

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

Manuscript Number:	GIGA-D-19-00006	
Full Title:	A chromosomal-level genome assembly for the giant African snail <i>Achatina fulica</i>	
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:	<p>Background: <i>Achatina fulica</i> (<i>A. fulica</i>), also called giant African snail, is the largest species in the reported terrestrial mollusks. Due to its greedy appetite, wide environmental adaptability, high growth rate and reproduction capacity, the species caused world-wide invasion, mainly in Southeast Asia, Japan, the western Pacific islands and China. <i>A. fulica</i> is a pest to damage the agricultural crops, as well as an intermediate host of many parasites to threaten human health. However, genomic information of <i>A. fulica</i> is still limited, hindering the genetic and genomic studies with the aim to invasion control and management of the species.</p> <p>Finding: Using Kmer-based method, we estimated the <i>A. fulica</i> genome size of 2.12 Gb with a high repeat content up to 71%. About 101.6 Gb genomic long-read data of <i>A. fulica</i> were generated from the PacBio sequencing platform and assembled to the first <i>A. fulica</i> 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 <i>A. fulica</i> separated from the common ancestor with <i>Biomphalaria glabrata</i> around 182 million years ago.</p> <p>Conclusion: As our best knowledge, the <i>A. fulica</i> genome was the first terrestrial mollusk genome reported so far. The chromosome sequences of <i>A. fulica</i> will not only provided the research community valuable genome resource for the population genetics and environmental adaptation studies for the species, as well as, for the chromosome level evolution investigation with other mollusks.</p> <p>Key Words: Giant African snail, <i>Achatina fulica</i>, PacBio, Hi-C, chromosome assembly</p>	
Corresponding Author:	ning xiao 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:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist . Information essential to interpreting the data presented should be made available in the figure legends. Have you included all the information requested in your manuscript?	Yes
Resources A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.	Yes

<p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

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 Jian³, Wu Zhongdao⁴, Li Shizhu^{1,2}, Zhou
Xiaonong^{1,2}, Hu Wei^{1,2,3,*}, Xiao Ning^{1,2,*}

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

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

³Department of Microbiology and Microbial Engineering , School of Life Sciences , Fudan University , Shanghai 200438 , China

⁴Department of Parasitology, Zhongshan School of Medicine, Sun Yat-sen University, Guangzhou 510080, China

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

26 Abstract

27 Background:

28 *Achatina fulica* (*A. fulica*), also called giant African snail, is the largest species in the
29 reported terrestrial mollusks. Due to its greedy appetite, wide environmental
30 adaptability, high growth rate and reproduction capacity, the species caused
31 world-wide invasion, mainly in Southeast Asia, Japan, the western Pacific islands and
32 China. *A. fulica* is a pest to damage the agricultural crops, as well as an intermediate
33 host of many parasites to threaten human health. However, genomic information of *A.*
34 *fulica* is still limited, hindering the genetic and genomic studies with the aim to
35 invasion control and management of the species.

36 Finding:

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

50 Conclusion:

51 As our best knowledge, the *A. fulica* genome was the first terrestrial mollusk genome
52 reported so far. The chromosome sequences of *A. fulica* will not only provided the
53 research community valuable genome resource for the population genetics and
54 environmental adaptation studies for the species, as well as, for the chromosome
55 level evolution investigation with other mollusks.

56
57
58 **Key Words:** Giant African snail, *Achatina fulica*, PacBio, Hi-C, chromosome
59 assembly

60 Data description

61 Introduction

62 The giant African snail, *A. fulica*, is a Gastropod species (**Figure 1**). It is the largest
63 interrestrial mollusks with greedy appetite, strong environmental adaptability, and high
64 growth and reproduction rate¹⁻³. Originating from East Africa, *A. fulica* gradually
65 invaded Southeast Asia, Japan and the western Pacific islands in the last century⁴⁻⁶
66 with the direct or indirect help from humans⁷⁻⁹. In mainland China, the first *A. fulica*
67 invasion event was reported in 1931¹⁰. At present, the snail's natural distribution in the
68 wild has been found in Guangdong, Hainan, Guangxi, southern parts of Yunnan
69 Province and Fujian Province, and a county of Guizhou Province¹¹. *A. fulica* was
70 included as the first 16 alien invasive species in China
71 (http://www.mee.gov.cn/gkml/zj/wj/200910/t20091022_172155.htm) in 2003, and was
72 also listed by International Union for Conservation of Nature (IUCN) as the 100 most
73 threatening alien invasive species¹². This snail has been recognized as an agricultural
74 and garden pest that has caused significant damages in both tropical and subtropical
75 regions^{9,12,13}. In addition, *A. fulica* is also the intermediate host of *Angiostrongyl*
76 *cantonensis*. Human infection with angiostrongyliasis, which occurs mainly through
77 consumption of snails carrying *A. cantonensis* larvae, causes eosinophilic
78 meningoencephalitis^{4,11,14-18}. As a consequence, *A. fulica* is attracting more and more
79 attention in fields of both agricultural crops protection and human disease control.

