3'-RACE amplifications with degenerate primers. *Insect*
672 *Biochem. Mol. Biol.* 2002; 32:1681–9.
- 673 28. Sakamoto K, Uji S, Kurokawa T, Toyohara H. Molecular cloning of endogenous
674 β -glucosidase from common Japanese brackish water clam *Corbicula japonica*. *Gene*. 2009;
675 435:72–9.
- 676 29. Zhu L, Wu Q, Dai J, Zhang S, Wei F. Evidence of cellulose metabolism by the giant panda gut
677 microbiome. *Proc. Natl. Acad. Sci. USA*. 2011; 108:17714–9.
- 678 30. Bird NC, Mabee PM. Developmental morphology of the axial skeleton of the zebrafish, *Danio*
679 *rerio* (Ostariophysi: Cyprinidae). *Dev. Dyn.* 2003; 228:337–57.
- 680 31. Ornitz D, Marie P. FGF signaling pathways in endochondral and intramembranous bone
681 development and human genetic disease. *Genes Dev.* 2002; 16:1446–65.
- 682 32. Ortega N, Behonick DJ, Werb Z. Matrix remodeling during endochondral ossification. *Trends*
683 *Cell Biol.* 2004; 14:86–93.
- 684 33. Even-Ram S, Doyle AD, Conti MA, Matsumoto K, Adelstein RS, Yamada KM. Myosin IIA
685 regulates cell motility and actomyosin-microtubule crosstalk. *Nat. Cell Biol.* 2007;
686 9:299–309.
- 687 34. Sheetz MP, Felsenfeld DP, Galbraith CG. Cell migration: Regulation of force on
688 extracellular-complexes. *Trends Cell Biol.* 1998; 8:51–4.
- 689 35. Gunst SJ, Zhang W. Actin cytoskeletal dynamics in smooth muscle: a new paradigm for the
690 regulation of smooth muscle contraction. *Am. J. Physiol. Cell Physiol.* 2008; 295:C576–87.
- 691 36. Webb RC. Smooth muscle contraction and relaxation. *Adv. Physiol. Educ.* 2003; 27:201–6.
- 692 37. Ulrich TA, De Juan Pardo EM, Kumar S. The mechanical rigidity of the extracellular matrix
693 regulates the structure, motility, and proliferation of glioma cells. *Cancer Res.* 2009;
694 69:4167–74.
- 695 38. Ridley AJ. Rho GTPases and cell migration. *J. Cell Sci.* 2001; 114:2713–22.
- 696 39. Etienne-Manneville S, Hall A. Rho GTPases in Cell Biology. *Nature*. 2002; 420:629–35.
- 697 40. Ridley AJ. Cell Migration: Integrating Signals from Front to Back. *Science*. 2003;
698 302:1704–9.
- 699 41. Chen G, Deng C, Li YP. TGF- β and BMP signaling in osteoblast differentiation and bone
700 formation. *Int. J. Biol. Sci.* 2012; 8:272–88.
- 701 42. Harada S, Rodan GA. Control of osteoblast function and regulation of bone mass. *Nature*.

