# GigaScience

## The genome of golden apple snail Pomacea canaliculata provides insight into stress tolerance and invasive adaptation

### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-18-00030R2 |
| Full Title: | The genome of golden apple snail Pomacea canaliculata provides insight into stress tolerance and invasive adaptation |
| Article Type: | Research |

| Abstract: | Background: The golden apple snail (Pomacea canaliculata) is a fresh water snail listed among the top-100 worst invasive species, worldwide and a noted agricultural and quarantine pest that causes great economic losses. It is characterized by fast growth, strong stress tolerance, a high reproduction rate, and adaptation to a broad range of environments.<br>Results: Here, we used long-read sequencing to produce a 440-Mb high-quality chromosome-level assembly for the P. canaliculata genome. In total, 50 Mb (11.4%) repeat sequences and 21,533 gene models were identified in the genome. The major findings of this study include the recent explosion of DNA/hAT-Charlie transposable elements (TEs), the expansion of the P450 gene family and the constitution of the cellular homeostasis system, which contributes to ecological plasticity in stress adaptation. In addition, the high transcriptional levels of perivitellin genes in the ovary and albumen gland promote the function of nutrient supply and defence ability in eggs. Furthermore, the gut metagenome also contains diverse genes for food digestion and xenobiotic degradation.<br>Conclusions: These findings collectively provide novel insights into the molecular mechanisms of the ecological plasticity and high invasiveness. |
|---|---|

| Corresponding Author: | Wei Fan<br>Chinese Academy of Agricultural Sciences<br>CHINA |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | Chinese Academy of Agricultural Sciences |
| Corresponding Author's Secondary Institution: | |
| First Author: | Conghui Liu |
| First Author Secondary Information: | |
| Order of Authors: | Conghui Liu |
| | Yan Zhang |
| | Yuwei Ren |
| | Hengchao Wang |
| | Shuqu Li |

| | Fan Jiang |
| --- | --- |
| | Lijuan Yin |
| | Xi Qiao |
| | Guojie Zhang |
| | Wanqiang Qian |
| | Bo Liu |
| | Wei Fan |

| Order of Authors Secondary Information: | |
| --- | --- |
| Response to Reviewers: | Reviewer reports:<br>Reviewer #1: Still there are some typos in the revised manuscript.<br>For example, in line 181 "Pintada" must be "Pinctada".<br>Please carefully check the manuscript again before submission.<br>Reply: We have revised all the typos in the manuscript. We revised "Lottia giganta" to "Lottia gigantea" in line 161, "Pintada fucata" to "Pinctada fucata" in line 162, "giganta" to "gigantea" in line 166, "orthfinder" to "orthoFinder" in line 169, "L. fortune" to "L. fortunei" in line 172, "L. giganta" to "L. gigantea" in line 259, "Lottia giganta" to "Lottia gigantea" in line 462, "Pintada fucata" to "Pinctada fucata" in line 463, "giganta" to "gigantea" in line 476, "L. fortune" to "L. fortunei" in line 543, "L. giganta, Lottia giganta" to "L. gigantea, Lottia gigantea" in line 543, "L. giganta" to "L. gigantea" in Table 1, "L. giganta" to "L. gigantea" in the legend of Figure 4.<br><br>Reviewer #3: Dear authors,<br><br>Thank you for providing a revised version of the manuscript and for addressing my suggestions. I think this manuscript will be a great contribution for the genomic studies of mollusks and invasive species. I, however, still have a few comments.<br><br>1-) The written English is much improved, but there are still a few persistent mistakes. Such as "L. giganta" where it should be 'L. gigantea' and the same with "L. fortune" which is actually 'L. fortunei'.<br>Reply: We have revised "L. giganta" to "L. gigantea", and "L. fortune" to "L. fortunei" in the manuscript. We revised "Lottia giganta" to "Lottia gigantea" in line 161, "giganta" to "gigantea" in line 166, "L. fortune" to "L. fortunei" in line 172, "L. giganta" to "L. gigantea" in line 259, "Lottia giganta" to "Lottia gigantea" in line 462, "giganta" to "gigantea" in line 476, "L. fortune" to "L. fortunei" in line 543, "L. giganta, Lottia giganta" to "L. gigantea, Lottia gigantea" in line 543, "L. giganta" to "L. gigantea" in Table 1, "L. giganta" to "L. gigantea" in the legend of Figure 4.<br><br>I've attached again a manuscript with some purple highlights of critical pieces of text that should be revised. For example, the sentence between lines 479-484 is too long and non-technical. The same for "With its easy acquisition" in line 377. The improvement of such sentences would greatly benefit the manuscript readers.<br>Reply: (1) The long and non-technical sentence between lines 479-484 "Raw reads were cleaned to exclude adapter sequences, low-quality sequences, and contaminated DNA. The adapter sequence was identified and trimmed from the reads by an ungapped dynamic programming algorithm; the low-quality part (head or tail) of the reads was trimmed off to ensure that the average error rate of the remaining reads was lower than 0.001; the reads that were mapped to contaminated DNA by BWA-MEM were filtered out…" has been revised to short sentences, with the non-technical description removed and the applied in-house program cited:<br>"The Illumina raw reads were filtered by trimming the adapter sequence and low-quality regions (https://github.com/fanagislab/common_use), resulting in high-quality reads with an average error rate < 0.001. Then, the reads mapped to the following genomes by BWA-MEM were filtered out (https://github.com/fanagislab/metagenome_analysis.git), to exclude the contaminated host, food, parasite, and human DNA sequences …"<br>(2) The "With its easy acquisition" in line 377 has been revised to "With wide distribution", and the whole sentence became: "With wide distribution, rapid growth and efficient reproduction, P. canaliculata possesses the potential to be a model organism |

of Mollusca."

(3) The "orthologue groups" in line 170 has been revised to "orthologous groups".

(4) The "maintains" in line 238 has been revised to "contributes to", and the whole sentence became: "Apoptosis is a process of cell death when sensing stress, and the regulation of apoptosis contributes to the dynamic homeostasis of the internal environment."

(5) The sentence between lines 319-322 "The gut microbiome is well known as the second genome of animals and plays important roles in food digestion, immune defence, and other processes that are essential to the animal host. To investigate whether the gut microbiome influences the invasive lifestyle" has been improved to: "The gut microbiome is regarded as the "second genome" of the host animal, due to the fact that the gut microbiota contributes to the food digestion, immune system development, and many other processes important to the host. To investigate the relationship between the gut microbiome and the invasive lifestyle of P. canaliculata."


Also, the final subtopic should not be "Conclusion and Discussion", at that point, I would say, its time to just conclude.
Reply: We have deleted "and Discussion" in the subtopic.

In the results sections, however, many paragraphs start with a discussion of the literature instead of presenting the results: I would advice to revise those, present results first in the paragraphs and then discuss them. Again, coherence benefit readers a great deal.
Reply: Yes, we agree that results should be presented in front of discussion. To make it easier to understand for the readers, the sentences in the head of these paragraphs are brief background information, not discussion on the results. Real discussions are put after the results, in the end of the paragraphs.

2-) The amount of data generated is one of the strongest points of the work presented. And specially because of that, a great deal of analysis can be performed. For example, as you have 60x coverage of PacBio data for the snail, I would suggest running the Falcon and Falcon-Unzip pipeline to actually phase the genome: separate the haplotypes, instead of trying to merge or just through away the variation, as described in lines 424-432. The high heterozygosity described for the species actually helps in the phasing of haplotypes: there are several manuscripts describing methods to do so. I would run FALCON and FALCON-unzip, then I would polish with Illumina and try filling gaps with it in the different haplotypes and then would use the Hi-C data. I know its a great deal of analysis and highly experimental, so I'll leave it as a suggestion. But I would be interested in having a supplementary material with the imperfect alternate contigs generated by the phasing. This is the kind of information that were almost impossible to obtain with the generation of short reads, but now the long-reads technologies allow us to phase some long genome portions, and this is a very valuable information to some of us. With that, we can start understanding how much variation there are - and what are their evolutionary implications - in coding and non-coding regions within a genome.