80 To date, a variety of mollusk genomes have been analyzed and published,
81 including two freshwater gastropods snails *Pomacea canaliculata*¹⁹ and *Biomphalaria*
82 *glabrata*²⁰. However, no genome has been reported for terrestrial mollusks. *A. fulica* is
83 considered to be one of the most serious threat and a destructive terrestrial gastropod
84 which poses a significant hazard to agriculture, the environment, biodiversity and
85 human health. In this work, we applied Illumina, PacBio and Hi-C techniques to
86 construct the chromosome of *A. fulica*. The genome is the first terrestrial mollusk

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

87 genome, providing an important reference for the molecular mechanism
88 investigations for its broad environmental adaptability and the development of control
89 strategy of the world-wide invasion.

90 **Sample and sequencing**

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

100 Using the DNA molecules from abdominal foot, a library with the insertion length
101 of 300 bp were constructed and sequenced for Illumina sequencing platform
102 according to the manufacturer's protocol. About 195.4 Gb short reads were obtained
103 from the Illumina X Ten sequencing technology (**Table 1**), which was used for the
104 following genome survey analysis, and for final base-level genome sequence
105 correction. Meanwhile, four 20 kb libraries were constructed for PacBio Sequel
106 sequencing. Using 16 sequencing SMRT cells, 101.6 Gb long reads were generated
107 (**Table 1**). The mean and N50 lengths of the polymerases for sequencing cells ranged
108 from 6.4 kb to 10.4 kb and from 12.3 kb to 20.3 kb for cells, respectively. Those long
109 genomic DNA reads were used for reference genome construction.

110 **Genome features estimation from Kmer method**

111 With sequencing data from the Illumina platform, several genome characters could be
112 evaluated from *A. fulica*. To ensure the quality of the analysis, ambiguous bases and
113 low-quality reads were trimmed and filtered using the HTQC package²¹. The following

114 quality control were performed under the framework of HTQC. First, the quality of
115 bases at two read ends were checked. Bases in sliding 5 bp windows were deleted if
116 the average quality of the window was below phred quality score of 20. Second, reads
117 were filtered if the average phred quality score were smaller than 20 or the read
118 length was shorter than 75 bp. Third, the mate reads were also removed if the
119 corresponding reads were filtered.

120 The quality-controlled reads were used for genome character estimation. We
121 calculated the number of each 17-mer from the sequencing data using the jellyfish
122 software²², and the distribution was analyzed with GCE software²³. We estimated the
123 genome size of 2.12 Gb with the heterozygosity of 0.47% and repeat content of 71%
124 in the genome. Previous studies revealed that repeat content varies in mollusks, and
125 that repeat content is correlated with genome size²⁴. The large genome size and high
126 proportion of repeat contents of *A. fulica* provided additional supporting data for the
127 statically analysis.

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

129 After removing adaptor sequences in polymerases, 101.6 Gb subreads were
130 generated for the following whole genome assembly. The average and N50 length of
131 subreads reached 5.25 kb and 8.80 kb, respectively. To optimize the genome
132 assembly using the PacBio sequencing data, we applied two packages in the
133 assembly process, Canu²⁵ and FALCON²⁶. Canu package was first applied for the
134 assembly with the default parameters. As a result, a 1.93 Gb genome was
135 constructed with 10,417 contigs and a contig N50 length of 662.40 kb. FALCON was
136 also employed using the length_cutoff and pr_length_cutoff parameters of 10 kb and
137 8 kb, respectively. We obtained 1.85 Gb genome with 8,585 contigs, with a contig N50
138 of 726.63 kb. We adopted the FALCON assembly as the reference genome for *A.*
139 *fulica* (**Table 2**). The genome sequences were subsequently polished by PacBio long
140 reads using arrow²⁷ and Illumina short reads by pilon²⁸ to correct base errors. The

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

141 corrected genome was further applied for the following chromosome assembly
142 construction using Hi-C data.

143 ***In situ* Hi-C library construction and chromosome assembly using Hi-C**
144 **data**

145 Liver pancreas tissue of *A. fulica* was used for library construction for Hi-C analysis
146 and the library was constructed using the identical method in previous studies²⁹.
147 Finally, the library was sequenced with 150 paired-end mode on the Illumina HiSeq X
148 Ten platform (San Diego, CA, United States). From the Illumina sequencing platform,
149 1,313 million paired-end reads were obtained for the Hi-C library (**Table 1**). The reads
150 were mapped to the above *A. fulica* genome with Bowtie³⁰, with two ends of paired
151 reads being mapped to the genome separately. To increase the interactive Hi-C reads
152 ratio, an iterative mapping strategy was performed as previous studies, and only read
153 pairs that both ends uniquely mapped were used for the following analysis. From the
154 alignment status of two ends, self-ligation, non-ligation and other sorts of invalid reads,
155 including StartNearRsite, PCR amplification, random break, LargeSmallFragments
156 and ExtremeFragments, were filtered out by Hi-Clib³¹. Through the recognition of
157 restriction sites in sequences, contact counts among contigs were calculated and
158 normalized.

