# GigaScience

## A chromosomal-level genome assembly for the giant African snail Achatina fulica

### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-19-00006R3 |
| Full Title: | A chromosomal-level genome assembly for the giant African snail Achatina fulica |
| Article Type: | Data Note |
| Funding Information: | This work was supported by the National Key Research and Development Program of China (No. 2016YFC1200500 and 2016YFC1202000) — Dr Ning Xiao |

**Abstract:**

Background:
Achatina fulica (A. fulica), also called the giant African snail, is the largest terrestrial mollusk species. Due to its voracious appetite, wide environmental adaptability, high growth rate and reproductive capacity, it has become an invasive species across the world, mainly in Southeast Asia, Japan, the western Pacific islands and China. A. fulica is a pest that is able to damage agricultural crops, as well as an intermediate host of many parasites that can threaten human health. However, genomic information of A. fulica is still limited, hindering genetic and genomic studies for invasion control and management of the species.
Finding:
Using a Kmer-based method, we estimated the A. fulica genome size to be 2.12 Gb with a high repeat content up to 71%. Roughly 101.6 Gb genomic long-read data of A. fulica were generated from the PacBio sequencing platform and assembled to produce a first A. fulica genome of 1.85 Gb with a contig N50 length of 726 kb. Using contact information from the Hi-C sequencing data, we successfully anchored 99.32% contig sequences into 31 chromosomes, leading to the final contig and scaffold N50 length of 721 kb and 59.6 Mb, respectively. The continuity, completeness and accuracy were evaluated by genome comparison with other mollusk genomes, BUSCO assessment and genomic read mapping. 23,726 protein-coding genes were predicted from the assembled genome, among which 96.34% of the genes were functionally annotated. The phylogenetic analysis using whole-genome protein-coding genes revealed that A. fulica separated from a common ancestor with Biomphalaria glabrata around 182 million years ago.
Conclusion:
To the best of our knowledge, the A. fulica genome is the first terrestrial mollusk genome published to date. The chromosome sequence of A. fulica will provide the research community with a valuable resource for population genetics and environmental adaptation studies for the species as well as investigations of the chromosome level of evolution within mollusks.

| | |
|---|---|
| Corresponding Author: | ning xiao<br><br>CHINA |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Yunhai Guo |
| First Author Secondary Information: | |
| Order of Authors: | Yunhai Guo |
| | Yi Zhang |
| | Qin Liu |
| | Yun Huang |

|  | Guangyao Mao |
|  | Zhiyuan Yue |
|  | Eniola M. Abe |
|  | Jian Li |
|  | Zhongdao Wu |
|  | Shizhu Li |
|  | Xiaonong Zhou |
|  | Wei Hu |
|  | Ning Xiao |

| Order of Authors Secondary Information: | |
|---|---|
| Response to Reviewers: | Reviewer #1: Regarding RNA samples, the authors commented that "we eventually selected high-quality samples for the sequencing." and showed Bioanalyzer report of only four samples in the supplementary. In the revised manuscript, however, they described that "RNA was extracted from the pallium, liver, foot, spleen, stomach, gut, heart...". Which four samples or tissues were actually used for RNA extraction and sequencing?<br>Reply: Thanks a lot for reminding. We have collected multi samples for pallium, liver, foot, spleen, stomach, gut and heart tissues. The qualities of RNA samples were summarized in the following table. We selected the high-quality samples highlighted by red for the PacBio library construction and sequencing in the Supplementary Table S1, covering all tissues that mentioned above. We have accordingly revised the manuscript and supplementary information. |

| Additional Information: | |
|---|---|
| **Question** | **Response** |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>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.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>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 | Yes |

| | |
|---|---|
| Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | |
| **Availability of data and materials**<br><br>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.<br><br>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist? | Yes |

# A chromosomal-level genome assembly for the giant African snail *Achatina fulica*

Guo Yunhai[1,2, #]，Zhang Yi[1,2, #]，Liu Qin[1,2]，Huang Yun[1,2]，Mao Guangyao[1,2]，

Yue Zhiyuan[1,2]，Eniola M. Abe[1,2]，Li Jian[3]，Wu Zhongdao[4]，Li Shizhu[1,2]，Zhou

Xiaonong[1,2]，Hu Wei[1,2,3,*]，Xiao Ning[1,2,*]


[1]National Institute of Parasitic Diseases, Chinese Center for Disease Control and

Prevention

[2]Key Laboratory of Parasite and Vector Biology, Ministry of Health, Shanghai, China

[3]Department of Microbiology and Microbial Engineering , School of Life Sciences , Fudan

University , Shanghai 200438 , China

[4]Department of Parasitology, Zhongshan School of Medicine, Sun Yat-sen University,

Guangzhou 510080, China

ORCID:

Ning Xiao, 0000-0002-1361-7013

Eniola Abe: 0000-0002-9573-5506

## Abstract

**Background**:

*Achatina fulica (A. fulica),* also called the giant African snail, is the largest terrestrial mollusk species. Due to its voracious appetite, wide environmental adaptability, high growth rate and reproductive capacity, it has become an invasive species across the world, mainly in Southeast Asia, Japan, the western Pacific islands and China. *A. fulica* is a pest that is able to damage agricultural crops, as well as an intermediate host of many parasites that can threaten human health. However, genomic information of *A. fulica* is still limited, hindering genetic and genomic studies for invasion control and management of the species.

**Finding**:

Using a *K*mer-based method, we estimated the *A. fulica* genome size to be 2.12 Gb with a high repeat content up to 71%. Roughly 101.6 Gb genomic long-read data of *A. fulica* were generated from the PacBio sequencing platform and assembled to produce a first *A. fulica* genome of 1.85 Gb with a contig N50 length of 726 kb. Using contact information from the Hi-C sequencing data, we successfully anchored 99.32% contig sequences into 31 chromosomes, leading to the final contig and scaffold N50 length of 721 kb and 59.6 Mb, respectively. The continuity, completeness and accuracy were evaluated by genome comparison with other mollusk genomes, BUSCO assessment and genomic read mapping. 23,726 protein-coding genes were predicted from the assembled genome, among which 96.34% of the genes were functionally annotated. The phylogenetic analysis using whole-genome protein-coding genes revealed that *A. fulica* separated from a common ancestor with *Biomphalaria glabrata* around 182 million years ago.

**Conclusion**:

To the best of our knowledge, the *A. fulica* genome is the first terrestrial mollusk genome published to date. The chromosome sequence of *A. fulica* will provide the research community with a valuable resource for population genetics and environmental adaptation studies for the species as well as investigations of the chromosome level of evolution within mollusks.