Reply: Assembly the two haplotype chromosomes with long-reads is a very good suggestion, and we agree that the phased chromosomal sequences have greater value than the current mosaic reference genome sequence. In fact, we have run both SMARTdenovo and Falcon/Falcon-unzip, and polished by Pilon with illumina reads. The biggest difference of SMARTdenovo from Falcon is that SMRTdenovo does not need to correct sequencing errors in the first step, but instead perform an overlap-layout-consensus algorithm directly. With algorithms improved in many aspects, SMARTdenovo can achieve good assembly results with moderate sequencing coverage (50 X), in contrast, Falcon usually needs higher sequencing coverage ( 100 X) to get a good assembly. In this study, using the 60 X apple snail Pacbio data, SMARTdenovo generates contigs with N50 length over 1 megabases, which is 4 times of that of Falcon/Falcon-unzip (240 Kb).

The comparison between SMARTdenovo and Falcon/Falcon-unzip assemblies showed that contigs assembled by SMARTdenovo had the assembly size of 473.04 Mb, N50 size of 1010.40 Kb and N90 size of 172.34 Kb; primary contigs assembled by Falcon had the assembly size of 475.28 Mb, N50 size of 241.14 Kb and N90 size of 54.29 Kb; alternate contigs assembled by Falcon had the assembly size of 54.10 Mb, N50 size of 43.88 Kb and N90 size of 22.68 Kb; primary contigs assembled by Falcon-unzip had the assembly size of 474.23 Mb, N50 size of 246.62 Kb and N90 size of 58.36 Kb; haplotigs assembled by Falcon-unzip had the assembly size of 173.15 Mb, N50 size of 48.98 Kb and N90 size of 17.44 Kb.

Considering that Hi-C contains extremely long-range linkage information, the larger contig length is an import factor for the success application of Hi-C data for scaffolding. Therefore, we adopted the SMARTdenovo contigs and then applied Hi-C to get a chromosomal-scale scaffold sequence.
To make the phasing information available to the public, we also uploaded the SMARTdenovo alternate sequences excluded from the reference haploid genome sequence, as well as the Falcon-unzip assembly of the apple snail, to the GigaDB and our institution's ftp-site, respectively.

3-) About the expansions found between the snail and L. fortunei, could you please describe the methodology used to consider genes expanded in these two groups? Was this done in a comparative manner with other species? Which ones? What was the criteria to consider gene families expanded?
Reply: we added the sentence at method part "To identify the common expanded gene families, we compared the P. canaliculata and L. fortunei with other seven species. The gene number of orthologous group in P. canaliculata and L. fortunei were two or more times than that in all of other species, respectively. Additionally, these gene families with P-value less than 0.01 were considered as expansion by z-test."

3a-) Have you identified CPYs expanded in both invasive species? I would suggest that L. fortunei should be included in figure 4.
Reply: We have identified the CYP genes in the L. fortunei in the revised manuscript, which were included in Figure 4a. There were 115 CYP genes found in L. fortunei, with no obvious expansion.

| 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 | Yes |

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.

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 |

1 **The genome of the golden apple snail *Pomacea canaliculata* provides insight into**

2 **stress tolerance and invasive adaptation**

3 Conghui Liu[1]*, Yan Zhang[1]*, Yuwei Ren[1]*, Hengchao Wang[1], Shuqu Li[1], Fan Jiang[1], Lijuan Yin[1],

4 Xi Qiao[1], Guojie Zhang[2], Wanqiang Qian[1], Bo Liu[1]†, Wei Fan[1]†

5 [1]Agricultural Genomic Institute, Chinese Academy of Agricultural Sciences,

6 Shenzhen, Guangdong, 518120, China.

7 [2]BGI-Shenzhen, Shenzhen, Guangdong, 518083, China

8 Conghui Liu: rapherlch@163.com; Yan Zhang: milrazhang@163.com; Yuwei Ren:

9 xiaoshudaxia@126.com; Hengchao Wang: wanghengchao000@qq.com; Shuqu Li:

10 lishuqu1234@163.com; Fan Jiang: greatjf@163.com; Lijuan Yin:

11 yinlijuan1005@163.com; Xi Qiao: qiaoxi@caas.cn; Guojie Zhang:

12 guojie.zhang@bio.ku.dk

13 *These authors contributed equally to this work.

14 †Correspondence should be addressed to Bo Liu (lb_bobo@aliyun.com) or Wei Fan

15 (fanwei@caas.cn).

16 **Abstract**

17 **Background:** The golden apple snail (*Pomacea canaliculata*) is a fresh water snail

18 listed among the top-100 worst invasive species worldwide and a noted agricultural

19 and quarantine pest that causes great economic losses. It is characterized by fast

20 growth, strong stress tolerance, a high reproduction rate, and adaptation to a broad

21 range of environments.

22 **Results:** Here, we used long-read sequencing to produce a 440-Mb high-quality,

23 chromosome-level assembly of the *P. canaliculata* genome. In total, 50 Mb (11.4%)

24 repeat sequences and 21,533 gene models were identified in the genome. The major

25 findings of this study include the recent explosion of DNA/hAT-Charlie transposable

26 elements (TEs), the expansion of the P450 gene family and the constitution of the

27 cellular homeostasis system, which contributes to ecological plasticity in stress

28 adaptation. In addition, the high transcriptional levels of perivitellin genes in the

29 ovary and albumen gland promote the function of nutrient supply and defence ability

30 in eggs. Furthermore, the gut metagenome also contains diverse genes for food

31 digestion and xenobiotic degradation.

32 **Conclusions:** These findings collectively provide novel insights into the molecular

33 mechanisms of the ecological plasticity and high invasiveness.

34 **Keywords:** golden apple snail, *Pomacea canaliculata*, genome, adaptive evolution,

35 stress tolerance, P450, reproduction, perivitelline, metagenome

36 **Background**

37 The golden apple snail *Pomacea canaliculata* (family Ampullariidae, order

38 Architaenioglossa) is a fresh water snail listed among the world's top 100 worst

39 invasive species [1] and is considered an agricultural and quarantine pest worldwide

40 [2]. Native to tropical and subtropical South America, *P. canaliculata* gradually

41 spread to non-indigenous regions, such as Southeast and East Asia [3], Africa [4],

42 North America [5], Oceania [6] and even Europe [7]. Its successful

43 biological invasion was closely related to its polyphagous feeding habits [8],

44 voracious appetite [9], broad environmental adaptability [10] and rapid growth and

45 high rate of reproduction [11]. In addition to its ecological impact, *P. canaliculata*

46 ravages a wide range of crops, including grains, fruits and vegetables [12], causing

47 severe economic losses each year as a result of yield loss, replanting cost and

48 expenditures on control (https://www.cabi.org/isc/datasheet/68490). More seriously, *P.*

49 *canaliculata* has been involved in the transmission of a fatal human disease,

50 eosinophilic meningitis, which first appeared in East Asia where people frequently

51 consume these snails [13]. During this pathophoresis, *P. canaliculata* acts as an

52 important intermediate host of the pathogenic parasite *Angiostrongylus cantonensis*,

53 and the range of infected regions is still expanding, creating a great challenge in terms

54 of human health [14, 15].

55 Molluscs are a highly diverse group, second only to arthropods in species number [16],

56 and their high biodiversity makes them an excellent model to address issues such as

57 biogeography, adaptability and evolutionary processes [17]. The worldwide invasive

58 species *P. canaliculata* provides valuable potential in these fields [18]. As a primitive

59 circumtropical species, *P. canaliculata* possesses strong ecological plasticity with

60 many advantages, including low-temperature resistance [19] and drought tolerance

61 [20], which has contributed to its competitive success in resource acquisition. *P.*

62 *canaliculata* has been reported to establish populations at temperatures ranging from

63 10 °C to 35 °C [19, 21]. Additionally, *P. canaliculata* tolerates heavy metal

64 contamination. When living in contaminated water, the gill is enriched with a high

65 concentration of heavy metals, and histopathological changes in the digestive tract are

66    detected; however, an extremely low mortality rate is observed [22]. The conspicuous

67    colouration and neurotoxic lectin could confer a survival advantage on the eggs,

68    defending the embryos against potential predators [23]. Moreover, an

69    immune-neuroendocrine system can also be detected in *P. canaliculata*, as

70    demonstrated by the existence of a specific immune memory after bacterial challenge

71    [24, 25], broadening the study of invertebrate immunology.