159 According to previous karyotype analyses, *A. fulica* had 31 chromosomes³². By
160 clustering the contigs using the contig contact frequency matrix, we were able to
161 correct some minor errors in the FALCON assembly results. Contigs with errors were
162 broken into shorter contigs. We obtained 8,701 contigs, slightly more than the 8,585
163 contigs in the FALCON assembly. We successfully clustered these contigs into 31
164 groups in Lachesis³³ using the agglomerative hierarchical clustering method (**Figure**
165 **2**). Lachesis was further applied to order and orient the clustered contigs according to
166 the contact matrix. As a result, 7,106 contigs were reliably anchored, ordered and
167 orientated on chromosomes, accounting for 99.32% of the total genome bases. Then,

168 we applied PBJelly³⁴ to fill the gap using PacBio long reads to merge the contig
169 sequences. Finally, the first chromosomal-level assembly of *A. fulica* was obtained
170 with 8,211 contigs, a contig N50 of 721.0 kb and a scaffold N50 of 59.59 Mb (**Table 2**
171 and **Figure 3**).

172 **Genome quality evaluation**

173 We assessed the quality of genome of *A. fulica* after the assembly process. The
174 quality evaluation was carried out in three aspects: continuity, completeness and base
175 level accuracy.

176 First of all, we compared the sequence number and N50 length of contig of *A.*
177 *fulica* with public genome of mollusks and found that our assembly has a high quality
178 on contig and scaffold N50 among mollusk genomes. (**Figure 3**) As previous studies,
179 genomic heterozygosity of mollusk was one of the biggest challenges for genome
180 assembly, both in terms of contig and scaffold assembly³⁵. Our work illustrated that
181 the genome assembly using PacBio long sequencing data was affordable and
182 effective to overcome the difficulty of mollusk genome assembly. Traditional
183 chromosomal genome assembly requires physical maps and genetic maps, which is
184 enormously time-and labor-consuming. With Hi-C data analysis, we successfully
185 assembled *A. fulica* genome into chromosome-level with just one individual.

186 Second, the assembled genome was subjected to the BUSCO (version 3.0)³⁶ to
187 assess the completeness of the genome. 91.7% of the BUSCO genes were identified
188 in *A. fulica* genome. More than 84.7% BUSCO gene were single-copy completed in
189 our genome, illuminating a high level of completeness of the genome.

190 Third, NGS short reads were aligned to the genome using BWA package³⁷. About
191 98.7% of paired reads were aligned to the genome, of which 98.24% were reads
192 paired aligned. From the NGS reads alignment, we detected 128,998 homologous
193 SNP loci using the GATK pipeline³⁸, demonstrating the high base-level accuracy of
194 99.33%.

195 **Repeat element and gene annotation**

1
2
3 196 Tandem Repeat Finder (TRF)³⁹ was used for repetitive element identification in the *A.*
4
5 197 *fulica* genome. A *de novo* method applying RepeatModuler was used to detect
6
7 198 transposable elements (TEs). The resulted *de novo* data, combined with known
8
9 199 repeat library from Repbase⁴⁰, were used to identify TEs in the *A. fulica* genome by
10
11 200 RepeatMasker⁴¹ software. All repetitive elements were masked in the genome for the
12
13 201 protein-coding gene prediction.

14
15
16 202 Protein-coding genes in the *A. fulica* genome were annotated using the *de novo*
17
18 203 program Augustus⁴². Protein sequences of the closely related species including
19
20 204 *Aplysia californica*, *Biomphalaria glabrata*, *Crassostrea gigas*, *Lottia gigantea* and
21
22 205 *Patinopecten yessoensis*, were downloaded from the Ensembl database, and aligned
23
24 206 to the *A. fulica* genome with TBLASTN. Full-length transcripts obtained using Iso-Seq
25
26 207 were mapped to the genome using Genewise⁴³. Finally, gene models predicted from
27
28 208 all above methods were combined by MAKER⁴⁴, resulting in 23,726 protein-coding
29
30 209 genes. The gene number, gene length, CDS length, exon length and intron length
31
32 210 distribution were all comparable with the related mollusks (**Figure 4**).

33
34
35
36
37 211 To functionally annotate protein-coding genes in the *A. fulica* genome, we
38
39 212 searched all predicted gene sequences to NCBI non-redundant nucleotide (NT) and
40
41 213 protein (NR), Swiss-Prot databases by BLASTN⁴⁵ and BLASTX⁴⁶ utility. Blast2GO⁴⁷
42
43 214 was also used to assign gene ontology (GO)⁴⁸ and Kyoto Encyclopedia of Genes and
44
45 215 Genomes (KEGG)⁴⁹ pathways. A threshold of e-value of 1e-5 was used for all BLAST
46
47 216 applications. Finally, 22,858 (96.34%) genes were functionally annotated (**Table 3**).