**Key Words:** Giant African snail, *Achatina fulica*, PacBio, Hi-C, chromosome assembly

## Data description

## Introduction

The giant African snail, *A. fulica* (NCBI: txid6530), is a Gastropod species (**Figure 1**). It is the largest terrestrial mollusk, with a voracious appetite, strong environmental adaptability, and high growth and reproduction rate[1-3]. Originating in East Africa, *A. fulica* over the last century it has gradually invaded Southeast Asia, Japan and the western Pacific islands [4-6] with the direct and indirect help from humans[7-9]. In mainland China, the first *A. fulica* invasion event was reported in 1931[10]. At present, the snail's natural distribution in the wild has been found in Guangdong, Hainan, Guangxi, southern parts of Yunnan Province and Fujian Province, and a county of Guizhou Province[11]. *A. fulica* was included as the first 16 alien invasive species in China (http://www.mee.gov.cn/gkml/zj/wj/200910/t20091022_172155.htm, in Chinese) in 2003, and was also listed by International Union for Conservation of Nature (IUCN) as among the 100 most threatening alien invasive species[12]. This snail has been recognized as an agricultural and garden pest causing significant damage in both tropical and subtropical regions[9, 12, 13]. In addition, *A. fulica* is also the intermediate host of the parasitic nematode *Angiostrongyl cantonensis*. Human infection with angiostrongyliasis, which occurs mainly through consumption of snails carrying *A. cantonensis* larvae, causes eosinophilic meningoencephalitis[4, 11, 14-18]. As a consequence, *A. fulica* is attracting more and more attention in the fields of agricultural crops protection and human disease control.

To date, a variety of mollusk genomes have been analyzed and published, including two freshwater gastropods snails *Pomacea canaliculata*[19] and *Biomphalaria glabrata*[20]. However, no genome has been reported for terrestrial mollusks. *A. fulica* is considered to be a destructive terrestrial gastropod which poses a significant hazard to agriculture, the environment, biodiversity and human health. A chromosome-level genome of *A. fulica*could would provide crucial resources in the

population genetics and evolution studies based on genomic sequencing data aiming to discover the invasion and adaptation history of *A. fulica*. Furthermore, the genome could also be used to probe gene expression during important biological processes, such as gene expression patterns in various developmental stages and the interaction of *Angiostrongylus* and *A. fulica*. In this study we applied Illumina, PacBio and Hi-C techniques to construct the chromosome of *A. fulica*. The genome is the first terrestrial mollusk genome, providing an important reference for the molecular mechanisms underlying its broad environmental adaptability and the development of control strategy of the world-wide invasion.

**Sample and sequencing**

An adult snail (**Figure 1**), which was collected in Pingxiang city, Guangxi Autonomous Region, was used for reference genome construction. The snail was dissected and abdominal foot (17.4 g) and liver pancreas (40.4 g) tissues were collected and quickly frozen in liquid nitrogen overnight before transferring to -80 °C for storage. DNA was extracted using the traditional phenol/chloroform extraction method and was quality checked using agarose gel electrophoresis, meeting the requirement for library construction for the Illumina X Ten (Illumina Inc., San Diego, CA, USA) and for the PacBio Sequel (Pacific Biosciences of California, Menlo Park, CA, USA) sequencing platforms.

RNA was extracted from the pallium, liver, foot, spleen, stomach, gut, heart using TRIzol® Reagent (Life Technologies, USA). The RNA quality was checked using the Nanodrop ND-1000 spectrophotometer (LabTech, USA) and 2100 Bioanalyzer (Agilent Technologies, USA) with RNA integrity number lager than 8 (Supplemental Table S1 and Supplemental Figure S1). The RNA from each samples were equally mixed for the RNA sequencing on PacBio Sequel platform. Firstly, mRNA molecules were reversely transcribed to cDNA using Clontech SMARTer cDNA synthesis kit. After cDNA amplification and purification, two SMRTbell libraries of 0-4 kb and 4-10

115 kb were generated using the size selection in BluePippin Size Selection System

116 (Pacific Biosciences of California, Menlo Park, CA, USA) and protocols suggested by

117 manufacturer. The finale libraries were sequenced in the PacBio SEQUEL platform

118 (Pacific Biosciences of California, Menlo Park, CA, USA), resulting 12,439,996

119 subreads totaling about 22.5 Gb PacBio long reads with average length longer than

120 1,801 bps. Subsequently, a total of 782,613 circular consensus sequences (CCS)

121 were generated based on the subreads, and a number of 553,889 Full-length

122 Non-chimeric sequences (FLNC) representing 23,726 gene loci were obtained,

123 eventually. All aforementioned data processing were performed using SMRT Link

124 v5.0 (www.pacb.com). Moreover, about 70.37% of the multi-exon FLNCs were really

125 full-length sequences embracing all the exons of the gene locus predicted from the

126 whole genome sequences.

127 Using the DNA from the abdominal foot, a library with the insertion length of 300

128 bp was constructed and sequenced using the Illumina sequencing platform according

129 to the manufacturer's protocol. About 202.23 Gb short reads were obtained using

130 Illumina X Ten sequencing technology (**Table 1**), which was used for the following

131 genome survey analysis, and for final base-level genome sequence correction.

132 Meanwhile, four 20 kb libraries were constructed for PacBio Sequel sequencing.

133 Using 16 sequencing SMRT cells, 104.6 Gb long reads were generated (**Table 1**).

134 The mean and N50 lengths of the polymerases for sequencing cells ranged from 6.4

135 kb to 10.4 kb and from 12.3 kb to 20.3 kb for cells, respectively. Those long genomic

136 DNA reads were then used for reference genome construction.

137

138 **Genome features estimation from *K*mer method**

139 With sequencing data from the Illumina platform, several genome characters could be

140 evaluated for *A. fulica.* To ensure the quality of the analysis, ambiguous bases and

141 low-quality reads were trimmed and filtered using the HTQC package (version

142  1.92.3)[21]. The following quality control were performed under the framework of

143  HTQC. First, the quality of bases at two read ends were checked. Bases in sliding 5

144  bp windows were deleted if the average quality of the window was below phred quality

145  score of 20. Second, reads were filtered if the average phred quality score were

146  smaller than 20 or the read length was shorter than 75 bp. Third, the mate reads were

147  also removed if the corresponding reads were filtered.