- 702 2003; 423:349–55.
- 1
2 703 43. Long F. Building strong bones: molecular regulation of the osteoblast lineage. *Nat. Rev. Mol.*
3 704 *Cell Biol.* 2011; 13:27–38.
- 4
5 705 44. Boyle WJ, Simonet WS, Lacey DL. Osteoclast differentiation and activation. *Nature.* 2003;
6 706 423:337–42.
- 7
8
9 707 45. Fakhry A, Ratisoontorn C, Vedhachalam C, Salhab I, Koyama E, Leboy P, et al. Effects of
10 708 FGF-2/-9 in calvarial bone cell cultures: Differentiation stage-dependent mitogenic effect,
11 709 inverse regulation of BMP-2 and noggin, and enhancement of osteogenic potential. *Bone.*
12 710 2005; 36:254–66.
- 13
14
15 711 46. Sato K, Suematsu A, Nakashima T, Takemoto-Kimura S, Aoki K, Morishita Y, et al.
16 712 Regulation of osteoclast differentiation and function by the CaMK-CREB pathway. *Nat. Med.*
17 713 2006; 12:1410–6.
- 18
19
20 714 47. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.
21 715 *Bioinformatics.* 2009; 25:1754–60.
- 22
23
24 716 48. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Trinity:
25 717 reconstructing a full-length transcriptome without a genome from RNA-Seq data . *Nat.*
26 718 *Biotechnol.* 2011; 29:644–52.
- 27
28
29 719 49. Benson G. Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids Res.*
30 720 1999; 27:573–80.
- 31
32
33 721 50. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. Repbase
34 722 Update, a database of eukaryotic repetitive elements. *Cytogenet. Genome Res.* 2005;
35 723 110:462–7.
- 36
37
38 724 51. Griffiths-Jones S, Moxon S, Marshall M, Khanna A, Eddy SR, Bateman A. Rfam: Annotating
39 725 non-coding RNAs in complete genomes. *Nucleic Acids Res.* 2005; 33:121–4.
- 40
41
42 726 52. Birney E, Clamp M, Durbin R. GeneWise and Genomewise. *Genome Res.* 2004; 14:988–95.
- 43
44
45 727 53. Trapnell C, Pachter L, Salzberg SL. TopHat: Discovering splice junctions with RNA-Seq.
46 728 *Bioinformatics.* 2009; 25:1105–11.
- 47
48
49 729 54. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, Van Baren MJ, et al. Transcript
50 730 assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform
51 731 switching during cell differentiation. *Nat. Biotechnol.* 2010; 28:511–5.
- 52
53
54 732 55. Elisk CG, Mackey AJ, Reese JT, Milshina NV, Roos DS, Weinstock GM. Creating a honey
55 733 bee consensus gene set. *Genome Biol.* 2007; 8:R13.
- 56
57
58 734 56. Bairoch A, Apweiler R. The SWISS-PROT protein sequence database and its supplement
59 735 TrEMBL in 2000. *Nucleic Acids Res.* 2000; 28:45–8.

- 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
- 736 57. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene ontology: tool
737 for the unification of biology. *Nat. Genet.* 2000; 25:25–9.
- 738 58. Kanehisa M, Goto S. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*
739 2000; 28:27-30.
- 740 59. Lowe TM, Eddy SR. tRNAscan-SE: A program for improved detection of transfer RNA genes
741 in genomic sequence. *Nucleic Acids Res.* 1997; 25:955–64.
- 742 60. Nawrocki EP, Kolbe DL, Eddy SR. Infernal 1.0: Inference of RNA alignments. *Bioinformatics.*
743 2009; 25:1335–7.
- 744 61. Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, et al. Rapid SNP
745 discovery and genetic mapping using sequenced RAD markers. *PloS One.* 2008; 3:e3376.
- 746 62. Li R, Li Y, Fang X, Yang H, Wang J, Kristiansen K, et al. SNP detection for massively
747 parallel whole-genome resequencing. *Genome Res.* 2009; 19:1124–32.
- 748 63. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, et al. SOAP2: An improved ultrafast tool
749 for short read alignment. *Bioinformatics.* 2009; 25:1966–7.
- 750 64. Grattapaglia D, Sederoff R. Genetic linkage maps of *Eucalyptus grandis* and *Eucalyptus*
751 *urophylla* using a pseudo-testcross: Mapping strategy and RAPD markers. *Genetics.* 1994;
752 137:1121–37.
- 753 65. Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, et al. TreeFam: a curated
754 database of phylogenetic trees of animal gene families. *Nucleic Acids Res.* 2006;
755 34:D572–80.
- 756 66. Edgar RC. MUSCLE: Multiple sequence alignment with high accuracy and high throughput.
757 *Nucleic Acids Res.* 2004; 32:1792–7.
- 758 67. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and
759 methods to estimate maximum-likelihood phylogenies: Assessing the performance of PhyML
760 3.0. *Syst. Biol.* 2010; 59:307–21.
- 761 68. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by
762 maximum likelihood. *Syst. Biol.* 2003; 52:696–704.
- 763 69. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 2007;
764 24:1586–91.
- 765 70. Yang Z, Rannala B. Bayesian estimation of species divergence times under a molecular clock
766 using multiple fossil calibrations with soft bounds. *Mol. Biol. Evol.* 2006; 23:212–26.
- 767 71. Hahn MW, Demuth JP, Han SG. Accelerated rate of gene gain and loss in primates. *Genetics.*
768 2007; 177:1941–9.