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

52
53 218 OrthoMCL⁵⁰ was used to cluster gene families. First, proteins from *A. fulica* and the
54
55 219 closely related mollusks, including *Aplysia californica*, *Biomphalaria glabrata*,
56
57 220 *Crassostrea gigas*, *Lingula anatina*, *Lottia gigantea*, *Patinopecten yessoensis*,
58
59 221 *Octopus bimaculoides*, *Helobdella robusta*, *Drosophila melanogaster* and *Pomacea*
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

222 *canaliculata* were all-to-all blasted by BLASTP⁴⁶ utility with an e-value threshold of
223 1e-5. Only proteins from the longest transcript were used for genes with alternative
224 splices. We identified 25,448 gene families for *A. fulica* and the related species,
225 among them 675 single-copy orthologs families were detected.

226 Using single-copy orthologs, we could probe the phylogenetic relationships for
227 the *A. fulica* and other mollusks. To this end, protein sequences of single-copy genes
228 were aligned using CLUSTALX⁵¹. Guided by the protein multi-sequence alignment,
229 the alignment of the coding DNA sequences (CDS) for those genes were generated
230 and concatenated for the following analysis. The phylogenetic relationships were
231 constructed using PhyML⁵² using the concatenated nucleotide alignment with the
232 JTT+G+F model. The PAML MCMCtree program⁵² was used to estimate the species
233 divergent time scales for the mollusks using approximate likelihood method. We found
234 that *A. fulica* was most closely related to *Biomphalaria glabrata*, and the two species
235 diverged from their common ancestor around 177.1-187.1 million years ago (MYA)
236 **(Figure 5)**.

237 **Conclusion**

238 We reconstructed the first chromosome level assembly for *A. fulica* using an
239 integrated strategy of PacBio, Illumina and Hi-C technologies. Using the long reads
240 from PacBio Sequel platform and short reads from the Illumina X Ten platform, we
241 successfully constructed contig assembly for *A. fulica*. Leveraging contact information
242 among contigs from Hi-C technology, we further improved the assembly to the
243 chromosome-level quality **(Figure 2 and Figure 3)**. We annotated 23,726
244 protein-coding genes in the *A. fulica* genome and 22,858 of genes were functionally
245 annotated. With 675 single-copy orthologs from *A. fulica* and other related mollusks,
246 we construct the phylogenetic relationship of these mollusks, and found that *A. fulica*
247 might have diverged from its common ancestor of *Biomphalaria glabrata* around
248 177.1-187.1 MYA. Given the increasing interests in mollusk genomic evolution and

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

249 the biological importance of *A. fulica* as an invasive animal, our genomic and
250 transcriptome data provide valuable genetic resource for the following functional
251 genomics investigations for the research community.

252

253 **Ethics Statement**

254 This study was approved by the Animal Care and Use committee of National Institute
255 of Parasitic Diseases, Chinese Center for Disease Control and Prevention. All
256 participants consent the study under the 'Ethics, consent and permissions' heading.
257 All participants consent to publish the work under the 'Consent to publish' heading.

258 **Availability of supporting data**

259 The Illumina, PacBio and Hi-C sequencing data are available from NCBI via the
260 accession number of SRR8369706, SRR8369311 and SRR8371669, respectively.
261 The Illumina transcriptome sequencing data were deposited to NCBI via the
262 accession number of SRR8371872 and SRR8371873. The genome, annotation and
263 intermediate files were uploaded to GigaScience FTP server.

264 **Competing interests**

265 The authors declare that they have no competing interests.

266 **Acknowledgement**

267 This work was supported by the National Key Research and Development Program of
268 China (No. 2016YFC1200500 and 2016YFC1202000). The authors thank Frasergen
269 Bioinformatics for providing technique supports for this work.

270 **Author Contributions**

271 Z.X, H.W and X.N conceived the project. G.Y, Z.Y, L.Q collected the samples and
272 extracted the genomic DNA. G.Y, Z.Y and L.Q performed the genome assembly and
273 data analysis. G.Y, Z.X, H.W and X.N wrote the paper.