148  The quality-controlled reads were used for genome character estimation. We

149  calculated the number of each 17-mer from the sequencing data using the jellyfish

150  software (Jellyfish, RRID:SCR_005491; version 2.0)[22], and the distribution was

151  analyzed with GCE software (GCE, RRID:SCR_017332; version 3)[23] and was

152  shown in Supplemental Figure S2. We estimated the genome size of 2.12 Gb with the

153  heterozygosity of 0.47% and repeat content of 71% in the genome. Previous studies

154  revealed that repeat content varies in mollusks, and that repeat content is correlated

155  with genome size[24]. The large genome size and high proportion of repeat contents

156  of *A. fulica* provided additional supporting data for the statistical analysis. Moreover,

157  10,000 pairs of short reads were extracted randomly and were compared to the nt

158  database and no obvious external contamination were found.

159  **Genome assembly by third-generation long reads**

160  After removing adaptor sequences in polymerases, 101.6 Gb subreads were

161  generated for the following whole genome assembly. The average and N50 length of

162  subreads reached 5.25 kb and 8.80 kb, respectively. To optimize the genome

163  assembly using the PacBio sequencing data, we applied two packages in the

164  assembly process, Canu v1.8 (Canu, RRID:SCR_015880) [25] and FALCON v0.2.2

165  (Falcon, RRID:SCR_016089) [26]. Canu package was first applied for the assembly

166  using default parameters. As a result, a 1.93 Gb genome was constructed with 10,417

167  contigs and a contig N50 length of 662.40 kb. FALCON was also employed using the

168  length_cutoff and pr_length_cutoff parameters of 10 kb and 8 kb, respectively. We

169  obtained 1.85 Gb genome with 8,585 contigs, with a contig N50 of 726.63 kb. We

170  adopted the FALCON assembly as the reference genome for *A. fulica* (**Table 2**).

171  Compared to the estimated genome size, the assembled version was relatively

172  smaller, which may have resulted from the following two possible scenarios: the high

173  repeat content of the genome, and the probably larger size estimated from the Kmer

174  analysis. The genome sequences were subsequently polished, PacBio long reads

175  utilizing arrow[27] and Illumina short reads using pilon[28] to correct base errors. The

176  corrected genome was further applied for the following chromosome assembly

177  construction using Hi-C data.

178  ***In situ* Hi-C library construction and chromosome assembly using Hi-C**

179  **data**

180  Liver pancreas tissue of *A. fulica* was used for library construction for Hi-C analysis

181  and the library was constructed using the identical method in previous studies[29].

182  Finally, the library was sequenced with 150 paired-end mode on the Illumina HiSeq X

183  Ten platform (San Diego, CA, United States). From the Illumina sequencing platform,

184  1,313.87 million paired-end reads were obtained for the Hi-C library (**Table 1**). The

185  reads were mapped to the above *A. fulica* genome with Bowtie2 [30], with two ends of

186  paired reads being mapped to the genome separately. To increase the interactive Hi-C

187  reads ratio, an iterative mapping strategy was performed as previous studies, and

188  only read pairs with both ends uniquely mapped were used for the following analysis.

189  From the alignment status of two ends, self-ligation, non-ligation and other sorts of

190  invalid reads, including StartNearRsite, PCR amplification, random break,

191  LargeSmallFragments and ExtremeFragments, were filtered out by Hi-Clib[31].

192  Through the recognition of restriction sites in sequences, contact counts among

193  contigs were calculated and normalized.

194  According to previous karyotype analyses, *A. fulica* has 31 chromosomes[32]. By

195  clustering the contigs using the contig contact frequency matrix, we were able to

196 correct some minor errors in the FALCON assembly results. Contigs with errors were
197 broken into shorter contigs. We obtained 8,701 contigs, slightly more than the 8,585
198 contigs in the FALCON assembly. We successfully clustered these contigs into 31
199 groups in Lachesis[33] using the agglomerative hierarchical clustering method
200 (**Figure 2**). Lachesis was further applied to order and orient the clustered contigs
201 according to the contact matrix. As a result, 7,106 contigs were reliably anchored,
202 ordered and orientated on chromosomes, accounting for 99.32% of the total genome
203 bases. The first near chromosomal-level assembly of *A. fulica* was obtained with
204 8,211 contigs, a contig N50 of 721.0 kb and a scaffold N50 of 59.59 Mb (**Table 2** and
205 **Table 3**).

206 **Genome quality evaluation**

207 We assessed the quality of genome of *A. fulica* after the assembly process. The
208 quality evaluation was carried out in three aspects: continuity, completeness and the
209 mapping rate of NGS data.

210 First of all, we compared the sequence number and contig N50 length of *A. fulica*
211 with public genome of mollusks and found that our assembly has a high quality on
212 contig and scaffold N50 among mollusk genomes. (**Table 3**) Traditional chromosomal
213 genome assembly requires physical maps and genetic maps, which is enormously
214 time- and labor-consuming. With Hi-C data analysis, we successfully assembled *A.*
215 *fulica* genome into near chromosome-level with just one individual.

216 Second, the assembled genome was subjected to the BUSCO (version 3.0,
217 metazoa_odb9)[34] to assess the completeness of the genome. About 91.7% of the
218 BUSCO genes were identified in *A. fulica* genome, and more than 84.7% of the
219 BUSCO genes were single-copy completed in our genome, illuminating a high level of
220 completeness of the genome.

221    Third, NGS short reads were aligned to the genome using BWA package (BWA,

222    RRID:SCR_010910; version 0.7.17)[35], and about 98.7% of paired reads were

223    aligned to the genome, of which 98.24% were reads paired aligned.

224    **Repeat element and gene annotation**

225    Tandem Repeat Finder4.09 (TRF)[36] was used for repetitive element identification in

226    the *A. fulica* genome. A *de novo* method applying RepeatModeler (RepeatModeler,

227    RRID:SCR_015027) was used to detect transposable elements (TEs). The resulted

228    *de novo* data, combined with known repeat library from Repbase[37], were used to

229    identify TEs in the *A. fulica* genome by RepeatMasker4-0-8 (RepeatMasker,

230    RRID:SCR_012954) [38] software. All repetitive elements were masked in the

231    genome before protein-coding gene prediction.

232    Protein-coding genes in the *A. fulica* genome were annotated using the *de novo*

233    program Augustus0.2.1 (Augustus, RRID:SCR_008417) [39]. Protein sequences of

234    the closely related species including *Aplysia californica*, *Biomphalaria glabrata*，

235    *Crassostrea gigas*，*Lottia gigantea* and *Patinopecten yessoensis*, were downloaded

236    from the Ensembl database, and aligned to the *A. fulica* genome with TBLASTN2.6.0

237    (TBLASTN, RRID:SCR_011822). Full-length transcripts obtained using Iso-Seq were

238    mapped to the genome using Genewise (GeneWise, RRID:SCR_015054) [40]. Finally,

239    gene models predicted from all above methods were combined by MAKERv2.31.10

240    (MAKER, RRID:SCR_005309) [41], resulting in 23,726 protein-coding genes. The

241    gene number, gene length, CDS length, exon length and intron length distribution

242    were all comparable with the related mollusks (**Figure 3**).