- 769 72. Löytynoja A, Goldman N. An algorithm for progressive multiple alignment of sequences with
770 insertions. *Proc. Natl. Acad. Sci. USA*. 2005; 102:10557–62.
- 771 73. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and
772 ambiguously aligned blocks from protein sequence alignments. *Syst. Biol.* 2007; 56:564–77.
- 773 74. Dawson AB. A note on the staining of the skeleton of cleared specimens with alizarin red S.
774 *Biotech. Histochem.* 1926; 1:123-4.
- 775 75. Gruber HE. Adaptations of Goldner’s Masson trichrome stain for the study of undecalcified
776 plastic embedded bone. *Biotech. Histochem.* 1992; 67:30–4.
- 777 76. Ott HC, Matthiesen TS, Goh SK, Black LD, Kren SM, Netoff TI, et al.
778 Perfusion-decellularized matrix: using nature’s platform to engineer a bioartificial heart. *Nat.*
779 *Med.* 2008; 14:213–21.
- 780 77. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat. Methods.* 2012;
781 9:357–9.
- 782 78. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or
783 without a reference genome. *BMC Bioinformatics.* 2011; 12:323.
- 784 79. Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying
785 mammalian transcriptomes by RNA-Seq. *Nat. Methods.* 2008; 5:621–8.
- 786 80. Robinson MD, Oshlack A. A scaling normalization method for differential expression analysis
787 of RNA-seq data. *Genome Biol.* 2010; 11:R25.
- 788 81. Tarazona S, García-Alcalde F, Dopazo J, Ferrer A, Conesa A. Differential expression in
789 RNA-seq: a matter of depth. *Genome Res.* 2011; 21:2213–23.
- 790 82. Audic S, Claverie JM. The significance of digital gene expression profiles. *Genome Res.* 1997;
791 7:986–95.
- 792 83. Niimura Y. Evolutionary dynamics of olfactory receptor genes in chordates: interaction
793 between environments and genomic contents. *Hum. Genomics.* 2009; 4:107–18.
- 794 84. Altschul S, Madden T, Schaffer A, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and
795 PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 1997;
796 25:3389–402.
- 797 85. Kumar S, Nei M, Dudley J, Tamura K. MEGA: A biologist-centric software for evolutionary
798 analysis of DNA and protein sequences. *Brief. Bioinform.* 2008; 9:299–306.
- 799 86. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ, et al.
800 Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. *Proc.*
801 *Natl. Acad. Sci. USA.* 2011; 108:4516–22.

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

802 87. Edgar, RC. UPARSE: highly accurate OTU sequences from microbial amplicon reads. Nat.
803 Methods. 2013; 10:996–8.

804 88. Liu, H; Chen, C; Gao, Z; Min, J; Gu, Y; Jian, J; Jiang, X; Cai, H; Ebersberger, I; Xu, M; Zhang,
805 X; Chen, J; Luo, W; Chen, B; Chen, J; Liu, H; Li, J; Lai, R; Bai, M; Wei, J; Yi, S; Wang, H;
806 Cao, X; Zhou, X; Zhao, Y; Wei, K; Yang, R; Liu, B; Zhao, S; Fang, X; Scharl, M; Qian, X;
807 Wang, W (2017): Supporting data for "The draft genome of blunt snout bream (*Megalobrama*
808 *amblycephala*) reveals the development of intermuscular bone and adaptation to herbivorous
809 diet" GigaScience Database. <http://dx.doi.org/10.5524/100305>

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

810 **Figure Legends**

811 **Figure 1** Image of an adult blunt snout bream (*Megalobrama amblycephala*).

812 **Figure 2** Global view of the *M. amblycephala* genome and syntenic relationship between *C.*
813 *idellus*, *M. amblycephala* and *D. rerio*. (A) Global view of the *M. amblycephala* genome. From
814 outside to inside, the genetic linkage map (a); Anchors between the genetic markers and the
815 assembled scaffolds (b); Assembled chromosomes (c); GC content within a 50-kb sliding window
816 (d); Repeat content within a 500-kb sliding window (e); Gene distribution on each chromosome (f);
817 Different gene expression of three transcriptomes (g). (B) Syntenic relationship between *C. idellus*
818 (a), *M. amblycephala* (b) and *D. rerio* (c) chromosomes.

819 **Figure 3** Phylogenetic tree and comparison of orthologous genes in *M. amblycephala* and other
820 fish species. (A) Phylogenetic tree of teleosts using 316 single copy orthologous genes. The color
821 circles at the nodes shows the estimated divergence times using *O. latipes*–*F. rubripes*
822 [96.9~150.9Mya], *F. rubripes*–*D. rerio* [149.85~165.2Mya], *F. rubripes*–*C. milii* [416~421.75Mya]
823 (<http://www.timetree.org/>) as the calibration time. Pentagon represents four cyprinid fish with
824 intermuscular bones. S, Silurian period; D, Devonian period; C, Carboniferous period; P, Permian
825 period in Paleozoic; T, Triassic period; J, Jurassic and k-cretaceous period in Mesozoic; Pg,
826 Paleogene in Cenozoic Era, N, Neogene. (B) Venn diagram of shared and unique orthologous gene
827 families in *M. amblycephala* and four other teleosts. (C) Over-represented GO annotations of
828 cyprinid-specific expansion genes.