274 **References**

- 1
2 275 1 Schreurs, J. Investigations on the biology, ecology and control of Giant African Snail 290 in
3 276 West New Guinea. (Manokwari Agricultural Research Station., 1963).
4
5 277 2 Albuquerque, F. S., Peso-Aguiar, M. C. & Assunção-Albuquerque, M. J. Distribution, feeding
6 278 behavior and control strategies of the exotic land snail *Achatina fulica*
7 279 (Gastropoda: Pulmonata) in the Northeast of Brazil. *Braz.J.Biol.* **68**, 6 (2008).
8
9 280 3 Thiengo, S. C., Fernandez, M. A., Torres, E. J., Coelho, P. M. & Lanfredi, R. M. First record of
10 281 an nematode Metastrongyloidea (Aelurostrongylus abstrusus larvae) in *Achatina* (Lissachatina)
11 282 *fulica* (Mollusca, Achatinidae) in Brazil. *J. Invertebr. Pathol.* **98**, 6 (2008).
12
13 283 4 Lv, S., Zhang, Y. & Liu, H. X. Invasive Snails and an Emerging Infectious Disease: Results from
14 284 the First National Survey on Angiostrongylus cantonensis in China. *BioOne*,
15 285 doi:10.1371/journal.pntd.0000368 (2009).
16
17 286 5 Cowie, R. H. Non-indigenous land and freshwater molluscs in the islands of the Pacific:
18 287 Conservation impacts and threats. 143-172 (2000).
19
20 288 6 Cowie, R. H. Can snails ever be effective and safe biocontrol agents? *Int J Pest Manage* **47**, 18
21 289 (2001).
22
23 290 7 Cowie, R. H. & Robinson, D. G. *Pathways of introduction of nonindigenous land and*
24 291 *freshwater snails and slugs.* 93-122 (Island Press, 2003).
25
26 292 8 Kotangale, J. P. Giant African snail (*Achatina fulica* Bowdich). *J Environ Sci Eng*, 6 (2011).
27 293 9 Raut, S. K. & Barker, G. M. *Achatina fulica* Bowdich and Other Achatinidae as Pests in Tropical
28 294 Agriculture. (CABI International, 2002).
29
30 295 10 Jarreit, V. H. C. The spread of the snail *Achatina fulica* to south China. *Hong Kong Nat* **2**, 3
31 296 (1931).
32
33 297 11 Shan, L., Yi, Z. & Peter, S. Emerging Angiostrongyliasis in Mainland China. *Emerging Infectious*
34 298 *Diseases* **14**, 4 (2008).
35
36 299 12 Lowe, S., Browne, S. M., Boudjrlas, S. & De Poorter, M. *100 of the world's worst invasive alien*
37 300 *species: A selection from the global invasive species database. The Invasive Species Specialists*
38 301 *Group of the Species Survival Commission of the World Conservation Union.*, (Hollands
39 302 Printing, 2000).
40
41 303 13 Mead, A. R. *Pulmonates volume 2B. Economic malacology with particular reference to*
42 304 *Achatina fulica.* (Academic Press, 1979).
43
44 305 14 Alicata, J. E. The discovery of Angiostrongylus cantonensis as a cause of human eosinophilic
45 306 meningitis. *Parasitol Today* **7**, 151-153 (1991).
46
47 307 15 Prociw, P., Spratt, D. M. & Carlisle, M. S. Neuro-angiostrongyliasis: unresolved issues. *Int J*
48 308 *Parasitol* **30**, 1295-1303 (2000).
49
50 309 16 Deng, Z. H., Zhang, Q. M., Huang, S. Y. & Jones, J. L. First provincial survey of Angiostrongylus
51 310 cantonensis in Guangdong Province, China. *Trop Med Int Health* **17**, 4 (2012).
52
53 311 17 Maldonado, J. A. *et al.* First report of Angiostrongylus cantonensis (Nematoda:
54 312 Metastrongylidae) in *Achatina fulica* (Mollusca: Gastropoda) from Southeast and South
55 313 Brazil. *Mem Inst Oswaldo Cruz* **105**, 4 (2010).
56
57 314 18 Vitta, A., Polseela, R., Nateeworanart, S. & Tattiyapong, M. Survey of Angiostrongylus
58 315 cantonensis in rats and giant African land snails in Phitsanulok Province, Thailand. *Asian Pac J*
59 316 *Trop Med* **4**, 3 (2011).
60
61 317 19 Liu, C. *et al.* The genome of the golden apple snail *Pomacea canaliculata* provides insight into

1 318 stress tolerance and invasive adaptation. *Gigascience* **7**, doi:10.1093/gigascience/giy101
2 319 (2018).
3 320 20 Adema, C. M. *et al.* Whole genome analysis of a schistosomiasis-transmitting freshwater
4 321 snail. *Nature communications* **8**, 15451, doi:10.1038/ncomms15451 (2017).
5 322 21 Neff, K. L. *et al.* Mojo Hand, a TALEN design tool for genome editing applications. *BMC*
6 323 *Bioinformatics* **14**, 1, doi:10.1186/1471-2105-14-1 (2013).
7 324 22 Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of
8 325 occurrences of k-mers. *Bioinformatics* **27**, 764-770, doi:10.1093/bioinformatics/btr011
9 326 (2011).
10 327 23 Binghang Liu, Y. S., Jianying Yuan, Xuesong Hu, Hao Zhang, Nan Li, Zhenyu Li, Yanxiang
11 328 Chen, Desheng Mu, Wei Fan. Estimation of genomic characteristics by analyzing k-mer
12 329 frequency in de novo genome projects. *Quantitative Biology* **35**, 62-67 (2013).
13 330 24 Murgarella, M. *et al.* A First Insight into the Genome of the Filter-Feeder Mussel *Mytilus*
14 331 *galloprovincialis*. *PLoS One* **11**, e0151561, doi:10.1371/journal.pone.0151561 (2016).
15 332 25 Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting
16 333 and repeat separation. *Genome Res* **27**, 722-736, doi:10.1101/gr.215087.116 (2017).
17 334 26 Chin, C. S. *et al.* Phased diploid genome assembly with single-molecule real-time sequencing.
18 335 *Nat Methods* **13**, 1050-1054, doi:10.1038/nmeth.4035 (2016).
19 336 27 Chin, C.-S. *et al.* Nonhybrid, finished microbial genome assemblies from long-read SMRT
20 337 sequencing data. *Nature methods* **10**, 563 (2013).
21 338 28 Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection
22 339 and genome assembly improvement. *PLoS one* **9**, e112963 (2014).
23 340 29 Gong, G. *et al.* Chromosomal-level assembly of yellow catfish genome using third-generation
24 341 DNA sequencing and Hi-C analysis. *Gigascience*, doi:10.1093/gigascience/giy120 (2018).
25 342 30 Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment
26 343 of short DNA sequences to the human genome. *Genome Biol* **10**, R25,
27 344 doi:10.1186/gb-2009-10-3-r25 (2009).
28 345 31 Burton, J. N. *et al.* Chromosome-scale scaffolding of de novo genome assemblies based on
29 346 chromatin interactions. *Nature biotechnology* **31**, 1119 (2013).
30 347 32 Sun, T. Chromosomal studies in three land snails. *Sinozoologia* **12**, 154-162 (1995).
31 348 33 Near, T. J. *et al.* Phylogeny and tempo of diversification in the superradiation of spiny-rayed
32 349 fishes. *Proceedings of the National Academy of Sciences of the United States of America* **110**,
33 350 12738 (2013).
34 351 34 English, A. C. *et al.* Mind the gap: upgrading genomes with Pacific Biosciences RS long-read
35 352 sequencing technology. *PLoS one* **7**, e47768 (2012).
36 353 35 Zhang, G. *et al.* The oyster genome reveals stress adaptation and complexity of shell
37 354 formation. *Nature* **490**, 49 (2012).
38 355 36 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO:
39 356 assessing genome assembly and annotation completeness with single-copy orthologs.
40 357 *Bioinformatics* **31**, 3210-3212 (2015).
41 358 37 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform.
42 359 *bioinformatics* **25**, 1754-1760 (2009).
43 360 38 McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
44 361 next-generation DNA sequencing data. *Genome research* (2010).