243    To functionally annotate protein-coding genes in the *A. fulica* genome, we

244    searched all predicted gene sequences to NCBI non-redundant nucleotide (NT) and

245    protein (NR), Swiss-Prot databases by BLASTN (BLASTN, RRID:SCR_001598) [42]

246    and BLASTX (BLASTX, RRID:SCR_001653) [43] utility. Blast2GO (Blast2GO,

247    RRID:SCR_005828) [44] was also used to assign gene ontology (GO)[45] and Kyoto

248　Encyclopedia of Genes and Genomes (KEGG)[46] pathways. A threshold of e-value

249　of 1e-5 was used for all BLAST applications. Finally, 22,858 (96.34%) genes were

250　functionally annotated (**Table 4**).

251　**Phylogenetic analysis of *A. fulica* with other mollusks**

252　OrthoMCLv1.2 [47] was used to cluster gene families. First, proteins from *A. fulica*

253　and the closely related mollusks, including *Aplysia californica*, *Biomphalaria glabrata*,

254　*Crassostrea gigas*, *Lingula anatina*, *Lottia gigantea*, *Patinopecten yessoensis*,

255　*Octopus bimaculoides*, *Helobdella robusta*, *Pomacea canaliculata*, and the outgroup,

256　*Drosophila melanogaster*, were all-to-all blasted by BLASTP (BLASTP,

257　RRID:SCR_001010) [43] utility with an e-value threshold of 1e-5. Only proteins from

258　the longest transcript were used for genes with alternative isoforms. We identified

259　25,448 gene families for *A. fulica* and the related species, among them 675

260　single-copy orthologs families were detected.

261　　　Using single-copy orthologs, we could probe the phylogenetic relationships for

262　the *A. fulica* and other mollusks. To this end, protein sequences of single-copy genes

263　were aligned using CLUSTALX2.0 (Clustal X, RRID:SCR_017055) [48]. Guided by

264　the protein multi-sequence alignment, the alignment of the coding DNA sequences

265　(CDS) for those genes were generated and concatenated for the following analysis.

266　The phylogenetic relationships were constructed using PhyML3.0 (PhyML,

267　RRID:SCR_014629) [49] using the concatenated nucleotide alignment with the

268　JTT+G+F model. The MCMCtree program in PAML4 [49] was used to estimate the

269　species divergent time scales for the mollusks using approximate likelihood method

270　and calibrated according to the fossil records. We found that *A. fulica* was most

271　closely related to *Biomphalaria glabrata*, and the two species diverged from their

272　common ancestor about 242 million years ago (MYA) (**Figure 4**).

273　**Conclusion**

274　We reconstructed the first chromosome level assembly for *A. fulica* using an

275    integrated sequencing strategy combining PacBio, Illumina and Hi-C technologies.

276    Using the long reads from the PacBio Sequel platform and short reads from the

277    Illumina X Ten platform, we successfully constructed contig assembly for *A. fulica*.

278    Leveraging contact information among contigs from Hi-C technology, we further

279    improved the assembly to near chromosome-level quality (**Table 3** and **Figure 2**). We

280    predicted 23,726 protein-coding genes in the *A. fulica* genome and 22,858 of genes

281    were functionally annotated with putative functions. With 675 single-copy orthologs

282    from *A. fulica* and other related mollusks, we constructed the phylogenetic

283    relationship of these mollusks, and found that *A. fulica* might have diverged from its

284    common ancestor of *Biomphalaria glabrata* around 177.1-187.1 MYA. Given the

285    increasing interest in mollusk genomic evolution and the biological importance of *A.*

286    *fulica* as an invasive animal, our genomic and transcriptome data will provide valuable

287    genetic resources for follow-on functional genomics investigations by the research

288    community.

289

290    **Ethics Statement**

291    This study was approved by the Animal Care and Use committee of National Institute

292    of Parasitic Diseases, Chinese Center for Disease Control and Prevention.

293    **Abbreviation**

294    CCS: circular consensus sequences; CDS: coding DNA sequences; FLNC:

295    Full-length Non-chimeric sequences; GO: gene ontology; KEGG: Kyoto Encyclopedia

296    of Genes and Genomes; MYA: million years ago; TE: transposable elements

297    **Availability of supporting data**

298    The Illumina, PacBio and Hi-C sequencing data are available from NCBI via the

299    accession number of SRR8369706, SRR8369311 and SRR8371669, respectively.

300    The Illumina transcriptome sequencing data were deposited to NCBI via the

accession number of SRR8371872 and SRR8371873. The genome, annotation and

intermediate files were uploaded to *GigaScience* GigaDB Database [66].

## Author Contributions

Z.X, H.W and X.N conceived the project. G.Y, Z.Y, L.Q collected the samples and

extracted the genomic DNA. G.Y, Z.Y and L.Q performed the genome assembly and

data analysis. G.Y, Z.X, H.W and X.N wrote the paper.

## References

1.    Schreurs J. *Investigations on the biology, ecology and control of Giant African Snail 290 in West New Guinea*.    1963. Manokwari Agricultural Research Station.

2.    Albuquerque FS, Peso-Aguiar MC and Assunção-Albuquerque MJ. Distribution,feeding behavior and control strategies of the exotic land snail Achatinafulica (Gastropoda:Pulmonata) in the Northeast of Brazil. BrazJBiol. 2008;68:6.

3.    Thiengo SC, Fernandez MA, Torres EJ, Coelho PM and Lanfredi RM. First record of anematode Metastrongyloidea (Aelurostrongylus abstrusus larvae) in Achatina (Lissachatina) fulica (Mollusca,Achatinidae) in Brazil. J Invertebr Pathol. 2008;98:6.

4.    Lv S, Zhang Y and Liu HX. Invasive Snails and an Emerging Infectious Disease: Results from the First National Survey on Angiostrongylus cantonensis in China. BioOne. 2009; doi:10.1371/journal.pntd.0000368.

5.    Cowie RH. *Non-indigenous land and freshwater molluscs in the islands of the Pacific: Conservation impacts and threats*.    2000.

6.    Cowie RH. Can snails ever be effective and safe biocontrol agents? Int J Pest Manage. 2001;47:18.

7.    Cowie RH and Robinson DG. Pathways of introduction of nonindigenous land and freshwater snails and slugs. Washington DC: Island Press; 2003.