72    The rich phenotypic and genetic diversity of molluscs makes them an excellent

73    species group for addressing many important issues in evolution, ecology and

74    function. However, the genomic resources on Mollusca are still insufficient compared

75    with those of other close phyla, such as Arthropoda and Nematoda, and few molluscs

76    can be employed as model organisms. *P. canaliculata*, however, possesses the

77    potential to be a model organism among molluscs because of several inherent

78    characteristics. For example, *P. canaliculata* is easy to acquire because it has a broad

79    global distribution originating from a primarily circumtropical environment.

80    Moreover, its high adaptability, rapid growth and efficient reproduction facilitate the

81    cultivation of *P. canaliculata* in the laboratory.

82    In recent years, the genomic features of *P. canaliculata* have been increasingly studied.

83    After the discovery of 14 pachytene bivalents in the karyotype [26], molecular

84    markers were identified to investigate the genetic diversity of the *P. canaliculata*

85    population, including 369 amplified fragment length polymorphism (AFLP) loci [27],

86    16,717 simple sequence repeats (SSR) [28, 29] and 15,412 single-nucleotide

87    polymorphisms SNPs [30]. In addition, multiple transcriptome analyses have been

4

88 performed to investigate the adaptation, invasion and immune mechanisms of *P.*

89 *canaliculata*. For instance, Sun et al. reported 128,436 unigenes based on a de novo

90 assembly of Illumina reads [30]; transcriptome changes in response to heat stress and

91 starving incubation were used to characterize its invasive and adaptive abilities [31,

92 32]; a transcriptome analysis comparing invasive *P. canaliculata* and indigenous

93 *Cipangopaludina cathayensis* provided insights into biological invasion [29]; and 402

94 immune-related differentially expressed genes (DEGs) in response to

95 lipopolysaccharide (LPyS) challenge were used to explore the mechanisms of defence

96 against pathogens [33]. Furthermore, proteomics tools such as isobaric tags for

97 relative and absolute quantitation (iTRAQ), and liquid chromatography-tandem mass

98 spectrometry (LC-MS/MS) were also applied in the study of protein expression

99 during estivation and oviposition [34,35], together providing plentiful omics- data for

100 the functional analysis of *P. canaliculata*.

101 However, research at the whole-genome level in *P. canaliculata* still lags far behind

102 that in other mollusc species due to the lack of a high-quality reference genome.

103 Multiple draft genomes of molluscs have been published, including the genomes of

104 the California sea hare [36], Pacific oyster [37], pearl oyster [38,39], owl limpet [40],

105 California two-spot octopus [41], golden mussel [42], and *Biomphalaria* snails [43],

106 greatly promoting research on mollusc genomics. In this study, we present a

107 chromosome-level genome assembly of *P. canaliculata* with high-quality gene

108 annotation, transcriptome data from several tissues and under various conditions, and

109 metagenomic data from the intestinal tracts, all of which were then applied to study

5

110 the species-specific environmental adaptation characteristics, such as the cellular

111 homeostasis system underlying strong stress and the colour and nutrient contents of

112 the eggs. Our data will not only strengthen the understanding of the evolutionary

113 mechanisms of molluscs and the molecular basis of biological invasion but also foster

114 the development of approaches to control the invasion of *P. canaliculata* and provide

115 a basis for interrupting the transmission of pathogenetic nematode parasites.


116 **RESULTS**


117 **Complete genome assembly at the chromosome level**


118 We generated 26.6 Gb (60.1 X) of PacBio SMRT raw reads with an average read

119 length of 10.1 kb, and 291 Gb (652.4 X) of Illumina HiSeq paired-end reads with an

120 average read length of 150-250 bp using DNA extracted from a single adult *P.*

121 *canaliculata* (Table S1). The 24.4 Gb (55.4 X) of clean PacBio SMRT reads that

122 passed quality filtering were assembled by smartdenovo

123 (https://github.com/ruanjue/smartdenovo), resulting in an assembly of 1,234 raw

124 contigs with a total length of 473.0 Mb and an N50 length of 1.0 Mb. After filtering of

125 alternatively heterozygous contigs, the 745 resulting contigs with a total length of

126 440.1 Mb and an N50 length of 1.1 Mb were taken as the final contigs. Previous

127 karyotype research has shown that the haploid *P. canaliculata* genome consists of 14

128 chromosomes [26]. Based on the Hi-C data, 439.5 Mb (99.9%) of final contigs were

129 anchored and oriented into 14 large scaffolds, each corresponding to a natural

130 chromosome (Figure 1a and Figure 1b), with the longest 45.4 Mb and the shortest

131    27.2 Mb. This assembly quality is much better than that of the other molluscan

132    genomes published thus far (Table 1). In addition to the length and continuity of the

133    assembled sequences, another important aspect for evaluating genome assembly is the

134    ratio of genome coverage. With an estimated genome size of 446 Mb and genome

135    heterozygosity between 1% and 2% based on the distribution of k-mer frequency [44]

136    (Figure S1), ~98.6 % of the *P. canaliculata* genome has been assembled. To further

137    confirm the accuracy and completeness of the assembly, we mapped the Illumina

138    shotgun reads to the assembled reference genome. Significantly, 97% and 95% of the

139    genome-derived and transcriptome-derived reads, respectively, could be aligned to the

140    reference genome, suggesting no obvious bias in sequencing and assembly.

141    Additionally, the mitochondrial genome of *P. canaliculata* was assembled as a single

142    contig 15,707 bp in length, which has 99.9% sequence identity to the published

143    mitochondrial genome (GenBank: KJ739609.1) (Figure S2). This high-quality

144    reference genome provides a good foundation for gene annotation.

145    The protein-coding genes were predicted on the reference genome by EVM,

146    integrating evidence from *de novo* prediction, transcriptome and homology data. In

147    total, 21,533 gene models were predicted as the reference gene set, with coding

148    regions spanning ~32.2 Mb (7.3 %) of the genome (Table 1 and Table S2). The

149    distribution of CDS length in *P. canaliculata* is similar to that in closely related

150    species (Figure 1c). Overall, 97.5% of the reference genes were supported by

151    transcriptome data, and 98.0% of eukaryote core genes from OrthoDB

152    (http://www.orthodb.org/) were identified in the reference gene set by BUSCO. These

153  results were comparable to those in other published molluscan genomes (Table 1). In

154  functional annotation, a total of 19,815 (91.9 %) reference genes were annotated by at

155  least one functional database. Specifically, 15,662 (72.7%), 13,769 (63.4%), 17,081

156  (79.3%), 18,847 (87.5%) and 17,003 (79.9%) reference genes were annotated with the

157  eggNOG, KEGG, NR, InterPro and UniProt databases, respectively (Figure S3).

**Signs of adaptive evolution in *P. canaliculata* genome**

159  To gain insight into the evolutionary perspective of *P. canaliculata*, a phylogenetic

160  tree was built based on 306 high-confidence single-copy orthologous genes from nine

161  related species (*P. canaliculata, Lottia gigantea, Aplysia californica, Biomphalaria*

162  *glabrata, Crassostrea gigas, Octopus bimaculoides, Pinctada fucata, Lingula anatina*

163  and *Limnoperna fortunei*) by PhyML [45] and the divergence time was estimated

164  using MCMCTree [46]. The results show that *P. canaliculata* diverged from the

165  ancestor of *B. glabrata* and *A. californica* 372 million years ago (Mya) and from *L.*

166  *gigantea* 491 Mya (Figure 2a).