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

362 39 Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res*
363 **27**, 573-580 (1999).

364 40 Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in
365 eukaryotic genomes. *Mob DNA* **6**, 11, doi:10.1186/s13100-015-0041-9 (2015).

366 41 Chen, N. Using RepeatMasker to identify repetitive elements in genomic sequences. *Current*
367 *protocols in bioinformatics* **5**, 4.10. 11-14.10. 14 (2004).

368 42 Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic acids*
369 *research* **34**, W435-W439 (2006).

370 43 Birney, E., Clamp, M. & Durbin, R. GeneWise and genomewise. *Genome research* **14**, 988-995
371 (2004).

372 44 Cantarel, B. L. *et al.* MAKER: an easy-to-use annotation pipeline designed for emerging model
373 organism genomes. *Genome research* **18**, 188-196 (2008).

374 45 Gertz, E. M., Yu, Y. K., Agarwala, R., Schaffer, A. A. & Altschul, S. F. Composition-based
375 statistics and translated nucleotide searches: improving the TBLASTN module of BLAST. *BMC*
376 *Biol* **4**, 41, doi:10.1186/1741-7007-4-41 (2006).

377 46 Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421,
378 doi:10.1186/1471-2105-10-421 (2009).

379 47 Conesa, A. *et al.* Blast2GO: a universal tool for annotation, visualization and analysis in
380 functional genomics research. *Bioinformatics* **21**, 3674-3676 (2005).

381 48 Consortium, G. O. The Gene Ontology (GO) database and informatics resource. *Nucleic acids*
382 *research* **32**, D258-D261 (2004).

383 49 Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids*
384 *research* **28**, 27-30 (2000).

385 50 Li, L., Stoeckert, C. J. & Roos, D. S. OrthoMCL: identification of ortholog groups for eukaryotic
386 genomes. *Genome research* **13**, 2178-2189 (2003).

387 51 Thompson, J. D., Gibson, T. J. & Higgins, D. G. Multiple sequence alignment using ClustalW
388 and ClustalX. *Current protocols in bioinformatics*, 2.3. 1-2.3. 22 (2003).

389 52 Guindon, S., Lethiec, F., Duroux, P. & Gascuel, O. PHYML Online—a web server for fast
390 maximum likelihood-based phylogenetic inference. *Nucleic acids research* **33**, W557-W559
391 (2005).

392
393

394

395 **Tables and Figures**

396

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

Source	Library type	Platform	Library size (bp)	Data size (Gb)	Application
Genome	Short reads	HiSeq X Ten	350	195.4	Genome survey and base correction
	Long reads	PacBio SEQUEL	20,000	101.6	Genome assembly
	Hi-C	HiSeq X Ten	300-500	208.9	Chromosome construction
Transcriptome	Long reads	PacBio SEQUEL	3000, 5000	22.5	Genome annotation