8.    Kotangale JP. Giant African snail (Achatina fulica Bowdich). 2011;J Environ Sci Eng 53:6.

9.    Raut SK and Barker GM. Achatina fulica Bowdich and Other Achatinidae as Pests in Tropical Agriculture. UK: CABI International; 2002.

10.    Jarreit VHC. The spread of the snail Achatina fulica to south China. Hong Kong Nat. 1931;2:3.

335    11.    Shan L, Yi Z and Peter S. Emerging Angiostrongyliasis in Mainland China. Emerging Infectious
336          Diseases. 2008;14 1:4.

337    12.    Lowe S, Browne SM, Boudjrlas S and De Poorter M. 100 of the world's worst invasive alien
338          species: A selection from the global invasive species database. The Invasive Species
339          Specialists Group of the Species Survival Commission of the World Conservation Union.
340          Auckland: Hollands Printing; 2000.

341    13.    Mead AR. Pulmonates volume 2B. Economic malacology with particular reference to
342          Achatina fulica. London: Academic Press; 1979.

343    14.    Alicata JE. The discovery of Angiostrongylus cantonensis as a cause of human eosinophilic
344          meningitis. Parasitol Today. 1991;7 6:151-3.

345    15.    Prociv P, Spratt DM and Carlisle MS. Neuro-angiostrongyliasis: unresolved issues. Int J
346          Parasitol. 2000;30 12-13:1295-303.

347    16.    Deng ZH, Zhang QM, Huang SY and Jones JL. First provincial survey of Angiostrongylus
348          cantonensis in Guangdong Province, China. Trop Med Int Health. 2012;17:4.

349    17.    Maldonado JA, Simoes RO, Oliveira AP, Motta EM, Fernandez MA, Pereira ZM, et al. First
350          report of Angiostrongylus cantonensis (Nematoda: Metastrongylidae) in Achatina fulica
351          (Mollusca: Gastropoda) from Southeast and South Brazil. Mem Inst Oswaldo Cruz.
352          2010;105:4.

353    18.    Vitta A, Polseela R, Nateeworanart S and Tattiyapong M. Survey of Angiostrongylus
354          cantonensis in rats and giant African land snails in Phitsanulok Province, Thailand. Asian Pac J
355          Trop Med. 2011;4:3.

356    19.    Liu C, Zhang Y, Ren Y, Wang H, Li S, Jiang F, et al. The genome of the golden apple snail
357          *Pomacea canaliculata* provides insight into stress tolerance and invasive adaptation.
358          GigaScience. 2018;7 9 doi:10.1093/gigascience/giy101.

359    20.    Adema CM, Hillier LW, Jones CS, Loker ES, Knight M, Minx P, et al. Whole genome analysis of
360          a schistosomiasis-transmitting freshwater snail. Nature communications. 2017;8:15451.
361          doi:10.1038/ncomms15451.

362    21.    Neff KL, Argue DP, Ma AC, Lee HB, Clark KJ and Ekker SC. Mojo Hand, a TALEN design tool for
363          genome editing applications. BMC Bioinformatics. 2013;14:1. doi:10.1186/1471-2105-14-1.

364    22.    Marcais G and Kingsford C. A fast, lock-free approach for efficient parallel counting of
365          occurrences of k-mers. Bioinformatics. 2011;27 6:764-70.
366          doi:10.1093/bioinformatics/btr011.

367    23.    Binghang Liu YS, Jianying Yuan,Xuesong Hu,Hao Zhang,Nan Li,Zhenyu Li,Yanxiang
368          Chen,Desheng Mu,Wei Fan. Estimation of genomic characteristics by analyzing k-mer
369          frequency in de novo genome projects. Quantitative Biology. 2013;35:62-7.

370    24.    Murgarella M, Puiu D, Novoa B, Figueras A, Posada D and Canchaya C. A First Insight into the
371          Genome of the Filter-Feeder Mussel Mytilus galloprovincialis. PloS one. 2016;11 3:e0151561.
372          doi:10.1371/journal.pone.0151561.

373    25.    Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH and Phillippy AM. Canu: scalable and
374          accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome
375          Res. 2017;27 5:722-36. doi:10.1101/gr.215087.116.

376    26.    Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, et al. Phased diploid
377          genome assembly with single-molecule real-time sequencing. Nat Methods. 2016;13
378          12:1050-4. doi:10.1038/nmeth.4035.

379  27.  Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid, finished
380      microbial genome assemblies from long-read SMRT sequencing data. Nature methods.
381      2013;10 6:563.

382  28.  Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated
383      tool for comprehensive microbial variant detection and genome assembly improvement.
384      PloS one. 2014;9 11:e112963.

385  29.  Gong G, Dan C, Xiao S, Guo W, Huang P, Xiong Y, et al. Chromosomal-level assembly of yellow
386      catfish genome using third-generation DNA sequencing and Hi-C analysis. Gigascience. 2018;
387      doi:10.1093/gigascience/giy120.

388  30.  Langmead B, Trapnell C, Pop M and Salzberg SL. Ultrafast and memory-efficient alignment of
389      short DNA sequences to the human genome. Genome Biol. 2009;10 3:R25.
390      doi:10.1186/gb-2009-10-3-r25.

391  31.  Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO and Shendure J. Chromosome-scale
392      scaffolding of de novo genome assemblies based on chromatin interactions. Nature
393      biotechnology. 2013;31 12:1119.

394  32.  Sun T. Chromosomal studies in three land snails. Sinozoologia. 1995;12:154-62.

395  33.  Near TJ, Dornburg A, Eytan RI, Keck BP, Smith WL, Kuhn KL, et al. Phylogeny and tempo of
396      diversification in the superradiation of spiny-rayed fishes. Proceedings of the National
397      Academy of Sciences of the United States of America. 2013;110 31:12738.

398  34.  Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing
399      genome assembly and annotation completeness with single-copy orthologs. Bioinformatics.
400      2015;31 19:3210-2.

401  35.  Li H and Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.
402      bioinformatics. 2009;25 14:1754-60.

403  36.  Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids Res.
404      1999;27 2:573-80.

405  37.  Bao W, Kojima KK and Kohany O. Repbase Update, a database of repetitive elements in
406      eukaryotic genomes. Mob DNA. 2015;6:11. doi:10.1186/s13100-015-0041-9.

407  38.  Chen N. Using RepeatMasker to identify repetitive elements in genomic sequences. Current
408      protocols in bioinformatics. 2004;5 1:4.10. 1-4.. 4.