167  Then, the molluscan orthologous genes were investigated for adaptive evolution.

168  Utilizing pairwise protein sequence similarities, gene family clustering was conducted

169  by orthoFinder [47]. A total of 239,541 reference genes from the nine species were

170  clustered into 69,582 orthologous groups, among which 14,766 orthologous groups

171  contained at least two genes each. We identified 66 orthologous groups that

172  underwent common expansion in both *P. canaliculata* and *L. fortunei* but not the other

173  seven species. The functions of these orthologous groups are mainly related to signal

8

174 transduction; replication and repair; translation, glycan biosynthesis and metabolism;

175 lipid metabolism; and the endocrine, immune and nervous systems (Figure S4). These

176 relations suggest that the gene families that underwent expansion may play important

177 roles in adaptation to the environment as invasive species.

178 The high-coverage genome assembly enables a comprehensive analysis of the

179 transposable elements (TEs), which play multiple roles in driving genome evolution

180 in eukaryotes [48]. In total, we identified 49.6 Mb TE sequences in the assembled *P.*

181 *canaliculata* genome (Table 1), including 3.4 Mb long terminal repeats (LTRs), 27.2

182 Mb long interspersed elements (LINEs), 17.5 Mb DNA transposons and 1.5 Mb short

183 interspersed elements (SINEs). Next, we analysed the divergence rate of each class of

184 TEs among the available sequenced mollusc genomes. Notably, the TE class of DNA

185 transposons showed a specific peak at a divergence rate of ~4% divergence rate for *P.*

186 *canaliculata* and *C. gigas* (Figure 2b), indicating a recent explosion of DNA

187 transposons in these two species. We analysed the expression of 709 genes, including

188 DNA elements restricted to the 4% peak inside the gene region, compared with that of

189 the other genes outside the 4% peak (Figure S5). DEGs were defined here by P-values

190 smaller than 0.05 for comparison of the treatment (heat, cold, heavy metal and air

191 exposure) and control data. The percentages of DEGs in the 4% peak were higher than

192 those of genes outside the peak (10.2% higher for heat, 8.6% higher for cold, 8.6%

193 higher for heavy metal, and 7.3% higher for air exposure). Among the DEGs in the 4%

194 peak, approximately half were up-regulated, and the other half were down-regulated.

195 Moreover, the DEGs in the 4% peak were mainly enriched in cellular metabolic

196    process, response to stimulus, localization and signaling according to GO annotation.

197    These results indicated that genes in the 4% peak were likely to be more active in the

198    response to stimulus, promoting potential plasticity in stress adaptation. TEs are

199    powerful facilitators of evolution that generate "evolutionary potential" to introduce

200    small adaptive changes within a lineage, and the importance of TEs in stress

201    responses and adaptation has been reported in numerous studies [49,50]. The recent

202    explosion of DNA TEs in *P. canaliculata* could also play an important role in

203    promoting the potential plasticity in stress adaptation.

204    **Investigation of cellular homeostasis system underlying strong stress adaptation**

205    The homeostasis system plays a crucial role in stress adaptability, providing the

206    molecular basis for re-establishing dynamic equilibrium after challenges by various

207    environmental stressors, including temperature, air exposure, anthropogenic pollution

208    and pathogens [51]. In this study, we addressed three constituent parts of the cellular

209    homeostasis system, which contributes to the successful ecological plasticity of *P.*

210    *canaliculata* (Figure 3). The transcriptomes of the hemocytes after different stimuli

211    (cold, heat, heavy metal and air exposure) were also sequenced and analysed to

212    address the potential roles of these genes in the cellular homeostasis system.

213    The unfolded protein response (UPR) system is the central component of protein

214    homeostasis [52]. Heat shock proteins (HSPs) act as molecular chaperones to

215    maintain correct folding, and heat shock transcription factor 1 (HSF1) is responsible

216    for the transcriptional induction of HSPs [53]. In the *P. canaliculata* genome, 13

217 HSP70s, 6 HSP90s, 7 HSP40s and 11 HSFs were identified (Table S3), and the

218 expression of HSP90s and HSFs was highly induced in response to heat, cold, heavy

219 metal and air exposure (Table S4 and Figure S6). Inositol-requiring protein 1 (IRE1),

220 protein kinase RNA-like ER kinase (PERK), and activating transcription factor 6

221 (ATF6) are three mediators recruited by the endoplasmic reticulum (ER) to regulate

222 the UPR [54]. We found putative coding genes of the three core mediators, their

223 respective downstream transcription factors, and the corresponding recognition

224 chaperones in the *P. canaliculata* genome (Table S3).

225 The xenobiotic biotransformation system helps the molluscs adapt to toxicants,

226 especially pesticides in aquatic environments [55]. Manual annotation of this genome

227 identified 157 cytochrome P450s (CYP450s), 15 flavin-containing monooxygenases

228 (FMOs), 53 glutathione S-transferases (GSTs) and 105 ATP binding cassette (ABC)

229 transporters, most of which showed up-regulated expression under stress (Table S3

230 and Table S4). These proteins have been shown to function in contaminant detection,

231 conjugative modification and expulsion for xenobiotic detoxification [56-58].

232 The massive production of reactive oxygen species (ROS) and reactive oxygen

233 intermediates (ROIs) induced by stress leads to many pathological conditions, and

234 antioxidant systems protect the organism from superoxide [59]. Four main antioxidant

235 enzyme classes, namely, superoxide dismutase (SOD), catalase (CAT), peroxidase

236 (Prx), and glutathione peroxidase (GPX), were found in *P. canaliculata* and showed

237 elevated global expression in response to stress (Table S3 and Table S4).

238 Apoptosis is a process of cell death when sensing stress and the regulation of

11

239 apoptosis contributes to the dynamic homeostasis of the internal environment. In *P.*

240 *canaliculata*, we propose the existence of both intrinsic and extrinsic apoptotic

241 signaling pathways, evidenced by the presence of homologous genes involved in both

242 pathways. These two pathways could be activated by cytochrome C and tumour

243 necrosis factor receptor (TNFR), respectively (Table S3). Inhibitors of apoptosis, such

244 as XIAP, Bcl2 and Bak, are also detected and show increased expression in response

245 to stress (Table S4), which is expected to delay the process of apoptosis and cell death

246 in the stress response.

247 **The expansion of the P450 gene family contribute to stress tolerance**

248 Cytochrome P450 (CYP) enzymes are a monooxygenase family with highly diverse

249 structures and functions that have been widely identified in all kingdoms of life [60].

250 P450s catalyse the reductive scission of molecular oxygen and are responsible for the

251 synthesis and metabolism of various molecules, including drugs, hormones,

252 antibiotics, pesticides, carcinogens and toxins [61]. The hormones they synthesize,

253 such as glucocorticoids, mineralocorticoids, progestins, and sex hormones, are critical

254 to stress response, growth and reproduction, and the endogenous and exogenous

255 chemical metabolism participate in combatting toxic compounds [62].

256 We found that the *P. canaliculata* CYP gene family had undergone an expansion

257 compared to that in the other molluscs. We identified 157 genes in the genome of *P.*

258 *canaliculata* and 128, 102, 135, 115, 78, 52 and 94 genes in *A. californica, B.*

259 *glabrata, C. gigas, L. fortunei, L. gigantea, O. bimaculoides* and *P. fucata* respectively,

260    using the same standard (Figure 4a). An expansive trend was also observed in

261    comparison with other model species, such as *Homo sapiens* (57), *Mus musculus*

262    (102), *Danio rerio* (94) and *Drosophila melanogaster* (94) [63]. Gene expansion was

263    mainly found in the CYP2U and CYP3A sub-families, whereas fewer genes were

264    expanded in CYP4F. In mammals, CYP2U participates in the metabolism of fatty

265    acids to generate bioactive eicosanoid derivatives, potentially regulating the

266    development of immune function [64]. In *P. canaliculata*, 40 genes formed the

267    CYP2U clade, mainly expressed in the hepatopancreas (Figure 4b and Table S5_a,

268    Table S5_b). CYP3A is a versatile enzyme that metabolizes a wide range of

269    xenobiotics, and its production promotes the growth of various cell types [65]. The 56

270    CYP3A genes are comprehensively expressed in the hepatopancreas, gill and kidney

271    (Figure 4b and Table S5_a, Table S5_b). CYP4F possesses epoxygenase activity,

272    metabolizing fatty acids to epoxides to suppress hypertension, pain perception and

273    inflammation [66]. Twenty genes were identified in CYP4F, and Pc06G011748,

274    Pc06G011460, Pc06G011458, Pc06G011459, Pc04G006708, Pc04G006710 and

275    Pc04G006707 exhibited highly induced expression levels under cold, heat, heavy

276    metal and air exposure stress, indicating their critical roles in the stress tolerance

277    (Figure 4b, Table S5_a and Table S5_b).