398

399

400

401 **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

402

403

404 **Table 3: 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	

405

406

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

407

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



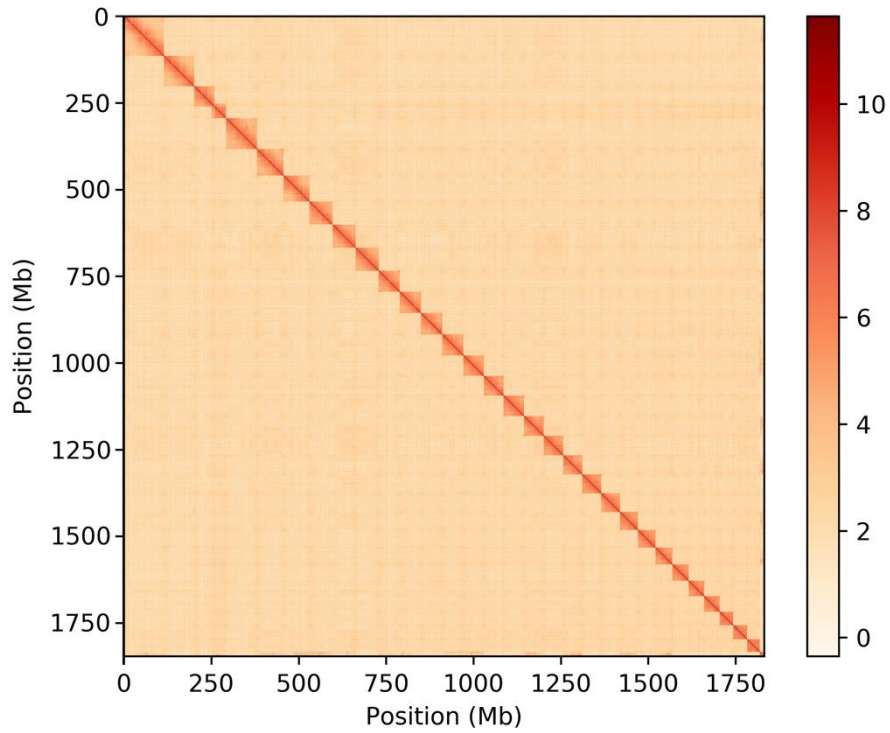
408

409 **Figure 1. A picture of *A. fulica* that used for genome sequencing and assembly.**

