

## Draft genome of *Glyptosternon maculatum*, an endemic fish from Tibet-plateau --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-18-00128R2	
<b>Full Title:</b>	Draft genome of <i>Glyptosternon maculatum</i> , an endemic fish from Tibet-plateau	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	the special finance of Tibet autonomous region (2017CZZX003)	Dr haiping liu
	National Natural Science Foundation of China (31560144)	Dr haiping liu
	National Natural Science Foundation of China (31602207)	Dr Shijun Xiao
<b>Abstract:</b>	<p><b>Abstracts</b></p> <p><b>Background</b></p> <p>Mechanisms for high altitude adaption have arisen widespread interest to evolution biologists. Several genome wide studies have been carried out for endemic vertebrates in Tibet, including mammals, birds and amphibians. However, little information was known about the adaptive evolution of highland fishes. <i>Glyptosternon maculatum</i> (<i>G. maculatum</i>, Regan, 1905), also known as Regan or barkley, is a fish endemic to the Tibetan plateau, which belongs to Sisoridae family, Siluriformes (catfishes) order. This species live within an elevation ranging from roughly 2800 m to 4200 m. Hence, a high-quality reference genome of <i>G. maculatum</i> provides an opportunity to address high altitude adaption mechanisms of fishes.</p> <p><b>Findings</b></p> <p>To get a high-quality reference genome of <i>G. maculatum</i>, we combined PacBio single-molecule real-time sequencing, Illumina paired-end sequencing, 10X Genomics linked-reads and BioNano optical map techniques. In total, 603.99 Gb sequencing data were generated. The assembled genome was about 662.34 Mb with scaffold and contig N50 sizes of 20.90 Mb and 993.67 kb, respectively, which captured 83% complete and 3.9% partial vertebrate Benchmarking Universal Single-copy orthologs (BUSCO). Repetitive elements account for 35.88% of the genome, and 22,066 protein-coding genes were predicted from the genome, of which 91.7% have been functionally annotated.</p> <p><b>Conclusions</b></p> <p>We provide the first comprehensive de novo genome of the <i>G. maculatum</i>. This genetic resource is fundamental for investigating the origin of the <i>G. maculatum</i> and will improve our understanding of high altitude adaption of fishes. The assembled genome can also be used as reference for future population genetic studies of <i>G. maculatum</i>.</p>	
<b>Corresponding Author:</b>	haiping liu CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	haiping liu	

<b>First Author Secondary Information:</b>	
<b>Order of Authors:</b>	<p>haiping liu</p> <p>Qiyong Liu</p> <p>Zhiqiang Chen</p> <p>Yanchao Liu</p> <p>Chaowei Zhou</p> <p>Qiqi Liang</p> <p>Caixia Ma</p> <p>Jianshe Zhou</p> <p>Yingzi Pan</p> <p>Meiqun Chen</p> <p>wangjiu wangjiu</p> <p>Wenkai Jiang</p> <p>Shijun Xiao</p> <p>Zhenbo Mou</p>
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Comments from the editor:</p> <p>Comment: Please make sure the data is in the SRA and the accession is included in the paper before resubmission.</p> <p>Reply: We have submit data used in this work to NCBI with SRA accession of SRR7279473-SRR7279474, SRR7268130-SRR7268162, SRR7350914-SRR7350921, SRR7351269-SRR7351265, SRR7403445-SRR7403454 under the Project accession number of PRJNA447978 for Illumina, PacBio, 10X Genomics and transcriptome sequencing data used in this work.</p> <p>Comment: Please include a Fishbase ID with the NCBI taxon ID, and also please use the more up-to-date BUSCO v3 rather than BUSCO v2. We also recommend adding RRIDs for the software tools, see:  <a href="ftp://penguin.genomics.cn/pub/10.5524/RRID/RRIDlist.pdf">ftp://penguin.genomics.cn/pub/10.5524/RRID/RRIDlist.pdf</a></p> <p>Reply: We have included the Fishbase ID and the NCBI taxon ID in the revised manuscript. For BUSCO version, we have double checked that the BUSCO v3 was used in our study. We are sorry for the mistake and corrected the description in our manuscript. We also have add the RRIDs of the software used in this work.</p> <p>Reviewer reports:</p> <p>Reviewer 1 comments:</p> <p>This manuscript describes a genome assembly for <i>Glyptosternon maculatum</i>, a Siluriform catfish endemic to the Tibetan plateau. The authors present sufficient data on sequence inputs and assembly methods, and also provide assembly data and analyses of assembly quality.</p> <p>The manuscript states that only kmers that occurred once were removed for the kmer-based estimation of genome size. Usually the kmer distribution from Illumina sequencing of vertebrate genomes is minimized somewhere beyond a single occurrence - perhaps 3 or 4 occurrences in the data (based on depth of coverage). If the authors are only removing single-occurrence kmers, they may be including untrusted kmers with low occurrence (two or three, perhaps?). This would lead to an overestimation of genome size. If so, then their assembly would be even closer to the estimated genome size.</p> <p>Reply: We thanks for the reviewer's comment. According to the suggestion, we removed all kmers with less than 3 occurrences and re-analyzed Kmer distribution data. Indeed, we observed that the estimated genome size was 763 Mb, slightly smaller than that (771 Mb) of our previous analysis. We have added the information into our manuscript, which was highlighted by red.</p>

How many iterations of Quiver and Pilon were performed? Current recommendations are to use Quiver to correct SNPs and indels in PacBio assemblies, then to use Pilon to only correct indels since short Illumina reads may be misaligned in repetitive regions.

Reply: The reviewer's question is quite important. We have performed one round of Quiver and Pilon correction using pacbio and NGS data, respectively. In Pilon correction process, because we have observed effects of indel correction using Pilon in our previous analysis (below figure); therefore both snp and indel were corrected in our analysis.

Please state whether the Illumina reads that were mapped with BWA were from the reference individual or another individual.

Reply: Thanks for the reminding. We used Illumina sequencing reads from the reference individual. We stated the detail in our revised manuscript.

Please provide discussion as to why some Trinity contigs only aligned at low coverages (75-85%).

Reply: Thanks for reviewer's suggestion. We have searched our mRNA sequencing reads to NT database and found that the top 5 hits were all from the closely related fish species, such as *Ictalurus punctatus* and Zebrafish. Therefore, the probability for external contamination was ruled out (SI Table 7).

We attribute the low coverage of some trinity contig to two fold reasons: 1) the potential chimeric transcript generated during the transcriptome assembly using trinity, especially for genes with various alternative splicing models; 2) the fragments of genomic contig sequences was also one reason for the low coverage alignment of some assembled transcripts. We have discussed the reason for the low coverage in our manuscript and the revision were highlighted by red.

There is confusion in the text when describing Figure 1b and Figure 1c. See first paragraph of Background information, First line of Sample collection and sequencing. Please clarify. Also, can the map be magnified to better locate the location of the reference sample?