278    **The identification of perivitellin genes and their high transcriptional levels in the**

279    **ovary and albumen gland**


280    *P. canaliculata* has eggs characterized by abundant nutrients, reddish or pinkish

281    colour, aerial oviposition and neurotoxicity [23, 67] due to the perivitelline Fluid

13

282 (PVF), which fills the space between the eggshell and the embryo and consists of

283 carbohydrates, lipids and proteins (Figure 5a). The PVF proteins in *P. canaliculata*,

284 include three major components, PcOvo, PcPV2, and PcPV3 [68], collectively named

285 perivitellins, which make up 90% of the total proteins, whereas most of the other

286 dozens of low-abundance components each account for less than 1% of the total

287 proteins [35]. The perivitellins are not only responsible for the major supply of

288 materials and energy during embryogenesis but also provide warning pigments and

289 deadly toxicants against predators [23, 69, 70].

290 We identified 28 candidate PVF genes in *P. canaliculata* by mapping each of the 59

291 fragmental PVF protein sequences derived from a previous proteomics study by Sun

292 [35] to its best hit in the reference gene set of *P. canaliculata*, using BLASTP with

293 requirements of over 85% identity and at least 50% alignment length (Table S6). Then,

294 the functional annotation of those fragmental proteins was also transferred to our

295 identified PVF genes. The transcriptome data show that 22 (79%) of the 28 candidate

296 PVF genes exhibit their highest expression in the ovary and albumen gland (PVF

297 protein synthesis factory) among all 7 tissues (Figure 5b and Table S7), confirming

298 that most of them are genuine functional PVF genes. Six of these 28 candidate PVF

299 genes are perivitellin genes, including two PcOvo genes, Pc09G015543 (PcOvo2) and

300 Pc09G015548 (PcOvo3); two PcPV2 genes, Pc07G012572 (PcPV2-31) and

301 Pc07G012571 (PcPV2-67); and two possible PcPV3 genes, Pc09G015546 and

302 Pc09G015547. The expression levels of these 6 genes in the ovary and albumen gland

303 are much higher than those of the other 22 candidate PVF genes.

14

304    By analysing the orthoFinder gene families that include orthologous and paralogous

305    genes from *P. canaliculata* and 8 other sequenced mollusc species, we found that

306    these 28 candidate PVF genes were classified into 20 multiple-gene families ($\geq 2$

307    genes) and 7 single-gene families (only one gene) (Table S8). Notably, 5 of the 6

308    perivitellin genes were classified into single-gene families, except for Pc07G012571

309    (PcPV2-67), which not only has homologous genes in other mollusc species but also

310    has three paralogous genes in *P. canaliculata* itself. However, none of these three

311    PcPV2-67 paralogous genes in *P. canaliculata* showed higher expression in the ovary

312    and albumen gland than in other tissues, indicating that they are likely not

313    PVF-related genes, i.e., only Pc07G012571 plays a role in PVF. The nearly unique

314    and single-copy nature of the 6 perivitellin genes in *P. canaliculata*, may be explained

315    by the long evolutionary distance, over 200 Mya for *P. canaliculata* and its most

316    closely related species, *A. californica*, as well as numerous differences in their living

317    characteristics and egg structures. Another possible explanation is that these 6 major

318    PVF genes may have experienced rapid evolution in their history to adapt to the

319    changing environment.

**The gut microbiome plays important roles in stress resistance and food digestion**

321    The gut microbiome is regarded as the "second genome" of the host animal, due to the

322    fact that gut microbiota contributes to the food digestion, immune system

323    development, and many other processes important to the host. To investigate the

324    relationship between the gut microbiome and the invasive lifestyle of *P. canaliculata*,

325 we collected gut digesta samples from 70 *P. canaliculata* snails and generated 31 Gb

326 of high-quality metagenomic data on the Illumina HiseqX10 platform. To our

327 knowledge, this study is the first in-depth sequencing of the snail gut microbiome. A

328 total of 1,142,095 non-redundant genes were obtained with an average open reading

329 frame (ORF) length of 604 bp (Table S9). The taxonomic composition analysis

330 showed that, at the phylum level, Proteobacteria was predominant, followed by

331 Verrucomicrobia, Bacteroidetes, Firmicutes, Spirochaetes, Actinobacteria, etc. (Table

332 S10_a). At the genus level, the most abundant genera included *Aeromonas*,

333 *Enterobacter*, *Desulfovibrio*, *Citrobacter*, *Comamonas*, *Klebsiella* and *Pseudomonas*

334 (Table S10_b), most of which were also present in *Achatina fulica* [71,72].

335 Interestingly, some of the most abundant genera, such as *Desulfovibrio*, *Citrobacter*

336 and *Pseudomonas*, were reported as having strong abilities to remove heavy metals by

337 bioprecipitation and bioabsorption [73-75]. For example, the sulfur-reducing bacteria

338 *Desulfovibrio* produce $H_2S$, which precipitates metals and therefore reduces the toxic

339 effects of dissolved metals [73]. Based on the KEGG pathway database, the complete

340 sulfate reduction metabolism pathway was identified in the *P. canaliculata* gut

341 microbiome. We suggested that these gut microbes might help *P. canaliculata* survive

342 the environmental stress of heavy metals in harsh conditions. In addition, a large

343 number of genes in xenobiotic biodegradation and metabolism pathways were

344 annotated, corresponding to 288 KEGG orthologous groups (KOs) and 21 pathways

345 (Table S11). As many of the pathways, such as benzoate degradation, toluene

346 degradation, xylene degradation and steroid degradation, could not be identified in the

16

347 host genome through KO analysis, we suggested that microbial detoxification abilities

348 may contribute to the ability *P. canaliculata* to resist stresses caused by xenobiotics

349 such as pesticides and environmental pollutants.

350 In digestion, the gut microbes are directly involved in the breakdown of the cellulose

351 portion of the diet, and previous studies have isolated cellulolytic bacteria and

352 evaluated the cellulolytic enzyme activities [76]. Our work found a broader range of

353 carbohydrate active enzymes (CAZymes). Of the 208 annotated CAZyme families, 99

354 were glycoside hydrolase (GH) families (Table S12). Enzymes that could be classified

355 as cellulases, endohemicelluloses, debranching enzymes, and

356 oligosaccharide-degrading enzymes were all identified. These findings indicate that

357 the gut microbiome provides assistance in digesting a broad range of food sources,

358 enabling *P. canaliculata* to grow rapidly and adapt to an invasive lifestyle.


359 **Conclusion**


360 Given its environmental invasiveness, broad stress adaptability and rapid reproduction,

361 the golden apple snail *P. canaliculata* has received a vast amount of attention

362 worldwide. However, the underlying genetic mechanisms of these properties have not

363 been comprehensively uncovered. The chromosome-level genome of *P. canaliculata*

364 presented in this study sheds the first light on into the genomic basis of its ecological

365 plasticity in response to various stressors. The major findings of this study include the

366 recent explosion of DNA/hAT-Charlie TEs, the expansion of the P450 gene family

367 and the constitution of the cellular homeostasis system, all of which contribute to the

17

368  plasticity of the organism in stress adaptation. Although the function of the recently

369  originated TEs could not be confirmed, TEs are considered powerful facilitators in

370  adaptive evolution, suggesting that their increased number plays an important role in

371  the stress resistance of *P. canaliculata*. The UPR system, xenobiotic biotransformation

372  system and ROS system are all major components of the cellular homeostasis system,

373  and the P450s in particular underwent expansion with specific functions. In addition,

374  exclusive perivitellin genes were identified in the *P. canaliculata* genome, and they

375  are believed to contribute to the high reproductive rate and the expansion of habitats.