409  39.  Stanke M, Keller O, Gunduz I, Hayes A, Waack S and Morgenstern B. AUGUSTUS: ab initio
410      prediction of alternative transcripts. Nucleic acids research. 2006;34 suppl_2:W435-W9.

411  40.  Birney E, Clamp M and Durbin R. GeneWise and genomewise. Genome research. 2004;14
412      5:988-95.

413  41.  Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B, et al. MAKER: an easy-to-use
414      annotation pipeline designed for emerging model organism genomes. Genome research.
415      2008;18 1:188-96.

416  42.  Gertz EM, Yu YK, Agarwala R, Schaffer AA and Altschul SF. Composition-based statistics and
417      translated nucleotide searches: improving the TBLASTN module of BLAST. BMC Biol.
418      2006;4:41. doi:10.1186/1741-7007-4-41.

419  43.  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
420      architecture      and      applications.      BMC      Bioinformatics.      2009;10:421.
421      doi:10.1186/1471-2105-10-421.

422  44.  Conesa A, Götz S, García-Gómez JM, Terol J, Talón M and Robles M. Blast2GO: a universal

423  tool for annotation, visualization and analysis in functional genomics research.
424  Bioinformatics. 2005;21 18:3674-6.

425  45.  Consortium GO. The Gene Ontology (GO) database and informatics resource. Nucleic acids
426  research. 2004;32 suppl_1:D258-D61.

427  46.  Kanehisa M and Goto S. KEGG: kyoto encyclopedia of genes and genomes. Nucleic acids
428  research. 2000;28 1:27-30.

429  47.  Li L, Stoeckert CJ and Roos DS. OrthoMCL: identification of ortholog groups for eukaryotic
430  genomes. Genome research. 2003;13 9:2178-89.

431  48.  Thompson JD, Gibson TJ and Higgins DG. Multiple sequence alignment using ClustalW and
432  ClustalX. Current protocols in bioinformatics. 2003; 1:2.3. 1-2.3. 22.

433  49.  Guindon S, Lethiec F, Duroux P and Gascuel O. PHYML Online—a web server for fast
434  maximum likelihood-based phylogenetic inference. Nucleic acids research. 2005;33
435  suppl_2:W557-W9.

436  50.  Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al. The oyster genome reveals stress adaptation
437  and complexity of shell formation. Nature. 2012;490 7418:49-54. doi:10.1038/nature11413.

438  51.  Takeuchi T, Kawashima T, Koyanagi R, Gyoja F, Tanaka M, Ikuta T, et al. Draft genome of the
439  pearl oyster Pinctada fucata: a platform for understanding bivalve biology. DNA research : an
440  international journal for rapid publication of reports on genes and genomes. 2012;19
441  2:117-30. doi:10.1093/dnares/dss005.

442  52.  Takeuchi T, Koyanagi R, Gyoja F, Kanda M, Hisata K, Fujie M, et al. Bivalve-specific gene
443  expansion in the pearl oyster genome: implications of adaptation to a sessile lifestyle.
444  Zoological letters. 2016;2:3. doi:10.1186/s40851-016-0039-2.

445  53.  Du X, Fan G, Jiao Y, Zhang H, Guo X, Huang R, et al. The pearl oyster *Pinctada fucata*
446  *martensii* genome and multi-omic analyses provide insights into biomineralization.
447  GigaScience. 2017;6 8:1-12. doi:10.1093/gigascience/gix059.

448  54.  Mun S, Kim YJ, Markkandan K, Shin W, Oh S, Woo J, et al. The Whole-Genome and
449  Transcriptome of the Manila Clam (Ruditapes philippinarum). Genome biology and evolution.
450  2017;9 6:1487-98. doi:10.1093/gbe/evx096.

451  55.  Wang S, Zhang J, Jiao W, Li J, Xun X, Sun Y, et al. Scallop genome provides insights into
452  evolution of bilaterian karyotype and development. Nature ecology & evolution. 2017;1
453  5:120. doi:10.1038/s41559-017-0120.

454  56.  Schell T, Feldmeyer B, Schmidt H, Greshake B, Tills O, Truebano M, et al. An annotated draft
455  genome for *Radix auricularia* (Gastropoda, Mollusca). Genome biology and evolution. 2017;
456  doi:10.1093/gbe/evx032.

457  57.  Albertin CB, Simakov O, Mitros T, Wang ZY, Pungor JR, Edsinger-Gonzales E, et al. The
458  octopus genome and the evolution of cephalopod neural and morphological novelties.
459  Nature. 2015;524 7564:220-4. doi:10.1038/nature14668.

460  58.  Simakov O, Marletaz F, Cho SJ, Edsinger-Gonzales E, Havlak P, Hellsten U, et al. Insights into
461  bilaterian evolution from three spiralian genomes. Nature. 2013;493 7433:526-31.
462  doi:10.1038/nature11696.

463  59.  Kenny NJ, Namigai EK, Marletaz F, Hui JH and Shimeld SM. Draft genome assemblies and
464  predicted microRNA complements of the intertidal lophotrochozoans *Patella vulgata*
465  (Mollusca, Patellogastropoda) and Spirobranchus (Pomatoceros) lamarcki (Annelida,
466  Serpulida). Marine genomics. 2015;24 Pt 2:139-46. doi:10.1016/j.margen.2015.07.004.

467     60.     Barghi N, Concepcion GP, Olivera BM and Lluisma AO. Structural features of conopeptide
468             genes inferred from partial sequences of the *Conus tribble*i genome. Molecular genetics and
469             genomics : MGG. 2016;291 1:411-22. doi:10.1007/s00438-015-1119-2.

470     61.     Uliano-Silva M, Dondero F, Dan Otto T, Costa I, Lima NCB, Americo JA, et al. A
471             hybrid-hierarchical genome assembly strategy to sequence the invasive golden mussel,
472             *Limnoperna fortunei*. GigaScience. 2018;7 2 doi:10.1093/gigascience/gix128.

473     62.     Sun J, Zhang Y, Xu T, Zhang Y, Mu H, Zhang Y, et al. Adaptation to deep-sea chemosynthetic
474             environments as revealed by mussel genomes. Nature ecology & evolution. 2017;1 5:121.
475             doi:10.1038/s41559-017-0121.

476     63.     Jiao W, Fu X, Dou J, Li H, Su H, Mao J, et al. High-resolution linkage and quantitative trait
477             locus mapping aided by genome survey sequencing: building up an integrative genomic
478             framework for a bivalve mollusc. DNA research : an international journal for rapid publication
479             of reports on genes and genomes. 2014;21 1:85-101. doi:10.1093/dnares/dst043.