Reply: Thanks for the reviewer's suggestion. Figure 1a and 1b described the *G. maculatum* and Figure 1c described the location of the sample collection. We have magnified the map in Figure 1c to Tibet-plateau according to the reviewer's suggestion.

Please provide more justification as to why the species in Figure 2 were chosen. Hopefully it is more than just because the data was available. If the purpose is to focus on the divergence between the two Siluriform catfish then are all the other species necessary?

Reply: The reviewer is correct. The species divergence analysis between *G. maculatum* and *I. punctatus* was the main purpose for phylogenetic analysis. The analysis could provide us useful information regarding to the species divergence time and relative relationships among fish species. For this purpose, we used other 12 fish genome to construct the phylogenetic tree, not only due to the availability of genome data of those evolutionarily close species, but also because more species (typically 10 or more species in previous studies) are needed to recalibrate the phylogenetic relationships and species divergence time.

Minor corrections:

Line 36 under Background Information - It is unclear what "causing by the unshift of Tibetan plateau" means. Please reword to clarify.

Reply: Thanks for the reviewer's suggestion. We have re-write this part to clarify the sentence. The revision were highlighted by red.

Line 52 under Background Information - Suggest "We thus chose *G. manulatum* to represent Glyptosternoid group fishes ..."

Reply: Thanks for the reviewer's suggestion. The revision were highlighted by red.

On Line 16 under Protein Coding Gene Prediction, *Sinocyclocheilus graham* should have the abbreviation Sga, and channel catfish should be listed as *Ictalurus punctatus* (Ipu).

Reply: Thanks for the reviewer's suggestion. The revision were highlighted by red.

On Line 7 under Functional annotation: "refers"

Reply: Thanks for the reviewer's reminding. We have corrected the sentence and the revision were highlighted by red in the manuscript.

On Line 20: "protein data"

Reply: We have corrected according to the reviewer's comment. The revision were highlighted by red in the manuscript.

Lines 43 and 44: "were" is used twice in the sentence and should only be used once

Reply: We have deleted "were" in the sentence. The revision were highlighted by red in the manuscript.

On Line 13 under Conclusion: "Glyptosternoids".

Reply: Thanks for reviewer's correction. The revision were highlighted by red in the manuscript.

Reviewer #2: This is a purely descriptive paper reporting the sequencing and genome annotation of *Glyptosternon maculatum*, an endemic catfish species from the Tibet plateau. This is a straightforward paper and a valuable resource, which deserves publication after minor revision.

- Was the fish individual sequenced a male or a female? Any hint of sex chromosomes in this species?

Reply: The fish individual used in this work was a female, and we have added the information in the manuscript. The sex determination and chromosome for the species were not identified so far, and no heterotropic chromosome was observed from the previous karyotype analysis (Wu Yunfei, Kang Bin, Men Qiang, et al. Chromosome diversity of tibetan fishes. *Zoological Research*, 1999, 20(4):258-264.) .

- It might be also of interest to predict long non-coding RNA genes (not presented in "all kinds" of non-coding RNA in Tab S5).

Reply: LncRNA is an important non-coding genes in gene expression regulation. However, the transcriptome used in this work were generated from the enrichment by oligo(dT), and it is not suitable for lncRNA prediction with a reference genome. Therefore, we did not annotate the lncRNA gene in this work. However, lncRNA regulation in high-altitude would be interesting direction, and related work will be performed in our following studies.

- To affirm that 228 genes are species-specific sounds always strange to me. More precise comparative analysis of their presence/absence in other (related) species should be performed to confirm this.

Reply: Thanks a lot for reviewer's suggestion. We further blast the 228 genes to NCBI NR (non-redundant) database, and found that 142 genes hit to database with e-value of  $1e-5$ ; however, there were still 86 genes failed to hit any protein sequences in the database. The function analysis of those genes is an interesting topic in our following studies.

We have corrected the term of "species-specific genes" to "genes without significant homologous hits" and added the additional analysis to the manuscript. The correction were highlighted by red.

- Repeats: did the authors check for the presence of MITE elements?

Reply: Thanks a lot for reviewer's suggestion. We checked for the presence of MITE elements using MITE-digger with default parameters. The results shown that there were 2,962 MITEs which account for 0.185 % of the whole genome. The resulted MITE sequences were attached as Supplemental\_file\_MITE to the revised manuscript.

- Species names should be italicized in Fig.2

Reply: We thanks reviewer for the reminding. Species names were italicized in Fig. 2.

- The paper should be edited for typos/grammatical errors

Reply: Thanks a lot for reviewer's suggestion. We have revised and corrected typos and grammatical errors through the manuscript. The corrections were highlighted by red.

	- Tab S3: satellites are not interspersed repeat (correct title) Reply: Thanks reviewer for the reminding. We have corrected the title in Table S3.
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>  Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . 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?	Yes
<b>Resources</b>  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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.  Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> ?	Yes
<b>Availability of data and materials</b>  All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials”	Yes

section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

[Click here to view linked References](#)

1 **Draft genome of *Glyptosternon maculatum*, an endemic**  
2 **fish from Tibet-plateau**

3  
4 Haiping Liu<sup>1,†</sup>, Qiyong Liu<sup>1,†</sup>, Zhiqiang Chen<sup>2,†</sup>, Yanchao Liu<sup>1,†</sup>, Chaowei Zhou<sup>1,†</sup>, Qiqi  
5 Liang<sup>2</sup>, Caixia Ma<sup>2</sup>, Jianshe Zhou<sup>1</sup>, Yingzi Pan<sup>1</sup>, Meiqun Chen<sup>1</sup>, Wangjiu<sup>1</sup>, Wenkai Jiang<sup>2,\*</sup>,  
6 Shijun Xiao<sup>3,\*</sup>, Zhenbo Mou<sup>1,\*</sup>

7  
8 <sup>1</sup> Institute of Fisheries Science, Tibet Academy of Agricultural and Animal Husbandry  
9 Sciences, Lhasa 850002, China

10 <sup>2</sup> Novogene Bioinformatics Institute, Beijing, China

11 <sup>3</sup> Department of Computer Science, Wuhan University of Technology, Wuhan, China

12 \* Correspondence: Zhenbo Mou ([mouzhenbo@163.com](mailto:mouzhenbo@163.com)), Wenkai Jiang  
13 ([jiangwenkai@novogene.com](mailto:jiangwenkai@novogene.com)) and Shijun Xiao([shijun\\_xiao@163.com](mailto:shijun_xiao@163.com))

14 † These authors contributed equally to this work.

## Abstracts

**Background:** Mechanisms for high altitude adaption have arisen widespread interest to evolution biologists. Several genome wide studies have been carried out for endemic vertebrates in Tibet, including mammals, birds and amphibians. However, little information was known about the adaptive evolution of highland fishes. *Glyptosternon maculatum* (*G. maculatum*, Regan, 1905), also known as Regan or barkley, is a fish endemic to the Tibetan plateau, which belongs to Sisoridae family, Siluriformes (catfishes) order. This species live within an elevation ranging from roughly 2800 m to 4200 m. Hence, a high-quality reference genome of *G. maculatum* provides an opportunity to address high altitude adaption mechanisms of fishes.