829 **Figure 4** Regulation of genes related to intermuscular bone formation and function identified from
830 developmental stages and adult tissues transcriptome data. (A) Gene expression pattern involved
831 in muscle contraction regulated genes in early developmental stages corresponds to intermuscular
832 bone formation of *M. amblycephala*, (alizarin red staining). M, myosepta; IB, intermuscular bone.
833 (B) Scanning electron microscope photos of muscle tissues, connective tissues, and intermuscular
834 bone. (C) Distribution of intermuscular bone specific genes in GO annotations indicative of
835 abundance in protein binding, calcium ion binding, GTP binding functions. (D) Several
836 developmental signals regulating key steps of osteoblast and osteoclast differentiation in the
837 process of intramembranous ossification. Colored boxes indicate significantly up-regulated genes

838 in these signals specifically occurred in intermuscular bone.

839 **Figure 5** Molecular characteristics of sensory systems and the composition of gut microbiota in *M.*
840 *amblycephala*. (A) Extensive expansion of olfactory receptor genes (ORs) in *M. amblycephala*
841 compared with other teleosts. (B) Phylogeny of ‘beta’ type ORs in eight representative teleost
842 species showing the significant expansion of ‘beta’ ORs in *M. amblycephala* and *C. idellus*. The
843 pink background shows cyprinid-specific ‘beta’ types of ORs. (C) Umami, sweet and bitter tastes
844 related gene families in teleosts with different feeding habits. (D) Structure of the umami receptor
845 encoding T1R1 gene in cyprinid fish. (E) Relative abundance of microbial flora and taxonomic
846 assignments in juvenile (LBSB), domestic adult (DBSB), wild adult (BSB) *M. amblycephala* and
847 wild adult *C. idellus* (GC) samples at the phylum level.

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863 **Table**

864 **Table 1 Features of the *M. amblycephala* whole genome sequence**

| | | |
|----|--|--------------|
| 4 | Total genome size (Mb) | 1,116 |
| 5 | N90 length of scaffold (bp) | 20,422 |
| 6 | N50 length of scaffold (bp) | 838,704 |
| 7 | N50 length of contig (bp) | 49,400 |
| 8 | Total GC content (%) | 37.30 |
| 9 | Protein-coding genes number | 23,696 |
| 10 | Average gene length (bp) | 15,797 |
| 11 | Content of transposable elements (%) | 34.18 |
| 12 | Number of chromosomes | 24 |
| 13 | Number of makers in genetic map | 5,317 |
| 14 | Scaffolds anchored on linkage groups (LGs) | 1,434 |
| 15 | Length of scaffolds anchored on LGs (Mb) | 779.54 (70%) |

865

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 2 May 2017

2 Dr. Hans Zauner

3 Journal: GigaScience

4
5
6
7 Dear Dr. Zauner,

8
9 **Manuscript No.: GIGA-D-16-00088R2**

10 Title: "**The draft genome of blunt snout bream (*Megalobrama amblycephala*) reveals the**
11 **development of intermuscular bone and adaptation to herbivorous diet"**

12
13 Author(s): Han Liu, Chunhai Chen, Zexia Gao, Jiumeng Min, Yongming Gu, Jianbo Jian, Xiewu
14 Jiang, Huimin Cai, Ingo Ebersberger, Meng Xu, Xinhui Zhang, Jianwei Chen, Wei Luo, Boxiang
15 Chen, Junhui Chen, Hong Liu, Jiang Li, Ruifang Lai, Mingzhou Bai, Jin Wei, Shaokui Yi,
16 Huanling Wang, Xiaojuan Cao, Xiaoyun Zhou, Yuhua Zhao, Kaijian Wei, Ruibin Yang, Bingnan
17 Liu, Shancen Zhao, Xiaodong Fang, Manfred Schartl, Xueqiao Qian, Weimin Wang
18
19
20
21
22
23

24 We have carefully read the referee's comments which you forwarded to us of 27 April 2017. We
25 would like to express our sincere thanks for the positive comments. We have addressed all the
26 suggestions. The common name of the species has been included in the title (Line 1 to 2). The
27 Fishbase number and image of an adult blunt snout bream have also been added in the revised
28 manuscript (Line 73 to 74). The amendments are highlighted in red in the revised manuscript.
29 Responses to the reviewer's comments are detailed below in this letter. We hope that with the
30 amendments made in response to you and the reviewer's comments, the manuscript is now
31 acceptable for publication in GigaScience.
32
33
34
35
36
37
38