410

411

412



413

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

417

418

419

420

421

422

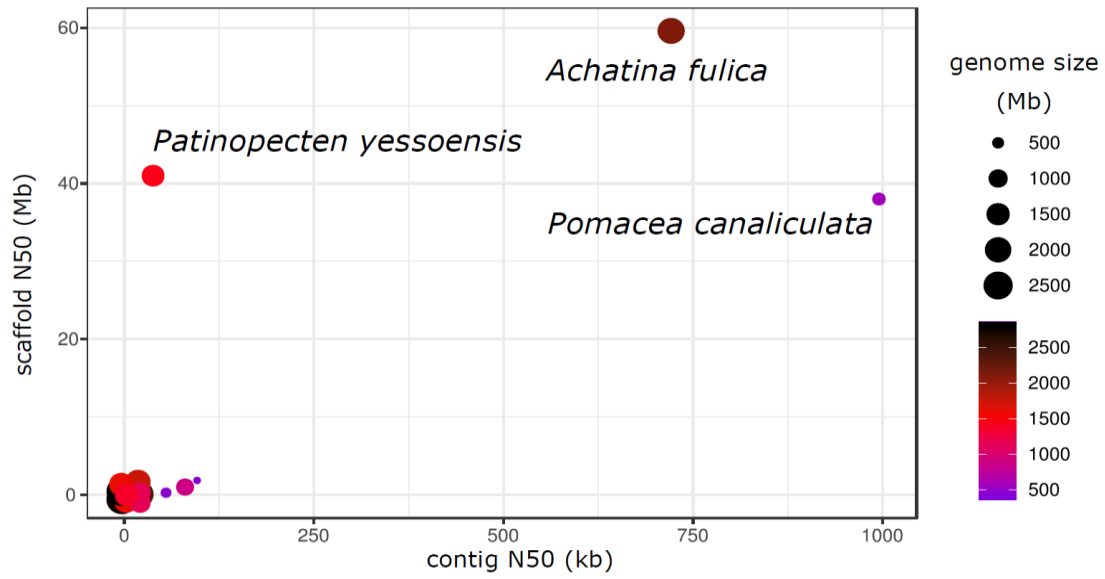
423

424

425

426

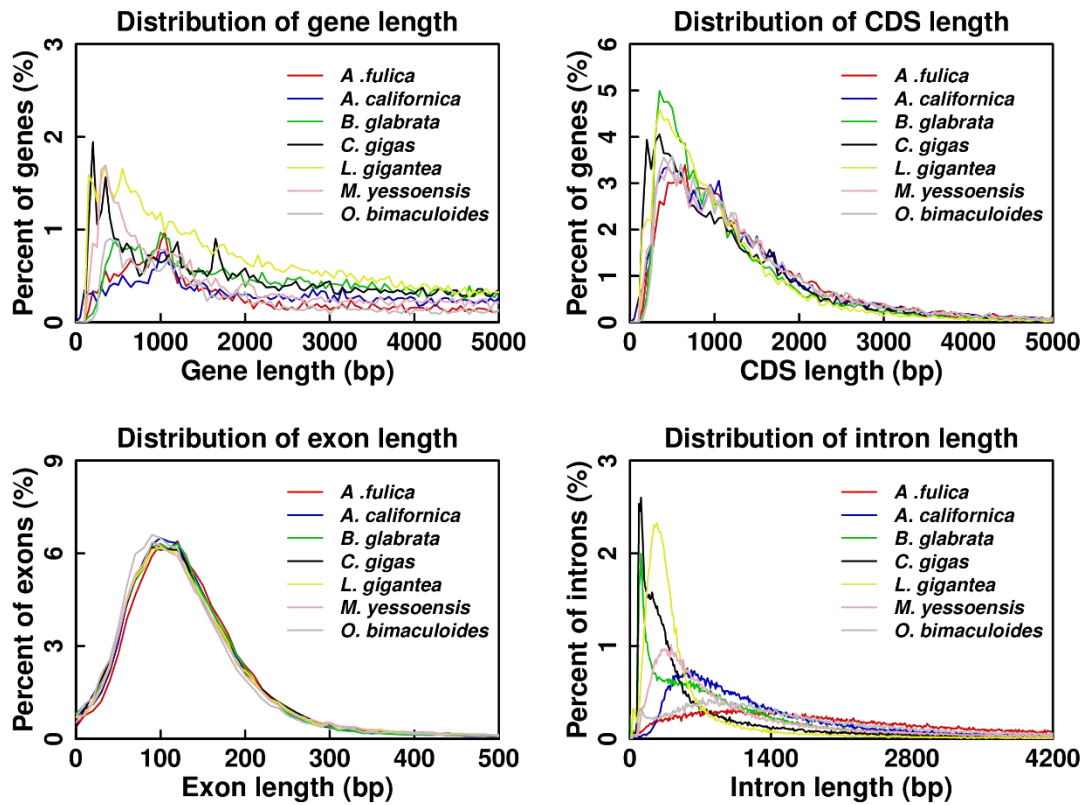
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



427
428
429
430
431

Figure 3: Genome assembly comparison of *A. fulica* with other sequenced mollusk genomes. The x- and y-axis represent the contig and scaffold N50s, respectively. The genomes assembled into chromosomal level are labeled with names.

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



433
 434
 435 **Figure 4. Length distribution comparison on total gene, CDS, exon, and intron of**
 436 **annotated gene models of *A. fulica* with other closely related species.** The
 437 comparison of length distribution of genes (A), CDS (B), exon (C) and intron (D) for *A.*
 438 *fulica* to those in *A.californica* , *B. glabrata* , *C. gigas* , *L. gigantea* , *P. yessoensis* and *O.*
 439 *bimaculoides*.

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

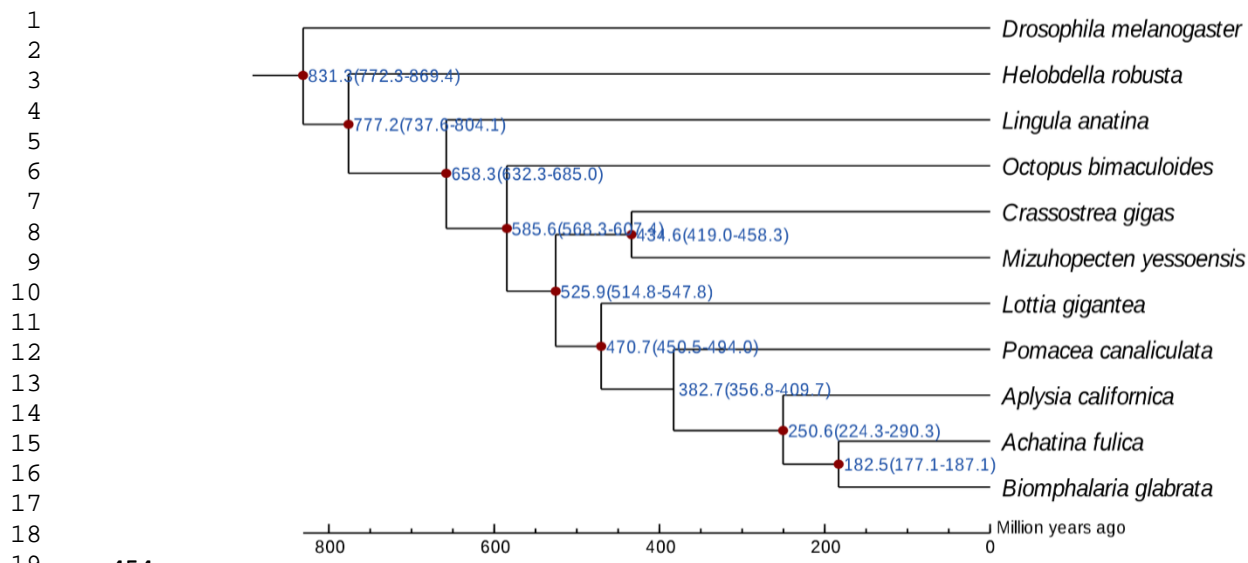
455

456

457

458

459



454

455 **Figure 5. Phylogenetic relationship between *A. fulica* and related species.** The
 456 divergence time (million years ago) and the 95% confidential intervals are labeled at
 457 branch sites and the red dots in the tree illuminated the speciation for the time
 458 recalibration.

459

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