**Findings:** To get a high-quality reference genome of *G. maculatum*, we combined PacBio single-molecule real-time sequencing, Illumina paired-end sequencing, 10X Genomics linked-reads and BioNano optical map techniques. In total, 603.99 Gb sequencing data were generated. The assembled genome was about 662.34 Mb with scaffold and contig N50 sizes of 20.90 Mb and 993.67 kb, respectively, which captured 83% complete and 3.9% partial vertebrate Benchmarking Universal Single-copy orthologs (BUSCO). Repetitive elements account for 35.88% of the genome, and 22,066 protein-coding genes were predicted from the genome, of which 91.7% have been functionally annotated.

**Conclusions:** We provide the first comprehensive *de novo* genome of the *G. maculatum*. This genetic resource is fundamental for investigating the origin of the *G. maculatum* and will improve our understanding of high altitude adaption of fishes. The assembled genome can also be used as reference for future population genetic studies of *G. maculatum*.

**Keywords:** *Glyptosternon maculatum*; Genome assembly; Annotation; Phylogeny



## Data description

### Background information on *G. maculatum*

The *G. maculatum* (Regan, 1905) (Fishbase ID: 24838, NCBI Taxon ID: 175778, [Figure 1a, 1b](#)), called barkley in Tibetan language, is a species in the *Glyptosternum* genus, Sisoridae family, Siluriformes order, Teleostei infraclass. The Sisoridae are the largest family of catfishes (Siluriformes) in China, consisting of 44 species divided into two natural groups, Glyptosternoids and non-Glyptosternoids [1, 2]. There are 8 Sisorids distributed in Yarlung Tsangpo (Brahmaputra) river. Of them, *G. maculatum* is the only one species that distributed at the middle section. Specifically, it is distributed at Niyang tributary, Tangjia to Zhaxue segment of the Lhasa tributary and Xietongmen segment of Yarlung Tsangpo, which across an elevation ranging from roughly 2800 m to 4200 m [3].

The karyotype of *G. maculatum* is a debated topic. Ren and Cui [4] reported a result of  $2n=48=28m+12sm+8st$ ,  $NF=88$ , based on specimens sampled at Quxur, and speculated it to be the most specialized karyotype among all sisoridae. Conversely, Wu et al. [5] reported a karyotype of  $2n=48=22m+12sm+10st+6t$ ,  $NF=80$  sampled at Xigaze, while  $2n=44$  and  $2n=42$  results were also found. They compared it to other Sisorid karyotypes and concluded that the karyotype of *G. maculatum* was not the most specialized. The genome assembly of *G. maculatum* might provide a route to resolve these debates.

Glyptosternoid group fishes distributed broadly at the south and southeast drainages of Tibetan plateau, providing a good model to study the speciation process causing by the up-shift of Tibetan plateau. He *et al.*[1] conducted a cladistical analysis of Glyptosternoid group based on 60 bone features, and found Glyptosternoids formed a monophyletic group, of which the *Glyptosternum* were the most primitive clade. He *et al.*[6] further analyzed the phylogeny of Glyptosternoids using 19 species distributed in four genera by their bone features, in combination with biogeographical analysis, they postulated the rise of the Tibetan Plateau had a direct influence on the diversification of Glyptosternoids, with *Glyptosternum* (particularly *G. maculatum*) as the most primitive clade, which was consistent to the conclusion of Hora and Silas [2]. Peng *et al.*[7] sequenced mitochondrial cytochrome b (CYTB) from 13 Glyptosternoids, with results also supporting them to be a monophyletic group, of which *Glyptosternum* and *Exostoma* were relatively primitive clades. We thus chose *G. maculatum* to represent Glyptosternoid group fishes, and its whole genome sequence would provide a foundation to explore the adaptive evolution process of highland fishes, also supplied as a starting point to study speciation mechanisms caused by rapid rising of Tibetan plateau.

*G. maculatum* had a specialized liver, which could be divided to two parts, one placed outside the abdominal cavity, connected to another part that located inside the cavity [8].

102 (Figure 1b) Several studies had also reported similar ectopic livers exist in other Sisords,  
103 suggesting that this specialized organ might be the result of adaptive evolution [9]. The  
104 genesis of livers in *G. maculatum* occurs in three stages: the ectopic liver is not present  
105 from the beginning till the end of the larva's exit from the egg envelope; a "bump" then  
106 develops, starting from day 17 till day 22; the ectopic liver appears starting from day 22 [9].  
107 Zhang [9] pointed out the expressions of Cu-Zn SOD, Mn SOD, and CAT mRNA were all  
108 higher in the primary liver relative to the secondary liver, suggesting that the two livers  
109 have different physiological roles in *G. maculatum*. However, the molecular mechanism  
110 for the liver development and their physiological functions in adaptive evolution were not  
111 fully understand; therefore genome assembly of the species could lay a solid foundation  
112 for the following investigations.

### 113 **Sample collection and sequencing**

114 The female fish individual used for genome sequencing came from Angren, Xizang  
115 Province (Figure 1c). Total genomic DNA was extracted from muscular tissue and kept at  
116 Novogene Bioinformatics Institute.

117 A combination of four technologies was applied: PacBio's single-molecule real-time  
118 sequencing, Illumina's paired-end sequencing, 10X Genomics link-reads and BioNano  
119 optical maps. Two paired-end Illumina sequence libraries were constructed with an insert  
120 size of 250 bp, and sequencing was carried out on the Illumina HiSeq 4000 platform  
121 according to the manufacturer's instructions, of which 147.16 Gb (191x coverage)  
122 sequencing data were produced. In addition, one 10X Genomics linked-read library was  
123 constructed and sequencing on Illumina HiSeq 4000 platform, which produced 157 Gb  
124 (203.5x coverage) sequencing data. Raw sequence data generated by Illumina platform  
125 were filtered by these criteria: (a) filtered reads with adapters; (b) filtered reads with N  
126 bases more than 10%; (c) filtered reads with low-quality bases ( $\leq 5$ ) more than 50%.  
127 Pacbio reads were sequenced by the Sequel platform, which gained 106.3 Gb (145.2x  
128 coverage) sequencing data. For the PacBio data, subreads were filtered with the default  
129 parameters. Finally, we obtained 106.32 Gb of long reads (polymerase reads) data. The  
130 average and the N50 length of long subreads reached 8.04 kb and 13.26 kb, respectively.  
131 An optical map was also constructed from Irys platform (BioNano Genomics), of which  
132 191.3 Gb (248x coverage) data were generated. All these sequence data were  
133 summarized in Table 1.