39 I look forward to hearing from you soon.

40
41 Yours sincerely,

42
43 Weimin Wang (PhD) (Correspondence author)

44
45 College of Fisheries

46
47 Huazhong Agricultural University

48
49 Wuhan 430070, P. R. China

50
51 E-mail address: wangwm@mail.hzau.edu.cn

52
53 Tel: +86-27-8728 4292; Fax: +86-27-8728 4292
54
55
56
57
58
59
60
61
62
63
64
65

Response to Reviewer

Reviewer Report

Reviewer #2:

1. Have a look at the sentence at line 104-107: “to assess the genome assembly quality”

Author response: This sentence has been modified as “To assess the quality of genome assembly, the short-insert size paired-end libraries reads and published ESTs [14] (Additional file 1: Tables S3 and S4) were mapped onto the genome. The results indicated that the assembled error is low.” (Line 105 to 107)

2. Line 151 “single –copy genes”?

Author response: This expression has been changed as “single-copy orthologous genes”. (Line 152)

3. 153: “outgroup” not “out group”

Author response: This has been corrected.

4. Rephrase line 552.

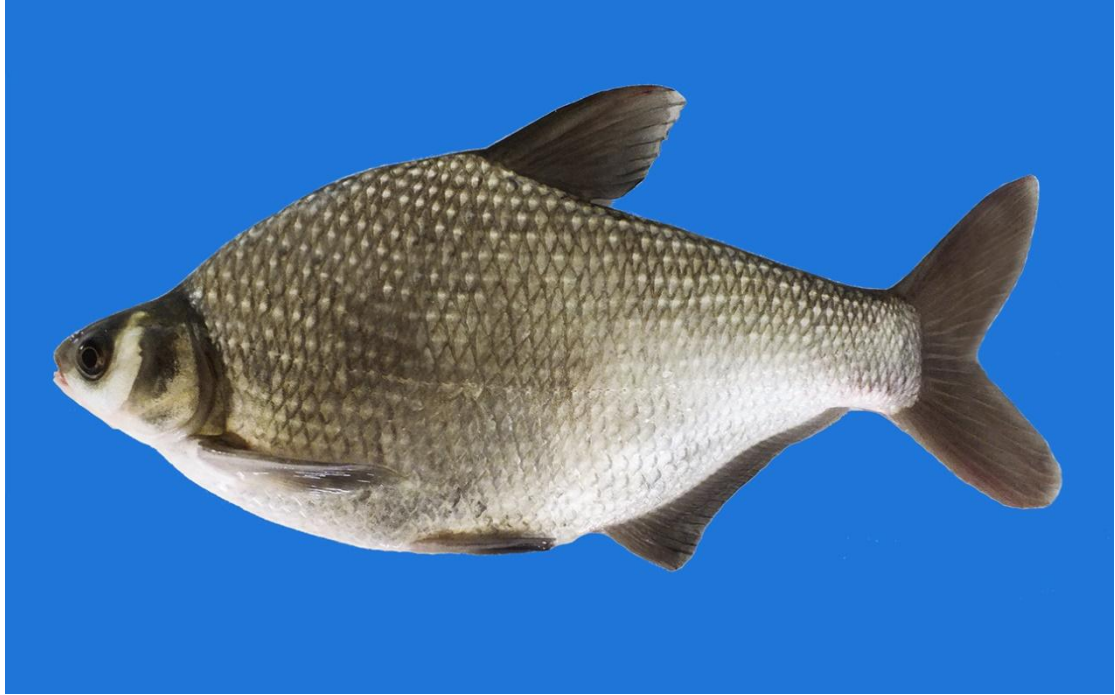
Author response: The SRP accession number has been added in the revised manuscript. This sentence has been rephrased as “Raw whole genome sequencing and RAD-Seq data have been deposited at NCBI in the SRA under accession number SRP090157 (BioProject Number: PRJNA343584)”. (Line 553 to 554)

5. 3 Additional file 2 Data note 1. Typo: “Expansion” instead of “Expasion”