376  Furthermore, the gut metagenome contains diverse genes for food digestion and

377  xenobiotic degradation. These findings collectively provide novel insight into the

378  molecular mechanisms of ecological plasticity and high invasiveness.

379  In this study, we report a fine reference genome of *P. canaliculata*, first

380  chromosome-level Mollusca genome published. With widespread distribution, rapid

381  growth and efficient reproduction, *P. canaliculata* possesses the potential to be a

382  model organism of Mollusca. As its cellular complexity and conservation of pathways

383  also make *P. canaliculata* a useful representative of Mollusca, the genome described

384  in this study can be used to advance our understanding of the molecular mechanisms

385  involved in various scientific questions regarding Mollusca.


386  **Methods**


387  **Samples collection and sequencing**


388  Adults of *P. canaliculata* were collected from a local paddy field in Shenzhen,

18

389 Guangdong province, China, and maintained in aerated freshwater at $15 \pm 2$ °C for a

390 week before processing. Genomic DNA was extracted from the foot muscles of a

391 single *P. canaliculata* for constructing PCR free Illumina 350-bp insert libraries and

392 PacBio 20-kb insert library, and sequenced on Illumina HiSeq 2500 and PacBio

393 SMRT platforms, respectively. The Hi-C library was prepared using the muscle tissue

394 of another single *P. canaliculata* by following methods: Nuclear DNA was

395 cross-linked in situ, extracted, and then digested with a restriction enzyme. The sticky

396 ends of the digested fragments were biotinylated, diluted, and then ligated to each

397 other randomly. Biotinylated DNA fragments were enriched and sheared again for

398 preparing the sequencing library, which was then sequenced on a HiSeq X Ten

399 platform (Illumina).

400 Seven tissues including embryos (2 days post fertilization), gill, hemocytes,

401 hepatopancreas, kidney, ovary and albumen gland and testis from six animals were

402 collected as parallel samples. Next, animals were cultivated in 37 °C and 10 °C for 24

403 hours heat and cold tolerance, in $Cr^{3+}$(2mg $L^{-1}$), $Cu^{2+}$(0.2mg $L^{-1}$) and $Pb^{2+}$(1mg $L^{-1}$)

404 for 24 hours heavy metal tolerance, and in waterless tank for 7 days air exposure.

405 Then the hemocytes were harvested and stored, with three replicates for each group.

406 In final, total RNAs were extracted from the stored tissues of *P. canaliculata*

407 materials, and then mRNAs were pulled out by beads with poly-T for constructing

408 cDNA libraries (insert 350-bp), and sequenced on an Illumina HiSeq 2500 sequencer.

409 The intestinal digesta from 70 adult snails of *P. canaliculata* were collected, pooled

410 into 6 samples and stored at $-20$ °C until microbial DNA was extracted. A

411  combination of cell lysis treatments was applied, including five freeze-thaw cycles

412  (alternating between 65 °C and liquid nitrogen for 5 min), repeated beads-beating in

413  ASL buffer (cat. no. 19082; Qiagen Inc.), and incubated at 95 °C for 15 min. DNA

414  was isolated following the protocol reported protocol [77]. Paired-end libraries of

415  metagenomic DNA were prepared with an insert size of 350 base pairs (bp) following

416  the manufacture's protocol (cat. no. E7645L; New England Biolabs). Sequencing was

417  performed on Illumina HiSeq X10.

418  **Genome assembly and annotation**

419  The Illumina raw reads were filtered by trimming the adapter sequence and

420  low-quality regions (https://github.com/fanagislab/common_use), resulting in clean

421  and high-quality reads with an average error rate < 0.001. For the PacBio raw data,

422  the short subreads (< 2 kb) and low-quality (error rate > 0.2) subreads were filtered

423  out, and only one representative subread was retained for each PacBio read. The clean

424  PacBio reads were assembled by the software smartdenovo

425  (https://github.com/ruanjue/smartdenovo), after which Illumina reads were aligned to

426  the contigs by BWA-MEM (BWA, RRID:SCR_010910), and single base errors in the

427  contigs were corrected by Pilon v1.16 (Pilon, RRID:SCR_014731) with the

428  parameters "-fix bases, -nonpf, -minqual 20". The *P. canaliculata* genome is highly

429  heterozygous, as illustrated by the double peaks on the distribution curve of k-mer

430  frequency, and the current assembly algorithm tends to collapse homozygous regions

431  and report heterozygous regions in alternative contigs. To obtain a haploid reference

432    contigs, we employed a whole-genome alignment (WGA) strategy with MUMmer

433    v3.23 to recognize and selectively remove alternative heterozygous contigs, which

434    were characterized by shorter length (less than 200 kb) and the ability of most regions

435    (more than 50%) to be aligned to another larger contig with confident identity (higher

436    than 80%). Next, Hi-C sequencing data were aligned to the haploid reference contigs

437    by BWA-MEM, and then these contigs were clustered into chromosomes with

438    LACH-ESIS (http://shendurelab.github.io/LACHESIS/).

439    A de novo repeat library for *P. canaliculata* was constructed by RepeatModeler v.

440    1.0.4                    (RepeatModeler,                    RRID:SCR_015027;

441    http://www.repeatmasker.org/RepeatModeler.html). TEs in the *P. canaliculata*

442    genome    were    also    identified    by    RepeatMasker    v4.0.6    (RepeatMasker,

443    RRID:SCR_012954; http://www.repeatmasker.org/) using both the Repbase library

444    and the de novo library. Tandem repeats in the *P. canaliculata* genome were predicted

445    using Tandem Repeats Finder v4.07b [78]. The divergence rates of TEs were

446    calculated between the identified TE elements in the genome and their consensus

447    sequence at the TE family level.

448    The gene models in the *P. canaliculata* genome were predicted by EVidence Modeler

449    v1.1.1 [79], integrating evidence from ab initio predictions, homology-based searches

450    and RNA-seq alignments. Then, these gene models were annotated by RNA-seq data,

451    UniProt database and InterProScan software (InterProScan, RRID:SCR_005829) [80].

452    Finally, the gene models were retained if they had at least one piece of supporting

453    evidence from the UniProt database, InterProScan domain and RNA-seq data. Gene

454  functional annotation was performed by aligning the protein sequences to the NCBI

455  NR, UniProt, COG and KEGG databases with BLASTP v2.3.0+ under an E-value

456  cutoff of $10^{-5}$ and choosing the best hit. Pathway analysis and functional classification

457  were conducted based on the KEGG database [81]. InterProScan was used to assign

458  preliminary GO terms, Pfam domains and IPR domains to the gene models.

459  **Evolutionary analysis**

460  Orthologous and paralogous groups were assigned from seven species (*P.*

461  *canaliculata, Lottia gigantea, Aplysia californica, Biomphalaria glabrata,*

462  *Crassostrea gigas, Octopus bimaculoides, Pinctada fucata, Limnoperna fortunei* and

463  *Lingula anatina*) by OrthoFinder [47] with default parameters. Orthologous groups

464  that contained only one gene for each species were selected to construct the

465  phylogenetic tree. The protein sequences of each gene family were independently

466  aligned by muscle v3.8.31 [82] and then concatenated into one super-sequence. The

467  phylogenetic tree was constructed by maximum likelihood (ML) using PhyML v3.0

468  (PhyML,          RRID:SCR_014629)          [45]          with          the

469  best-fit model (LG+I+G) estimated by ProtTest3   [83].    The    Bayesian    relaxed

470  molecular clock (BRMC) approach was adopted to estimate the neutral evolutionary

471  rate and species divergence time using the program MCMCTree, implemented in the

472  PAML v4.9 package (PAML, RRID:SCR_014932) [46]. The tree was calibrated with

473  the following time frames to constrain the age of the nodes between the species:

474  minimum = 260 Ma and maximum = 290 Ma for *P. fucata* and *C. gigas* [84];

475    minimum = 450 Ma and maximum = 480 Ma for *A. californica* (or *B. glabrata*) and *L.*

476    *gigantea* [85]. The calibration time (fossil record time) interval (550-610 Mya) of *O.*

477    *bimaculoides* was adopted from previous results [86]. To identify the common

478    expanded gene families, we compared the *P. canaliculata* and *L. fortunei* with other

479    seven species. The gene number of orthologous group in *P. canaliculata* and *L.*

480    *fortunei* were two or more times than that in all of other species, respectively.