### 134 **De novo assembly of *G. maculatum* genome**

135 The genome size was estimated based on the  $k$ -mer spectrum :  $G = (K_{total} - K_{error})/D$ ,  
136 where  $K_{total}$  is the total count of  $k$ -mers,  $K_{error}$  is the total count of low-frequency (frequency  
137  $\leq 3$ )  $k$ -mers that were probably caused by sequencing errors,  $G$  is the genome size and  $D$   
138 is the  $k$ -mer depth. Using Jellyfish [10] (v2.1.3), 17-mers were counted as 54,676,846,244

139 from short clean reads. The total count of error *k*-mers was 1,980,028,579 and the *k*-mer  
140 depth was 69 (Figure. S1). Therefore, the genome size of *G. maculatum* was estimated to  
141 be approximately 763.7 Mb.

142 The contig assembly of the *G. maculatum* genome was carried out using the  
143 FALCON assembler [11], followed by two rounds of polishing with Quiver [12]. FALCON  
144 implements a hierarchical assembly process, which include these steps: 1) subread error  
145 correction through aligning all reads to each other using daligner [13], the overlap data  
146 were then processed to generate error-corrected consensus reads; After error correction,  
147 we got 28 Gb (35x coverage) of error-corrected reads; 2) second round of overlap  
148 detection using error-corrected reads; 3) construction of a directed string graph from  
149 overlap data; 4) resolving contig path from the string graph. After FALCON assembly, the  
150 genome was polished by Quiver. Initial assembly of the PacBio data alone resulted in a  
151 contig N50 (the minimum length of contigs accounting for half of the haploid genome size)  
152 of 697.4 Kb. Then PacBio contigs were first scaffolded using optical map data, and the  
153 resulting scaffolds were further connected to super-scaffolds by 10X Genomics  
154 linked-read data using the fragScaff software [14]. Finally, we used Illumina-derived short  
155 reads to correct any remaining errors by pilon [15]. These processes yielded a final draft *G.*  
156 *maculatum* genome assembly with a total length of 662.34 Mb, contig N50 of 993.67 kb,  
157 and scaffold N50 of 20.90 Mb (Table 2).

158 To evaluate the accuracy of the genome at single base level, we mapped short  
159 sequence reads generated by Illumina platform to the *G. maculatum* genome with BWA  
160 (RRID:SCR\_010910) [16] and performed variant calling with SAMtools  
161 (RRID:SCR\_002105) [17]. We obtained a total of 3,632 homozygous SNPs (Table S2),  
162 reflecting a low homozygous rate (0.0007%) and a high accuracy of genome assembly at  
163 the single base level.

164 To assess the completeness of the assembled *G. maculatum* genome, we performed  
165 BUSCO (RRID:SCR\_015008) analysis [18] by searching against the vertebrate universal  
166 benchmarking single-copy orthologs (BUSCOs, version 3.0). Overall, 83% complete and  
167 3.9% partial of the 970 vertebrate BUSCOs were identified in the assembled genome. We  
168 also assessed the completeness of *G. maculatum* genome by CEGMA (Core Eukaryotic  
169 Genes Mapping Approach, RRID:SCR\_015055) [19]. According to CEGMA, 211 (85.08%)  
170 conserved genes were identified in the *G. maculatum* genome.

171 The muscle transcriptome *de novo* assembled by Trinity (RRID:SCR\_013048) [20]  
172 were also mapped to the genome assembly using BLAT [21] with default parameters,  
173 showing that the alignment coverage of expressed sequences ranged from 75 to 99% in  
174 the genome assembly. To answer the question why some contig has a low coverage (85%)  
175 on genome sequence alignment. We first searched mRNA sequencing reads to NT  
176 database and found that the top 5 hits were all from the closely related fish species, such

177 as *Ictalurus punctatus* and *Danio rerio* (SI Table 7). Therefore, the probability for external  
178 contamination was ruled out. We therefore attributed the low coverage of some trinity  
179 contig to two fold reasons: 1) the potential chimeric transcript generated during the  
180 transcriptome assembly using trinity, especially for genes with various alternative splices;  
181 2) the fragments of genomic contig sequences was also one reason for the low coverage  
182 alignment of some assembled transcripts.

### 183 **Annotation of repetitive sequences in *G. maculatum* genome**

184 The repetitive sequences in *G. maculatum* genome were identified by a combination of  
185 homology searching and *ab initio* prediction. For homology-based prediction, we used  
186 RepeatMasker (RRID:SCR\_012954) [22] and RepeatProteinMask to search against  
187 Repbase. For *abinitio* prediction, we used Tandem Repeats Finder (TRF) [23],  
188 LTR\_FINDER (RRID:SCR\_015247) [24], PILER [25] and RepeatScout  
189 (RRID:SCR\_014653) [26] with default parameters. We found that 33.96% of the *G.*  
190 *maculatum* assembly was composed of repetitive elements (Table S3 and Figure S2).  
191 Additionally, we predicted MITE elements through the genome using MITE-digger [27]  
192 with default parameters. As a result, we identified 2,962 MITEs accounting for 0.185% of  
193 the whole genome (Supplemental\_file\_MITE).

### 194 **Protein coding gene prediction and ncRNA prediction**

195 Gene prediction was conducted through a combination of homology-based prediction, *ab*  
196 *initio* prediction and transcriptome-based prediction methods. Protein repertoires of  
197 vertebrates including *Takifugu rubripes* (Tru, GCF\_000180615.1), *Ctenopharyngodon*  
198 *idellus* (Cid) [28], *Cyprinus carpio* (Cca, GCF\_000951615.1), *Danio rerio* (Dre,  
199 GCF\_000002035.5), *Sinocyclocheilus graham* (Sga, GCF\_001515645.1), channel catfish  
200 (Ipu, GCF\_001660625.1), *Homo sapiens* (Hom, GCF\_000001405.37) and *Mus musculus*  
201 (Mmu, GCF\_000001635.26) were used as queries to search against *G. maculatum*  
202 genome using TBLASTN (RRID:SCR\_011822) [29]. The BLAST hits were conjoined by  
203 Solar software [30]. GeneWise (RRID:SCR\_015054) [31] was used to predict the exact  
204 gene structure of the corresponding genomic region on each BLAST hit. Homology  
205 predictions were denoted as “Homology-set” (Table S4). RNA-seq data derived from 10  
206 tissues which obtained about 77.29 Gb clean data were assembled by Trinity [20]. The  
207 Trinity assembly included 572,416 contigs with an average length of 1,075 bp. These  
208 assembled sequences were aligned against the *G. maculatum* genome by PASA  
209 (Program to Assemble Spliced Alignment). Valid transcript alignments were clustered  
210 based on genome mapping location and assembled into gene structures. Gene models  
211 created by PASA [32] were denoted as PASA-T-set (PASA Trinity set). Besides, RNA-seq  
212 reads were directly mapped to the genome using Tophat (RRID:SCR\_013035) [33] to  
213 identify putative exon regions and splice junctions; Cufflinks (RRID:SCR\_014597) [34]  
214 was then used to assemble the mapped reads into gene models (Cufflinks-set). Augustus