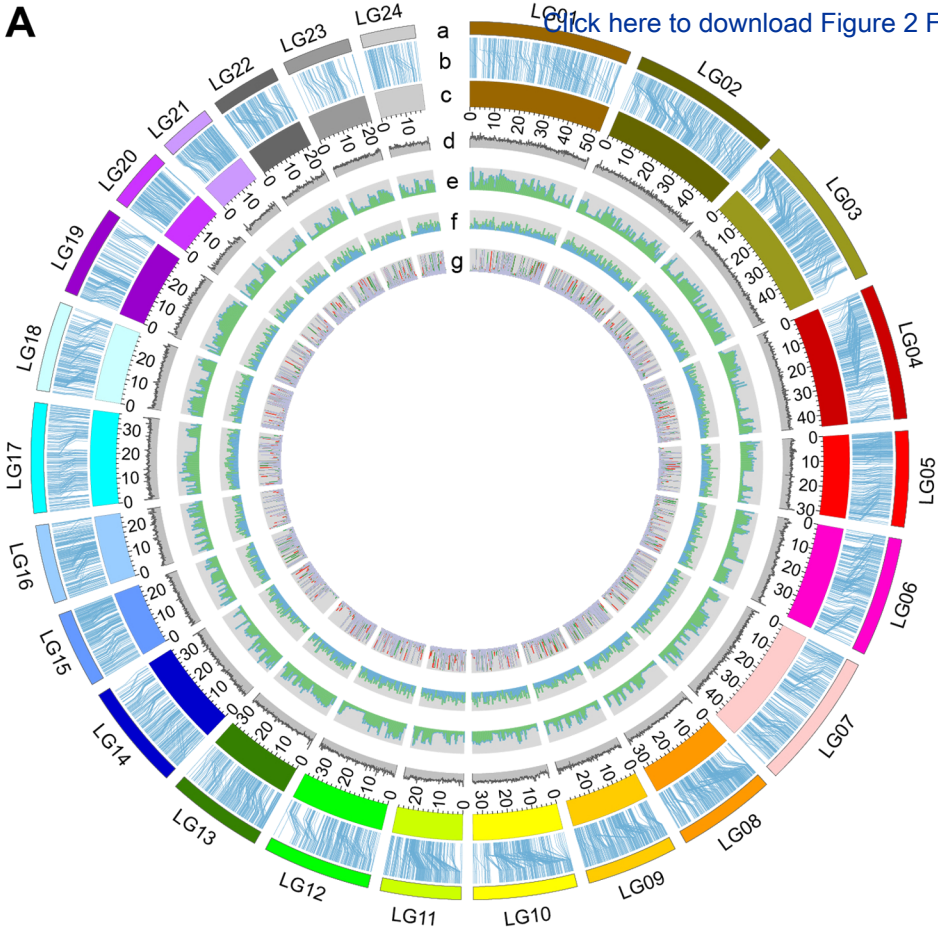
Author response: This has been corrected.

6. Figure 2B not mentioned in the text.

Author response: We have now mentioned it in the text and expressed as “We found 9349 orthologous gene families shared among five fish species. 246 are specific in the *M. Amblycephala* (Figure 3B)”. (Line 155 to 156)