480     64.     Luo YJ, Takeuchi T, Koyanagi R, Yamada L, Kanda M, Khalturina M, et al. The Lingula genome
481             provides insights into brachiopod evolution and the origin of phosphate biomineralization.
482             Nature communications. 2015;6:8301. doi:10.1038/ncomms9301.

483     65.     Li C, Liu X, Liu B, Ma B, Liu F, Liu G, et al. Draft genome of the Peruvian scallop Argopecten
484             purpuratus. GigaScience. 2018;7 4 doi:10.1093/gigascience/giy031.

485     66. Guo Y; Zhang Y; Liu Q; Huang Y; Mao G; Yue Z; Abe EM; Li J; Wu Z; Li S; Zhou X; Hu W; Xiao N
486             (2019): Supporting data for "A chromosomal-level genome assembly for the giant African
487             snail *Achatina fulica*" GigaScience Database. http://dx.doi.org/10.5524/100647

488

489

## Tables and Figures

**Table 1: Sequencing data generated for *A.fulica* genome assembly and annotation**

| Library type | Platform | Library size (bp) | Data size (Gb) | Application |
|---|---|---|---|---|
| Short reads | HiSeq X Ten | 350 | 202.24 | Genome survey and genomic base correction |
| Long reads | PacBio SEQUEL | 20,000 | 101.63 | Genome assembly |
| Hi-C | HiSeq X Ten | 300-500 | 199.73 | Chromosome construction |

**Table 2: Statistics for genome assembly of *A. fulica***

| Sample ID | Length | | Number | |
|---|---|---|---|---|
| | Contig** (bp) | Scaffold (bp) | Contig** | Scaffold |
| Total | 1,852,282,574 | 1,855,883,074 | 8,211 | 1,010 |
| Max | 5,947,392 | 116,558,012 | - | - |
| N50 | 721,038 | 59,589,303 | 697 | 13 |
| N60 | 538,883 | 58,013,356 | 995 | 16 |
| N70 | 399,612 | 53,672,006 | 1,396 | 20 |
| N80 | 268,901 | 50,673,968 | 1,957 | 23 |
| N90 | 141,756 | 44,109,545 | 2,888 | 27 |

The two stars (**) means the ultimate contigs since they were probably modified during the Hic step.

516    Table 3 Summary of the genome of *A. fulica* and other published mollusk genomes.

| Species | Size* (Mb) | Contig N50(kb) | Scaffold N50(kb) |
|---|---|---|---|
| *Achatina fulica* (this study)** | 2,120 | 721 | 59,590 |
| *Pomacea canaliculata*[19]** | 570 | 995 | 38,000 |
| *Crassostrea gigas*[50] | 545 | 7.5 | 401 |
| *Pinctada fucata*[51] | 1,150 | 1.6 | 14.5 |
| *Pinctada fucata new*[52] | 1,150 | 21 | 324 |
| *Pinctada fucata* V2[53] | 1,150 | 21 | 167 |
| *Biomphalaria glabrata*[20] | 931 | 7.3 | 48 |
| *Ruditapes philippinarum*[54] | 1,370 | 3.3 | 32.7 |
| *Patinopecten yessoensis*[55]** | 1,430 | 38 | 41,000 |
| *Radix auricularia*[56] | 1,600 | 0.324 | 578 |
| *Octopus bimaculoides*[57] | 2,800 | 5.4 | 470 |
| *Mytilus galloprovincialis*[24] | 1,600 | 2.6 | 2.9 |
| *Lottia gigantea*[58] | 420 | 96 | 1,870 |
| *Patella vulgata*[59] | 1,460 | 3.1 | 3.1 |
| *Aplysia californica* | 1,760 | 9.6 | 917 |
| *Conus tribblei*[60] | 2,760 | 0.85 | 215 |
| *Limnoperna fortunei*[61] | 1,600 | 10 | 312 |
| *Bathymodiolus platifrons*[62] | 1,640 | 13.2 | 343 |
| *Modiolus philippinarum*[62] | 2,380 | 19.7 | 100.2 |
| *Chlamys farreri*[63] | 1,200 | 1.2 | 1.5 |
| *Lingula anatina*[64] | 463 | 55 | 294 |
| *Argopecten prupruatus*[65] | 885 | 80.1 | 1,020 |

517    * Estimated size of the genome

518    ** Genomes assembled into near chromosomal level

519

520    **Table 4: Statistics for genome annotation of *A. fulica***

| Database | Number | Percent |
|---|---|---|
| InterPro | 16,252 | 68.50 |
| GO | 12,101 | 51.00 |
| KEGG ALL | 21,325 | 89.88 |
| KEGG KO | 10,161 | 42.83 |
| Swissprot | 17,050 | 71.86 |
| TrEMBL | 22,403 | 94.42 |
| NR | 22,553 | 95.06 |
| Total | 23,726 | |