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

215 (RRID:SCR\_008417) [35], GeneID [36], GeneScan [37], GlimmerHMM  
216 (RRID:SCR\_002654) [38], and SNAP [39] were also used to predict coding regions in the  
217 repeat-masked genome. Of these, Augustus, SNAP and GlimmerHMM were trained by  
218 PASA-H-set gene models. Gene models generated from all the methods were integrated  
219 by EvidenceModeler (EVM) [40]. Weights for each type of evidence were set as follows:  
220 PASA-T-set > Homology-set > Cufflinks set > Augustus > GeneID = SNAP =  
221 GlimmerHMM = GeneScan. The gene models were further updated by PASA2 to  
222 generate UTRs, alternative splicing variation information. In total, we have identified  
223 22,066 protein coding genes with a mean of 8.5 exons per gene (Table 3). The lengths of  
224 genes, coding sequence (CDS), introns, and exons in *G. maculatum* were  
225 comparable to those of close-related genomes (Table S4 and Figure S3). In addition, we  
226 predicted non-coding RNA genes in the *G. maculatum* genome. The rRNA fragments were  
227 predicted by searching against Human rRNA database using BLAST with an E-value of  
228 1E-10. The tRNA genes were identified by tRNAscan-SE (RRID:SCR\_010835) software  
229 [41]. The miRNA and snRNA genes were predicted by INFERNAL (RRID:SCR\_011809)  
230 [42] using Rfam database [43]. We found a total of 3,117 ribosomal RNA (rRNA), 3,512  
231 transfer RNA (tRNA), 1,235 microRNAs (miRNA), and 781 snRNA genes in the *G.*  
232 *maculatum* genome (Table S5).

### 233 **Functional annotation of protein-coding genes**

234 Gene function of predict protein-coding gene were annotated by searching functional  
235 motifs, domains, and possible biological process of genes to known databases such as  
236 SwissProt [44], Pfam [45], NR database (from NCBI), GeneOntology (GO) [46], and Kyoto  
237 Encyclopedia of Genes and Genomes (KEGG) [47]. In total, 20,234 protein-coding genes  
238 (91.7%) were successfully annotated for at least one function terms (Table S6, Figure S4).

### 239 **Phylogenetic analysis and species divergence time estimation**

240 To investigate the phylogenetic position of *G. maculatum*, we retrieved nucleotide and  
241 protein data for *Cyprinus carpio* (GCF\_000951615.1), *Sinocyclocheilus rhinoceros*  
242 (GCF\_001515625.1), *Sinocyclocheilus anshuiensis* (GCF\_001515605.1), *Astyanax*  
243 *mexicanus* (GCF\_000372685.2), *Pygocentrus nattereri* (GCF\_001682695.1),  
244 *Sinocyclocheilus grahami* (GCF\_001515645.1), *Ictalurus punctatus* (GCF\_001660625.1),  
245 *Danio rerio* (GCF\_000002035.6), *Amazon molly* (GCF\_000485575.1), *Oreochromis*  
246 *niloticus* (GCF\_001858045.1), *Takifugu rubripes* (GCF\_000180615.1) and  
247 *Ctenopharyngodon idellus* [28] from public databases. To remove redundancy caused by  
248 alternative splicing variations, we retained only gene models at each gene locus that  
249 encoded the longest protein sequence. To exclude putative fragmented genes, genes  
250 encoding protein sequences shorter than 50 amino acids were filtered out. All-against-all  
251 BLASTP (RRID:SCR\_001010) [29] was employed to identify the similarities among  
252 filtered protein sequences in these species with an E-value cut off of 1e<sup>-7</sup>. The OrthoMCL

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

253 (RRID:SCR\_007839) [48] method was used to cluster genes from these different species  
254 into gene families with the parameter of “-inflation 1.5”.

255 A total of 26,588 gene family clusters were constructed. There were 101 gene families  
256 and 228 genes in *G. maculatum* without significant homologous hits to other teleosts. We  
257 further searched the 228 genes to NCBI NR database by BLASTP (RRID:SCR\_001010)  
258 and found that 142 genes hit to database with e-value of 1e-5, and 86 genes still failed to  
259 hit any protein sequences in the database. The function of those genes lacking significant  
260 homology is an interesting topic in the following studies.

261 Protein sequences from 247 single copy gene families were used for phylogenetic  
262 tree reconstruction. MUSCLE (RRID:SCR\_011812) [49] was used to generate multiple  
263 sequence alignment for protein sequences in each single-copy family with default  
264 parameters. Then, the alignments of each family were concatenated to a super alignment  
265 matrix. The super alignment matrix was used for phylogenetic tree reconstruction through  
266 Maximum likelihood (ML) methods. Divergence time between species was estimated  
267 using MCMCtree in PAML [50] with the options ‘correlated molecular clock’ and ‘JC69’  
268 model. A Markov Chain Monte Carlo analysis was run for 20,000 generations, using a  
269 burn-in of 1000 iterations. Divergence time for the common ancestor of *C. idellus*, *S.*  
270 *rhinocerosus* and *P. nattereri* obtained from the TimeTree database  
271 (<http://www.timetree.org/>) was used as the calibrate point. These phylogenetic analyses  
272 indicated that *G. maculatum* diverged from the common ancestral of *I. punctatus* at  
273 approximately 48.3million years ago (Figure 2).

## 274 Conclusion

275 We have constructed a *de novo* assembly of the *G. maculatum* genome and describe its  
276 genetic attributes. To our knowledge, this is the first *de novo* genome for Glyptosternoids  
277 group fishes. The *G. maculatum* genome will support investigations concerning the origin  
278 and evolutionary history of Glyptosternoid. This resource was also important for the future  
279 conservation of this endangered plateau species. In addition, the *G. maculatum* genome  
280 laid a solid foundation to investigate molecular mechanism of high altitude adaption of  
281 fishes and the speciation process during the rising of Tibetan plateau.

## 282 Availability of supporting data

283 The raw sequencing and physical mapping data were deposited into The National  
284 Omics Data Encyclopedia (NODE) (<http://www.biosino.org/node/index>) with the project ID  
285 of OEP000007 (<http://www.biosino.org/node/project/detail/OEP000007>) and SRA  
286 accession of SRR7279473-SRR7279474, SRR7268130-SRR7268162,  
287 SRR7350914-SRR7350921, SRR7351269-SRR7351265, SRR7403445-SRR7403454  
288 under the Project accession number of PRJNA447978 in NCBI. The genome, annotation

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 and intermediate files were uploaded to GigaScience FTP server. All supplementary  
290 figures and tables are provided in Supplemental File.