A



B

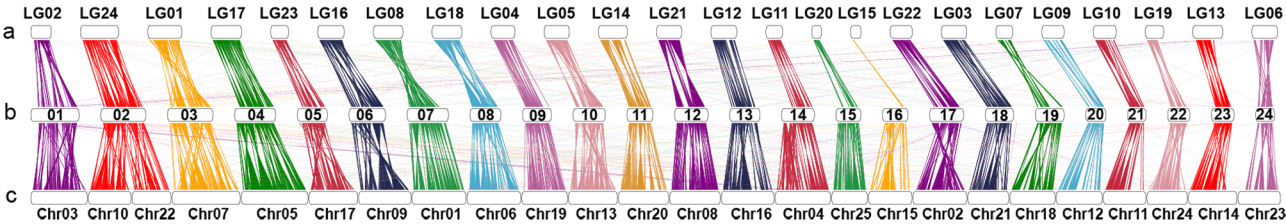
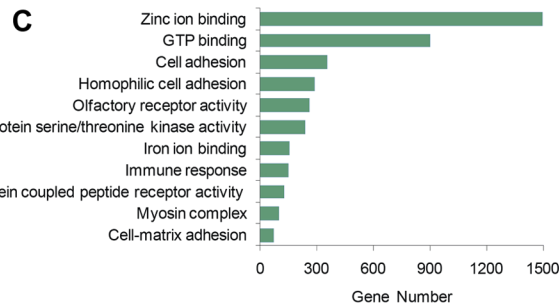
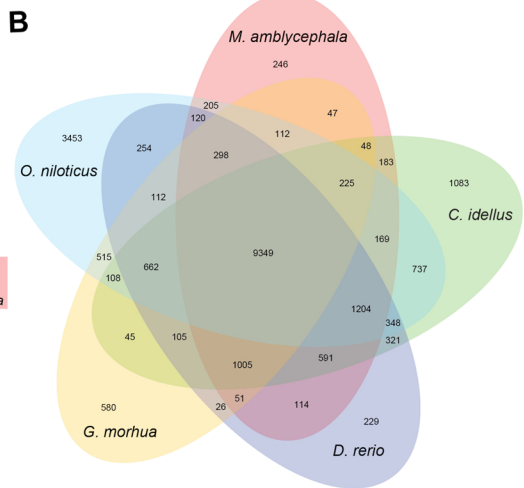
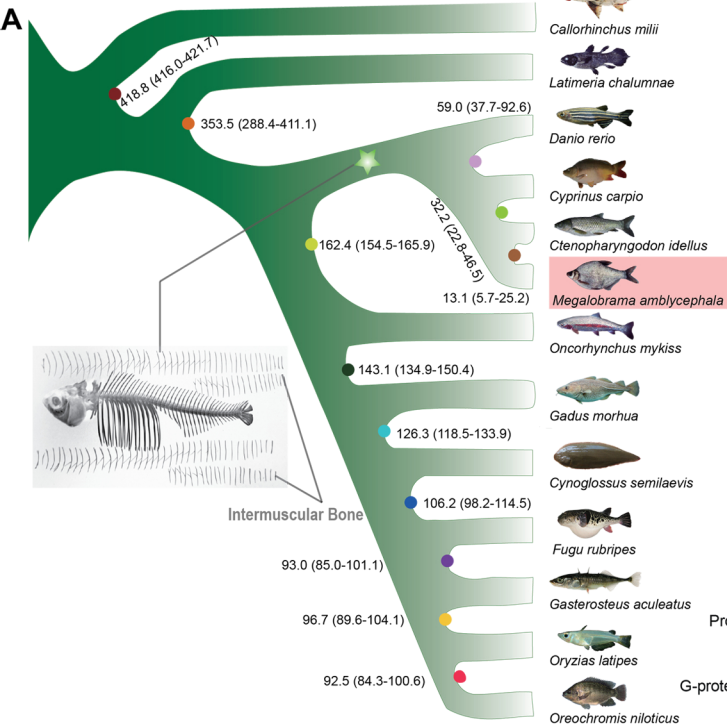
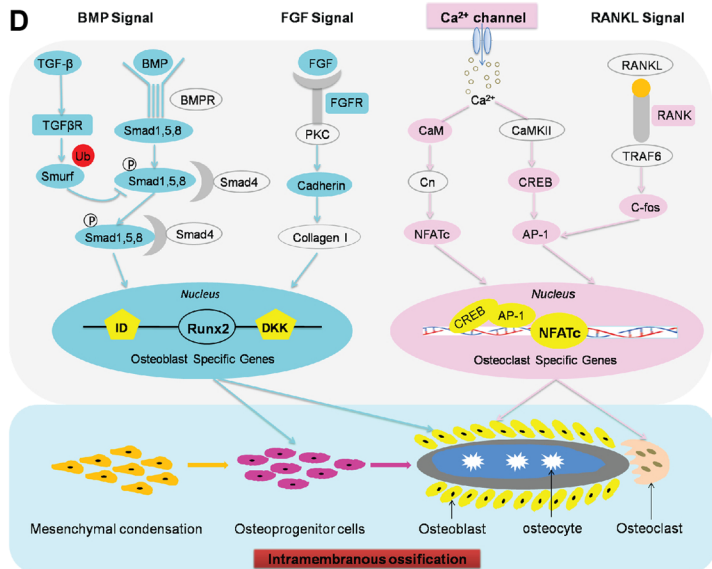
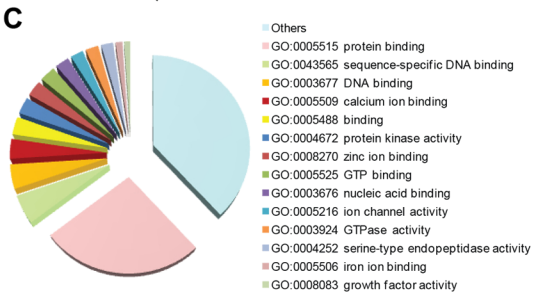
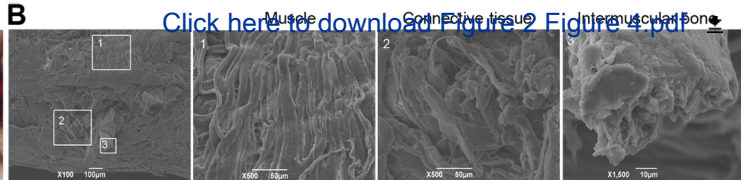
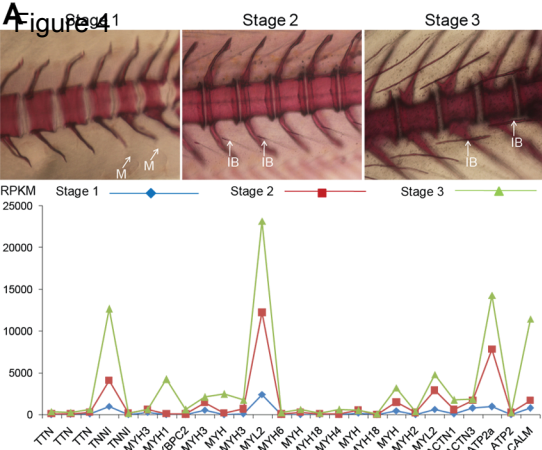
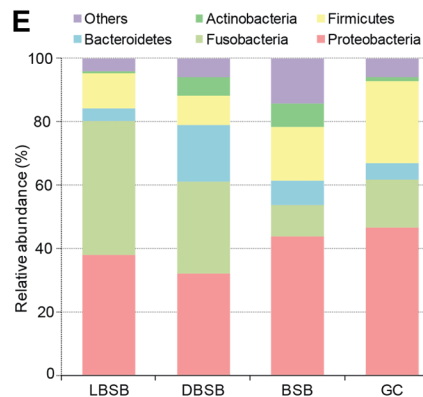
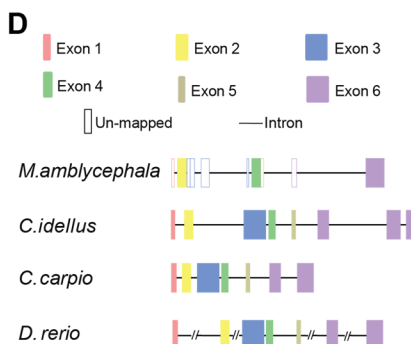
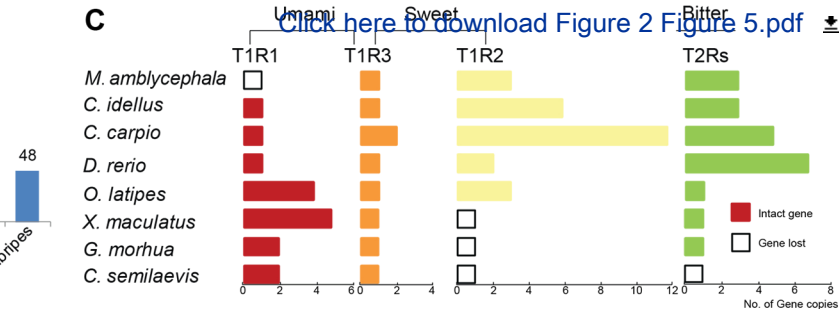