481    Additionally, these gene families with P-value less than 0.01 were considered as

482    expansion by z-test.

483    **Transcriptome data analysis**

484    Transcriptome reads were trimmed with the same method for genomic reads

485    (https://github.com/fanagislab/common_use), and then mapped to the reference

486    genome of *P. canaliculata* using TopHat v. 2.1.0 (TopHat, RRID:SCR_013035) with

487    default settings. The expression level of each reference gene in terms of FPKM was

488    computed by cufflinks v2.2.1 (cufflinks, RRID:SCR_014597). A gene was considered

489    to be expressed if its FPKM > 0. Differential gene expression analysis was conducted

490    using cuffdiff v2.2.1.

491    **Metagenome data analysis**

492    The Illumina raw reads were filtered by trimming the adapter sequence and

493    low-quality regions (https://github.com/fanagislab/common_use), resulting in

494    high-quality reads with an average error rate < 0.001. Then, the reads mapped to the

495    following genomes by BWA-MEM were filtered out

23

496    (https://github.com/fanagislab/metagenome_analysis.git)[87],    to    exclude    the

497    contaminated host, food, parasite, and human DNA sequences. The genomes include:

498    the *P. canaliculata* genome, the *Brassica rapa* genome, the *Oryza sativa* genome, 2

499    *Angiostrongylus cantonensis* genomes, the *Caenorhabditis elegans* genome, the

500    *Schistosoma mansoni* genome, the *Clonorchis sinensis* genome, the *Fasciola hepatica*

501    genome, the *Danio rerio* genome, and the *human hg38* genome. Finally, short reads

502    (length < 75 bp) and unpaired reads were excluded to form a set of clean reads.

503    The clean reads were assembled by metaSPAdes (v3.11.1) [88] in paired-end mode

504    for each sample. Then, gene prediction was performed on contigs longer than 500 bp

505    by Prodigal v2.6.3 (Prodigal, RRID:SCR_011936) [89] with the parameter "-p meta",

506    and gene models with cds length less than 102 bp were filtered out. A non-redundant

507    (NR) gene set (539,344 genes) was constructed using the gene models predicted from

508    each sample by cd-hit-est (v4.6.6) [90] with the parameter "-c 0.95 -n 10 -G 0 –a S

509    0.9", which adopts a greedy incremental clustering algorithm and the criteria of

510    identity > 95% and overlap > 90% of the shorter genes. Then, the clean reads were

511    mapped onto this NR gene set by BWA-MEM with the criteria of alignment length

512    $\geqslant$ 50 bp and identity > 95%. The unmapped reads from all samples were assembled

513    together, and the genes were predicted again. The newly predicted genes were

514    combined with the previous gene set by cd-hit-est to obtain a new NR gene set

515    (1,147,339 genes). After the taxonomic assignments to the new NR gene set, 5244

516    genes classified as Eukaryota but not fungi were removed, and the final NR gene set

517    (1,142,095 genes) was obtained.

24

518 The taxonomic assignments of the final NR genes were made on the basis of

519 DIAMOND (DIAMOND, RRID:SCR_016071) [91] protein alignment against the

520 NCBI -NR database by CARMA3 [92]. Functional annotation was performed by

521 aligning all the protein sequences to the KEGG [93] database (release 79) using

522 DIAMOND and taking the best hit with the criteria of E-value < 1e-5. CAZymes were

523 annotated with dbCAN (release 5.0) [94] using Hmmer v3.0 hmmscan (Hmmer,

524 RRID:SCR_005305) [95] by taking the best hit with an E-value < 1e-18 and

525 coverage > 0.35.

526 The clean reads from each sample were aligned against the gene catalogue (1,142,095

527 genes) by BWA-MEM with the criteria of alignment length $\geqslant$ 50 bp and identity >

528 95%. Sequence-based gene abundance profiling was performed as previously

529 described [96]. The taxonomic profiles of the samples were calculated by summing

530 the gene abundance according to the taxonomic assignment result.


531 **Abbreviations**


532 *A. californica, Aplysia californica; B. glabrata, Biomphalaria glabrata; C. gigas,*

533 *Crassostrea gigas; O. bimaculoides, Octopus bimaculoides; L. anatina, Lingula*

534 *anatina; L. fortunei, Limnoperna fortunei; L. gigantea, Lottia gigantea; P.*

535 *canaliculata, Pomacea canaliculata; P. fucata, Pinctada fucata;* Hem, hemocytes; Te,

536 testis; Ov, ovary and albumen gland; Kn, kidney; GI, gill; Hp, hepatopancreas, Em,

537 embryo; SSR, simple sequence repeats; mya, million years ago*; BLAST, basic local*

538 *alignment search tool;* SNP, single nucleotide polymorphism; PVF, Pervitelline Fluid;

539 Ovo, ovorubin; AFLP, amplified fragment length polymorphism; DEGs, differentially

540 expressed genes; LPyS, Lipopolysaccharide; iTRAQ, Isobaric Tags For Relative,

541 Absolute Quantitation; LC-MS/MS, Liquid Chromatography-tandem Mass

542 Spectrometry; TEs, transposable elements; LTR, long terminal repeats; LINE, long

543 interspersed elements; SINE, short interspersed elements; UPR, Unfolded protein

544 response; HSPs, heat shock proteins; HSF1, heat shock transcription factor 1; PERK,

545 protein kinase RNA-like ER kinase; ATF6,activating transcription factor 6; ER,

546 endoplasmic reticulum; CYP450s, cytochrome P450s; FMOs, flavin-containing

547 monooxygenases; GSTs, glutathione S-transferases; ABC, ATP binding cassette; ROS,

548 reactive oxygen species; ROI, reactive oxygen intermediates; SOD, superoxide

549 dismutase; CAT, catalase; Prx, peroxidase; GPX, glutathione peroxidase; TNFR,

550 tumor necrosis factor receptor; NR, non-redundant genes; ORF, open reading frame;

551 Kos, orthologous groups; CAZymes, carbohydrate active enzymes; GH, Glycoside

552 Hydrolase.

553 **Availability of data and materials**

554 Tables S1 to S12 and Figures S1 to S6 are available in the supplementary information

555 file. The raw sequencing data has been deposited in DDBJ/EMBL/GenBank under

556 project accession PRJNA427478, SRR6425828 for genomic Illumina_PE125

557 sequencing data, SRR6425829 for genomic Illumina_PE150 sequencing data,

558 SRR6425827 for genomic PacBio sequencing data, SRR6429132~SRR6429164 for

559 transcriptome sequencing data, and SRR6472920~SRR6472925 for gut microbiome

560 data. Other supporting data, including genome assemblies, annotations, phylogenetic

561 tree files and BUSCO results, are available via the *GigaScience* repository GigaDB

562 [97].

### Authors' contributions

564 WF, WQ and CL conceived the study and designed the experiments. CL performed

565 the genome sequencing and assembly, BL performed annotation and evolutionary

566 analysis. CL performed the stress tolerance analysis, YR performed the reproduction

567 analysis, YZ performed the metagenome analysis. HW, SL, FJ, LY, XQ provided

568 suggestions and helped checking. CL, WF, BL, YR, YZ wrote the manuscript, and GZ

569 helped revising the manuscript. All authors read and approved the final manuscript.

### Competing interests

571 The authors declare that they have no competing interests.

### Acknowledgements

27

578 Elite Youth Program of Chinese Academy of Agricultural Sciences. We thank

579 Fanghao Wan, Jue Ruan, Yutao Xiao for providing constructive suggestions to this

580 project.