## 291 **Competing interests**

292 All authors declare that they have no competing interests.

## 293 **Authors' contributions**

294 Haiping Liu, Wenkai Jiang and Zhenbo Mou conceived the study. Haiping Liu and Wenkai  
295 Jiang designed the scientific objectives. Qiyong Liu and Zhenbo Mou managed the project;  
296 Yanchao Liu and Chaowei Zhou collected the samples and extracted the genomic DNA;  
297 Zhiqiang Chen estimated the genome size and assembled the genome; Qiqi Liang and  
298 Caixia Ma assessed the assembly quality; Jianshe Zhou and Yingzi Pan carried out the  
299 repeat annotation and gene annotation. Zhiqiang Chen carried out comparative genomics  
300 analysis, Haiping Liu, Shijun Xiao, Zhiqiang Chen and Wenkai Jiang wrote the manuscript.  
301 And all authors read, edited, and approved the final manuscript.

## 302 **Acknowledgements**

303 This work was supported by the special finance of Tibet autonomous region, Grant No.  
304 2017CZZX003 and the National Natural Science Foundation of China (NSFC), Grant No.  
305 31560144 and 31602207. Thanks should go to anonymous reviewers for their helpful  
306 comments and constructive suggestions.

## 308 **References**

- 309 1. He, S., The phylogeny of the glyptosternoid fishes (Teleostei: Siluriformes, Sisoridae). *Cybium*,  
310 1996. **20**(2): p. 115-159.
- 311 2. Hora, S.L. and E.G. Silas, Evolution and distribution of Glyptosternoid fishes of the family  
312 Sisoridae (Order: Siluroidea). *Proceedings of the National Institute of Sciences of India*, 1952. **18**.
- 313 3. T, C.Q. and Z.B. S, *Systematic Index of Fish Species in China*. China Science Publishing, 1987.
- 314 4. Ren, X., The karyotype and haploidy NOR of *Glyptosternum maculatum*. *Hereditas*, 1992.
- 315 5. WU Yun-fei, K.B., MEN Qiang, WU Cui-zhen, Chromosome Diversity of Tibetan Fishes.  
316 *Zoological Research*, 1999. **20**(4): p. 258-264.
- 317 6. He, S., W. Cao, and Y. Chen, The uplift of Qinghai-Xizang (Tibet) Plateau and the vicariance  
318 speciation of glyptosternoid fishes (Siluriformes: Sisoridae). *Science China* 2001. **44**(6): p. 644.
- 319 7. Peng, Z., S. He, and Y. Zhang, Phylogenetic relationships of glyptosternoid fishes (Siluriformes:  
320 Sisoridae) inferred from mitochondrial cytochrome b gene sequences. *Molecular Phylogenetics &*  
321 *Evolution*, 2004. **31**(3): p. 979-987.
- 322 8. X, X.C., L.H. J, and L.D. P. *Special Organ in Glyptosternum maculatum:exo-celiac liver*. 2006.
- 323 9. Huijuan, Z., *Genesis of Liver in Glyptosternum maculatum and Related Bioadaptive Studies*, in  
324 *Library of Huazhong Agricultural University*. 2011, Huazhong Agricultural University: WuHan,  
325 China.
- 326 10. Marçais, G. and C. Kingsford, A fast, lock-free approach for efficient parallel counting of