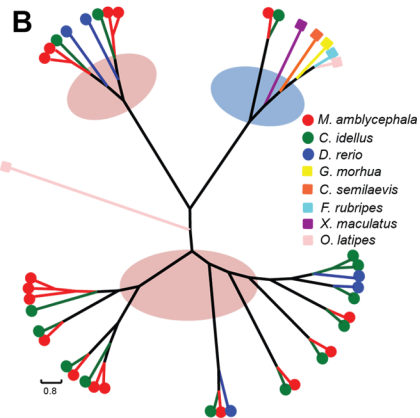
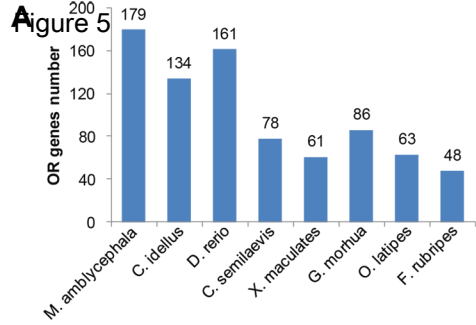


Figure 3

[Click here to download Figure 2 Figure 3.pdf](#)









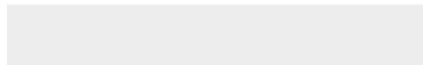
Click here to access/download

Supplementary Material

3 Additional file 1 Tables S1 to S17 and Figures S1 to
S28.pdf



Click here to access/download
Supplementary Material
3 Additional file 2 Data note1.xlsx





Click here to access/download
Supplementary Material
3 Additional file 3 Data note2.xlsx