521

522
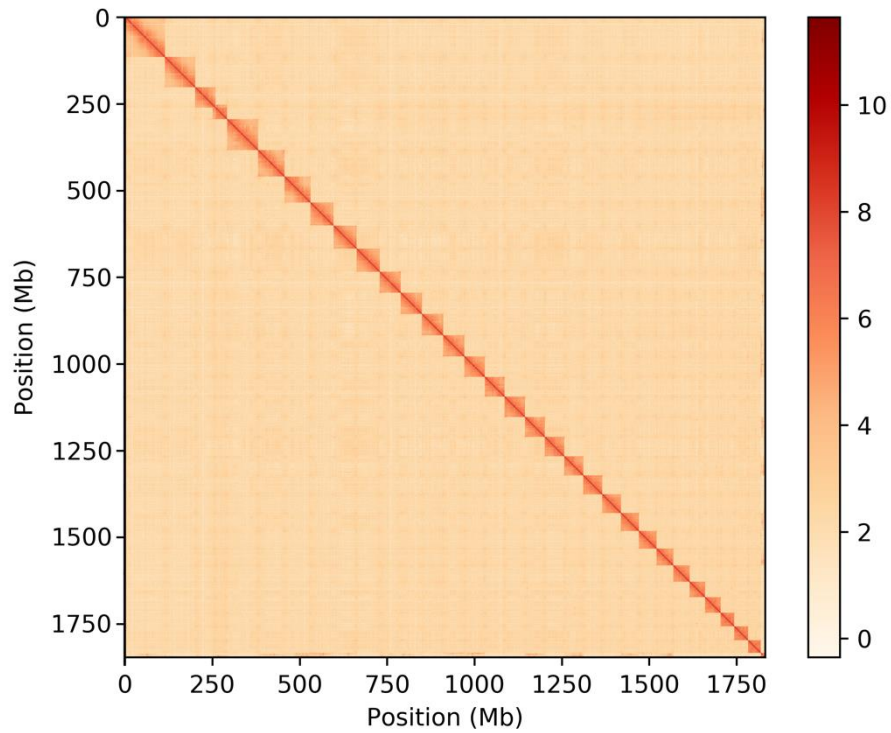
523



524

Figure 1. *A. fulica* individual **used for genome sequencing and assembly.**

526

527



528

**Figure 2. Contact matrix generated from the Hi-C data analysis showing sequence interactions in chromosomes.** The logarithm of the contact density were showed in the color bar.
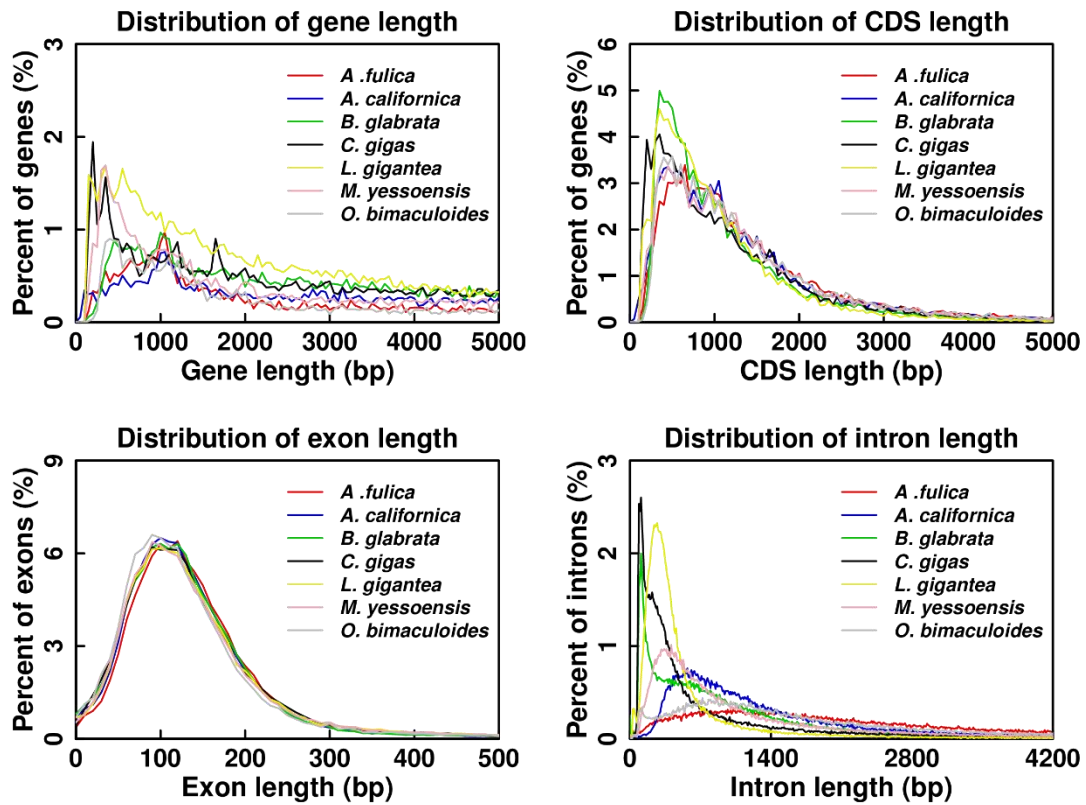
532

533

534

535

536



537
538

**Figure 3. Length distribution comparison on total gene, CDS, exon, and intron of annotated gene models of *A. fulica* with other closely related insect species.** The comparison of length distribution of genes (A), CDS (B), exon (C) and intron (D) for *A. fulica* to those in *A.californica* , *B. glabrata* , *C. gigas* , *L. gigantea* , *P. yessoensis* and *O. bimaculoides*.
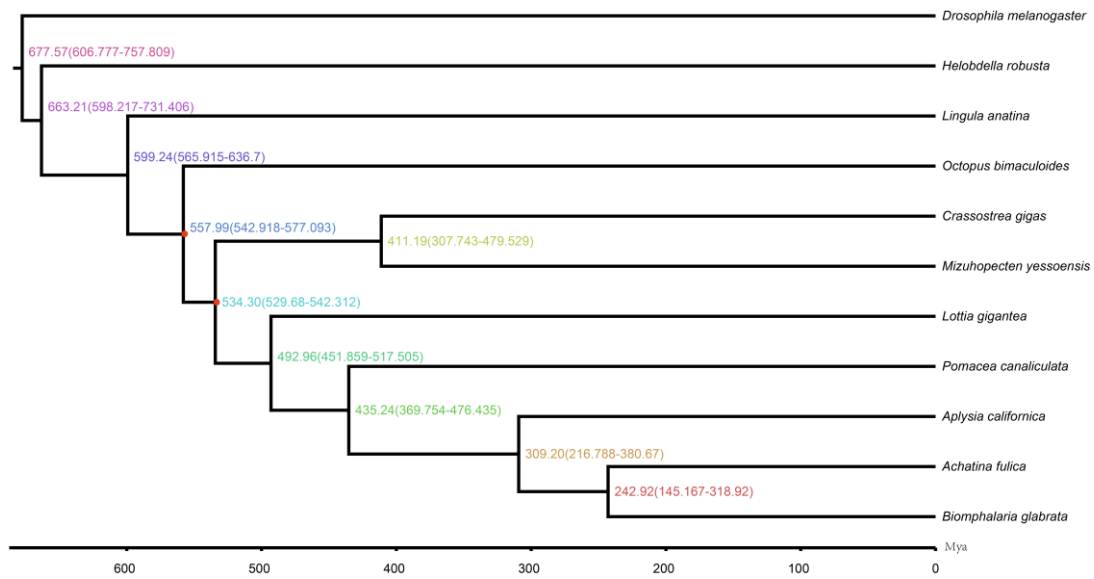
544
545
546
547
548

549

**Figure 4. Phylogenetic relationship between *A. fulica* and related species.**

The divergence time (million years ago, Mya) and the 95% confidential intervals are labeled at branch sites and the red dots in the tree denotes the fossil recalibration sites with the maximum and minimum age of Bivalve/gastropod divergence were 543 and 530 Mya, and the maximum age of Mollusk crown group divergence was 549 Mya.

Click here to access/download
**Supplementary Material**
supplementary_information.docx