- 327 occurrences of k-mers. *Bioinformatics*, 2011. **27**(6): p. 764-70.
- 1 328 11. Pendleton, M., R. Sebra, A.W. Pang, et al., Assembly and diploid architecture of an individual  
2 329 human genome via single-molecule technologies. *Nat Methods*, 2015. **12**(8): p. 780-6.
- 3 330 12. Chin, C.S., D.H. Alexander, P. Marks, et al., Nonhybrid, finished microbial genome assemblies  
4 331 from long-read SMRT sequencing data. *Nat Methods*, 2013. **10**(6): p. 563-9.
- 5 332 13. Myers, G. *Efficient Local Alignment Discovery amongst Noisy Long Reads*. in *Algorithms in*  
6 333 *Bioinformatics*. 2014. Berlin, Heidelberg: Springer Berlin Heidelberg.
- 7 334 14. Adey, A., J.O. Kitzman, J.N. Burton, et al., In vitro, long-range sequence information for de novo  
8 335 genome assembly via transposase contiguity. *Genome Research*, 2014. **24**(12): p. 2041-9.
- 9 336 15. Walker, B.J., T. Abeel, T. Shea, et al., Pilon: an integrated tool for comprehensive microbial variant  
10 337 detection and genome assembly improvement. *PLoS One*, 2014. **9**(11): p. e112963.
- 11 338 16. Li, H., *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*. Vol. 1303.  
12 339 2013.
- 13 340 17. Li, H., A statistical framework for SNP calling, mutation discovery, association mapping and  
14 341 population genetical parameter estimation from sequencing data. *Bioinformatics*, 2011. **27**(21): p.  
15 342 2987-93.
- 16 343 18. Simao, F.A., R.M. Waterhouse, P. Ioannidis, et al., BUSCO: assessing genome assembly and  
17 344 annotation completeness with single-copy orthologs. *Bioinformatics*, 2015. **31**(19): p. 3210-2.
- 18 345 19. Parra, G., K. Bradnam, and I. Korf, CEGMA: a pipeline to accurately annotate core genes in  
19 346 eukaryotic genomes. *Bioinformatics*, 2007. **23**(9): p. 1061-1067.
- 20 347 20. Grabherr, M.G., B.J. Haas, M. Yassour, et al., Full-length transcriptome assembly from RNA-Seq  
21 348 data without a reference genome. *Nat Biotechnol*, 2011. **29**(7): p. 644-52.
- 22 349 21. Kent, W.J., BLAT--the BLAST-like alignment tool. *Genome Res*, 2002. **12**(4): p. 656-64.
- 23 350 22. Bergman, C.M. and H. Quesneville, Discovering and detecting transposable elements in genome  
24 351 sequences. *Brief Bioinform*, 2007. **8**(6): p. 382-92.
- 25 352 23. Benson, G., Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*, 1999.  
26 353 **27**(2): p. 573-80.
- 27 354 24. Xu, Z. and H. Wang, LTR\_FINDER: an efficient tool for the prediction of full-length LTR  
28 355 retrotransposons. *Nucleic Acids Res*, 2007. **35**(Web Server issue): p. W265-8.
- 29 356 25. Edgar, R.C. and E.W. Myers, PILER: identification and classification of genomic repeats.  
30 357 *Bioinformatics*, 2005. **21 Suppl 1**: p. i152-8.
- 31 358 26. Price, A.L., N.C. Jones, and P.A. Pevzner, De novo identification of repeat families in large  
32 359 genomes. *Bioinformatics*, 2005. **21 Suppl 1**: p. i351-8.
- 33 360 27. Yang, G., MITE Digger, an efficient and accurate algorithm for genome wide discovery of miniature  
34 361 inverted repeat transposable elements. *BMC Bioinformatics*, 2013. **14**: p. 186.
- 35 362 28. Wang, Y., Y. Lu, Y. Zhang, et al., The draft genome of the grass carp (*Ctenopharyngodon idellus*)  
36 363 provides insights into its evolution and vegetarian adaptation. *Nat Genet*, 2015. **47**(6): p. 625-31.
- 37 364 29. Altschul, S.F., W. Gish, W. Miller, et al., Basic local alignment search tool. *J Mol Biol*, 1990. **215**(3):  
38 365 p. 403-10.
- 39 366 30. Yu, X.J., H.K. Zheng, J. Wang, et al., Detecting lineage-specific adaptive evolution of  
40 367 brain-expressed genes in human using rhesus macaque as outgroup. *Genomics*, 2006. **88**(6): p.  
41 368 745-751.
- 42 369 31. Birney, E., M. Clamp, and R. Durbin, GeneWise and Genomewise. *Genome Res*, 2004. **14**(5): p.  
43 370 988-95.
- 44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
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
- 371 32. Haas, B.J., A.L. Delcher, S.M. Mount, et al., Improving the Arabidopsis genome annotation using  
372 maximal transcript alignment assemblies. *Nucleic Acids Res*, 2003. **31**(19): p. 5654-66.
  - 373 33. Kim, D., G. Pertea, C. Trapnell, et al., TopHat2: accurate alignment of transcriptomes in the  
374 presence of insertions, deletions and gene fusions. *Genome Biol*, 2013. **14**(4): p. R36.
  - 375 34. Trapnell, C., A. Roberts, L. Goff, et al., Differential gene and transcript expression analysis of  
376 RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc*, 2012. **7**(3): p. 562-78.
  - 377 35. Stanke, M. and S. Waack, Gene prediction with a hidden Markov model and a new intron submodel.  
378 *Bioinformatics*, 2003. **19 Suppl 2**: p. ii215-25.
  - 379 36. Guigo, R., Assembling genes from predicted exons in linear time with dynamic programming. *J*  
380 *Comput Biol*, 1998. **5**(4): p. 681-702.
  - 381 37. Burge, C. and S. Karlin, Prediction of complete gene structures in human genomic DNA. *J Mol Biol*,  
382 1997. **268**(1): p. 78-94.
  - 383 38. Majoros, W.H., M. Pertea, and S.L. Salzberg, TigrScan and GlimmerHMM: two open source ab  
384 initio eukaryotic gene-finders. *Bioinformatics*, 2004. **20**(16): p. 2878-9.
  - 385 39. Korf, I., Gene finding in novel genomes. *BMC Bioinformatics*, 2004. **5**: p. 59.
  - 386 40. Haas, B.J., S.L. Salzberg, W. Zhu, et al., Automated eukaryotic gene structure annotation using  
387 EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome Biol*, 2008. **9**(1): p.  
388 R7.
  - 389 41. Lowe, T.M. and S.R. Eddy, tRNAscan-SE: a program for improved detection of transfer RNA genes  
390 in genomic sequence. *Nucleic Acids Res*, 1997. **25**(5): p. 955-64.
  - 391 42. Nawrocki, E.P., D.L. Kolbe, and S.R. Eddy, Infernal 1.0: inference of RNA alignments.  
392 *Bioinformatics*, 2009. **25**(10): p. 1335-7.
  - 393 43. Li, Y.-h., G. Zhou, J. Ma, et al., De novo assembly of soybean wild relatives for pan-genome  
394 analysis of diversity and agronomic traits. *Nature Biotechnology*, 2014. **32**: p. 1045.
  - 395 44. UniProt: the universal protein knowledgebase. *Nucleic Acids Research*, 2017. **45**(D1): p.  
396 D158-D169.
  - 397 45. Finn, R.D., P. Coghill, R.Y. Eberhardt, et al., The Pfam protein families database: towards a more  
398 sustainable future. *Nucleic Acids Res*, 2016. **44**(D1): p. D279-85.
  - 399 46. Expansion of the Gene Ontology knowledgebase and resources. *Nucleic Acids Res*, 2017. **45**(D1): p.  
400 D331-d338.
  - 401 47. Kanehisa, M., S. Goto, Y. Sato, et al., Data, information, knowledge and principle: back to  
402 metabolism in KEGG. *Nucleic Acids Res*, 2014. **42**(Database issue): p. D199-205.
  - 403 48. Li, L., C.J. Stoeckert, Jr., and D.S. Roos, OrthoMCL: identification of ortholog groups for  
404 eukaryotic genomes. *Genome Res*, 2003. **13**(9): p. 2178-89.
  - 405 49. Edgar, R.C., MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
406 *Nucleic Acids Res*, 2004. **32**(5): p. 1792-7.
  - 407 50. Yang, Z., PAML: a program package for phylogenetic analysis by maximum likelihood. *Comput*  
408 *Appl Biosci*, 1997. **13**(5): p. 555-6.

410 **Table 1. Sequencing data used for the *G. maculatum* genome assembly. The**  
 411 **coverage was calculated using an estimated genome size of 771.19 Mb.**

Pair-end libraries	Insert size (bp)	Raw data (Gb)	Clean data (Gb)	Read length (bp)	Sequence coverage (X)
<b>Illumina reads</b>	250 bp	148.16	147.16	150	191
<b>Pacbio reads</b>	20 Kb	106.32	106.05	11,745	145.2
<b>10X Genomics</b>	500 bp	157.21	157.02	150	203.5
<b>BioNano</b>	--	192.30	191.30	--	248
<b>Total</b>	--	603.99	601.53	--	787.7

413 [Back to main body](#)

414  
 415 **Table 2. Assembly statistics of *G. maculatum***

Sample ID	Length		Number	
	Contig** (bp)	Scaffold (bp)	Contig**	Scaffold
<b>Total</b>	637,133,884	662,339,741	3,281	531
<b>Max</b>	5,772,991	47,179,384	-	-
<b>Number&gt;=2000</b>	-	-	3,161	531
<b>N50</b>	993,673	20,902,354	161	11
<b>N60</b>	668,112	17,328,106	239	14
<b>N70</b>	418,057	12,288,896	359	19
<b>N80</b>	211,596	6,320,921	575	27
<b>N90</b>	77,392	1,017,220	1,067	50

416 \*\* Contig after scaffolding

417 [Back to main body](#)

418  
 419

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

420 **Table 3. General statistics of predicted protein-coding genes**

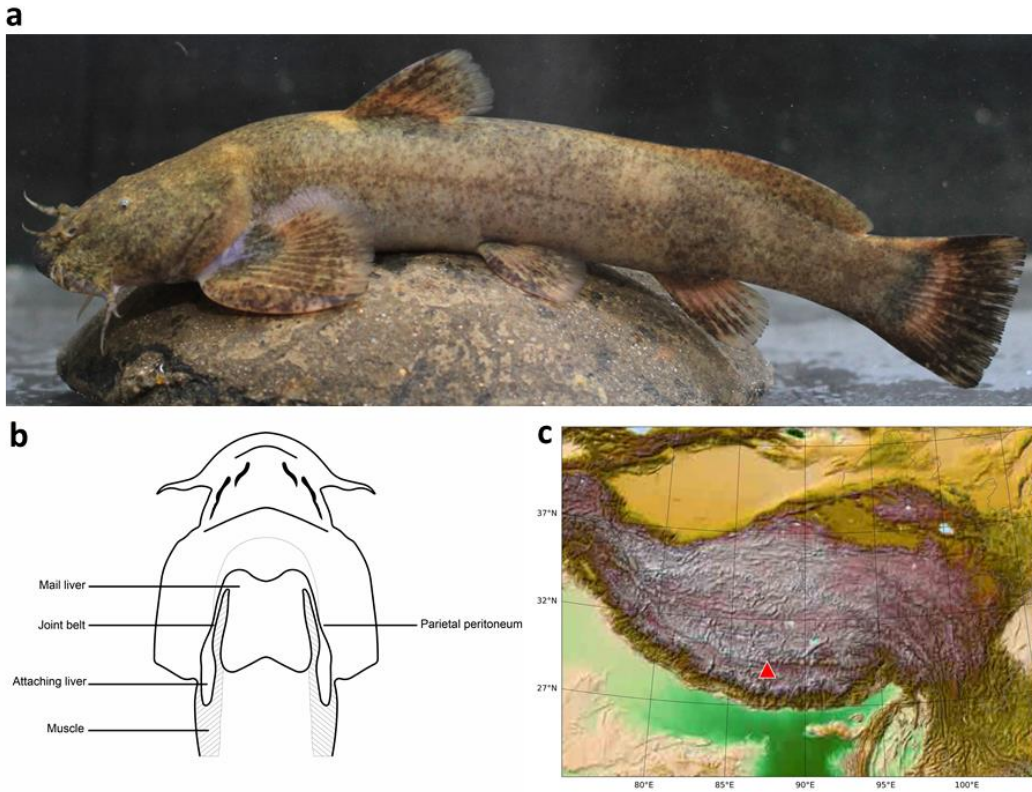
	Gene set	Number	Average transcript length (bp)	Average CDS length (bp)	Average exons per gene	Average exon length (bp)	Average intron length (bp)
	<b>Augustus</b>	14,910	9,534	1,241	6.93	179	1,399
	<b>GlimmerHMM</b>	73,320	7,896	574	3.87	148	2,551
<b>De novo</b>	<b>SNAP</b>	43,247	15,950	847	6.04	140	2,996
	<b>Geneid</b>	23,523	16,924	1,323	6.29	210	2,948
	<b>Genscan</b>	24,037	19,024	1,514	8.14	186	2,451
	<b>Sga</b>	32,364	6,413	1,142	5.12	223	1,279
	<b>Cca</b>	27,208	6,326	1,252	5.36	234	1,165
	<b>Cid</b>	30,336	5,615	1,048	4.87	215	1,181
	<b>Dre</b>	19,458	9,935	1,507	7.58	199	1,280
<b>Homolog</b>	<b>Hom</b>	16,090	10,844	1,432	7.83	183	1,379
	<b>Tru</b>	23,120	8,191	1,225	6.12	200	1,362
	<b>Mmu</b>	16,164	10,803	1,417	7.74	183	1,392
	<b>Ipu</b>	37,610	6,704	1,155	5.22	221	1,315
	<b>PASA</b>	97,309	9,419	1,201	7.09	169	1,348
<b>RNASeq</b>	<b>Cufflin ks</b>	92,180	19,478	4,707	10.13	465	1,618
	<b>EVM</b>	25,365	11,517	1,323	7.66	173	1,531
	<b>PASA-update*</b>	38,086	13,009	1,521	8.79	173	1,475
	<b>Final set*</b>	22,066	12,913	1,458	8.48	172	1,531

421 [back](#) to main body

422

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

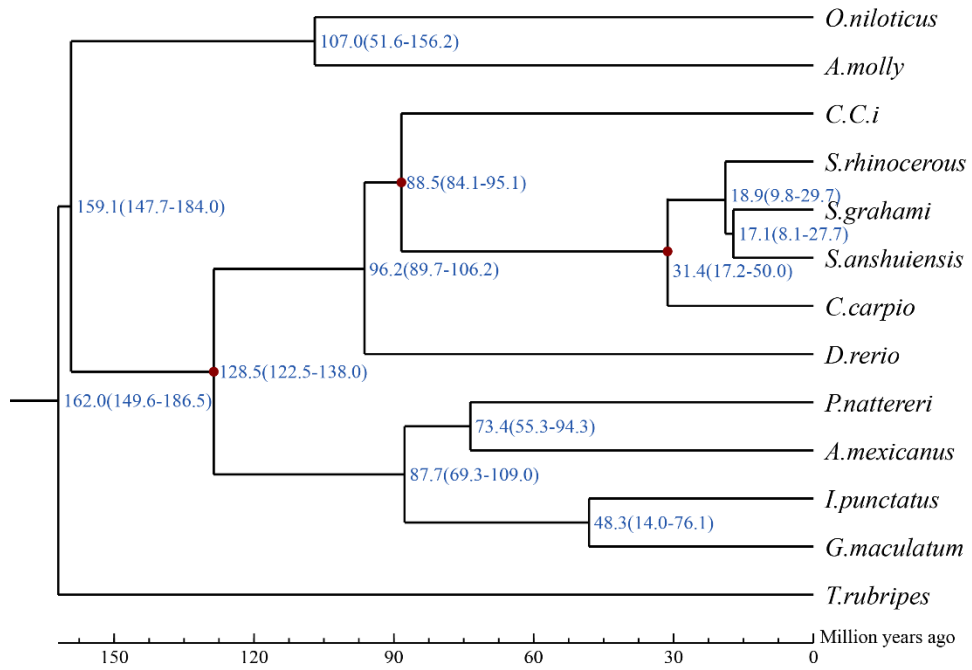
423 **Figure 1. A picture showing about *G. maculatum*.** (a) The appearance of *G.*  
424 *maculatum*. (b) Distributed localization (red triangle) of *G. maculatum* for sequencing.  
425 (c) The liver of *G. maculatum* was divided to two parts, one placed outside the  
426 abdominal cavity (attaching liver), connected to another part that located inside the  
427 cavity (mail liver) (Figure schematic drawings (ventral view) of *G. maculatum* (imaged  
428 from Zhang[9]).



34  
35 429  
36 430  
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

431

432 **Figure 2. Divergence time estimated between *G. maculatum* and other species**



433

434

435 [back](#) to main body

436

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



Click here to access/download  
**Supplementary Material**  
Supplemental\_file.docx





Click here to access/download

**Supplementary Material**

Supplemental\_file\_MITE.docx