## 581 Legends of tables and figures

### 582 Tables

583 **Table 1. Summary of assembly and annotation of mollusc genomes**

| Genome feature | *P. canaliculata* | *L. gigantea* | *A. californica* | *B. glabrata* | *C. gigas* | *O. bimaculoides* |
|---|---|---|---|---|---|---|
| Assembled sequences (bp) | 440,071,717 | 359,505,668 | 927,310,431 | 916,377,450 | 557,735,934 | 23,381,887,882 |
| Contig N50 size (bp) | 1,072,857 | 94,165 | 9,817 | 18,978 | 37,218 | 5,982 |
| Contig N90 size (bp) | 303,904 | 10,180 | 1,626 | 5,132 | 11,109 | 1,606 |
| Scaffold N50 size (bp) | 31,531,291 | 1,870,055 | 917,541 | 48,059 | 401,685 | 475,182 |
| Scaffold N90 size (bp) | 23,662,357 | 74,480 | 207,390 | 817 | 68,181 | 79,088 |
| GC content (%) | 40.3 | 33.3 | 40.3 | 36.0 | 33.4 | 36 |
| No. of gene models | 21,533 | 23,824 | 19,909 | 14,224 | 28,402 | 15,814 |
| Avg. CDS length (bp) | 1,497 | 1,136 | 1,568 | 1,066 | 1,472 | 1,535 |
| BUSCO (%) | 98.9 | 98.4 | 98.7 | 72.8 | 99.4 | 98.7 |
| Transposable elements (bp) | 49,579,006 | 37,369,817 | 202,174,499 | 189,550,886 | 103,381,274 | 737,398,096 |
| Tandem repeat (bp) | 873,801 | 257,674 | 8,263,822 | 2,145,821 | 590,907 | 62,633,792 |

### 584 Figures

585 **Figure 1. The genome characteristics of *P. canaliculata*.** (a) Circos plot showing the

586 genomic features. Track 1: 14 linkage groups of the genome; Track 2: distribution of

587 transposon elements in chromosomes; Track 3: protein-coding genes located on

588 chromosomes; Track 4: distribution of GC contents. (**b**) A genome-wide contact

589 matrix from Hi-C data between each pair of the 14 chromosomes using a 100 kb

590 window size. The colour value indicates the base 2 logarithm of the number of valid

591 reads ($\log_2$(valid reads)). (**c**) Distribution of CDS length in six closely related species.

592 **Figure 2. Evolutionary genomic analysis of *P. canaliculata*. (a)** Phylogenetic

593 placement of *P. canaliculata* within the dated tree of molluscs. The estimated

594 divergence time is shown at each branching point, and *P. canaliculata* is shown in red.

595 **(b)** Distribution of divergence rate for the class of DNA transposons in molluscs

596 genomes. The divergence rate was calculated by comparing all TE sequences

597 identified in the genome to the corresponding consensus sequence in each TE

598 subfamily. The red arrow indicates that *P. canaliculata* and *C. gigas* had a recent

599 explosion of TEs at a divergence rate of ~4%.

600 **Figure 3. The cellular homeostasis system in *P. canaliculata*.** The unfolded protein

601 response (UPR) system includes HSPs and HSF in the heat shock response and CNX,

602 NEF, GRP94, BIP, HSP40, ATF6, IRE1, PERK, COP2, XBP, ATF4, TRAM and

603 Derlin in the endoplasmic reticulum unfolded protein response (UPR-ERAD).

604 Apoptotic pathways include XIAPs, Bcl2, caspases, TNFR, and FADD. The

605 antioxidant systems include PRX, SOD, CAT and GPX. The xenobiotic

606 biotransformation system includes EPHX3, P450, FMO and ABC transporter. The

607 colours of the boxes for gene families represent the degree of upregulation

608 (FPKM-stimulus/FPKM-control) as an overall result of stress, including heat, cold,

609 heavy metal and air exposure. Pathways and genes were obtained based on KEGG

610 annotation.

611 **Figure 4. The expansion of the P450 gene family in *P. canaliculata*.** (a)

612 Phylogenetic tree demonstrating orthologous and paralogous relationships of all P450

613 genes from eight species including *P. canaliculata*, *A. californica*, *B. glabrata*, *C.*

29

614 *gigas*, *L. fortunei*, *L. gigantea*, *O. bimaculoides* and *P. fucata*. P450 genes from eight

615 species were obtained based on Pfam annotation (Interpro) with an E-value of $10^{-5}$.

616 Clades are labeled by P450 subfamily names. The tree was constructed using the

617 maximum likelihood method in MEGA7, and the branch length scale indicates the

618 average number of residue substitutions per site. (b) Phylogenetic tree of P450 genes

619 in *P. canaliculata*, which is a subset of the phylogenetic tree for the 7 species, and

620 their heat map of expression (FPKM) in seven tissues (Hem, hemocytes; Te, testis; Ov,

621 ovary and albumen gland; Kn, kidney; Gl, gill; Hp, hepatopancreas; Em, embryo) and

622 heat map of induced expression (FPKM-stimulus/FPKM-control) under stress (Con:

623 control; heat; cold; Hm: heavy metal; Exp: air exposure).

624 **Figure 5. The composition and expression of the *P. canaliculata* perivitellins in**

625 **different tissues.** (a) Perivitelline fluid (PVF) lies under the eggshell and surrounds

626 the embryo. It contains carbohydrates, lipids, and proteins. The proteins are also

627 known as perivitellins and are classified into three categories, PcOvo, PcPV2, and

628 PcPV3. (b) The displayed expression value of PVF proteins is the base 10 logarithm

629 of FPKM ($\log_{10}$FPKM). The genes marked in red encode perivitellins. The 7 tissues

630 examined are abbreviated as follows: Hem, hemocytes; Te, testis; Ov, ovary and

631 albumen gland; Kn, kidney; Gl, gill; Hp, hepatopancreas; Em, embryo.

632 **References**

633 1.      Lowe S, Browne M, Boudjelas S, de Poorter M. 100 of the World's Worst Invasive Alien

634         Species: A selection fromthe Global Invasive Species Database. Auckland, New Zeland:

635    World Conservation Union (IUCN); 2000.

636    2.    Ranamukhaarachchi SL, Wickramasinghe S. Golden apple snails in the world:

637          introduction, impact, and control measures. Global advances in ecology and management

638          of golden apple snails. 2006:133-52.

639    3.    Naylor R. Invasions in Agriculture: Assessing the Cost of the Golden Apple Snail in Asia.

640          Royal Swedish Academy of Sciences. 1996;25:443-8.

641    4.    Berthold T. Vergleichende Anatomie, Phylogenie und historische Biogeographie der

642          Ampullariidae: (Mollusca, Gastropoda). 1991.

643    5.    Howells RG, Burlakova LE, Karatayev AY, Marfurt RK, Burks RL. Native and

644          introduced Ampullaridae in North America: History, status, and ecology. 2006:73-112.

645    6.    Halwart M, Bartley DM. International mechanisms for the control and responsible use of

646          alien species in aquatic ecosystems, with special reference to the golden apple snail. Los

647          Baños, Philippines: Philippine Rice Research Institute (PhilRice); 2006.

648    7.    López MA, Altaba CR, Andree KB, López V. First invasion of the Apple snail Pomacea

649          insularum in Europe. Tentacle. 2010;18:26-8.

650    8.    Estebenet AL, Martín PR. Pomacea canaliculata (Gastropoda: Ampullariidae): life-history

651          traits and their plasticity. Biocell 2002;26:83-9.

652    9.    Lach L. The spread of the introduced freshwater apple snail *Pomacea canaliculata*

653          (Lamarck) (Gastropoda Ampullariidae) on Oahu, Hawaii. Bishop Museum Occasional

654          Papers. 1999;58:66-71.

655    10.   Yusa Y，Sugiura N, Wada T. Predatory Potential of Freshwater Animals on an Invasive

656          Agricultural Pest, the Apple Snail *Pomacea canaliculata* (Gastropoda: Ampullariidae), in

657        Southern Japan. Biol Invasions. 2006;8:137-47.

658    11.    Lach L, Britton DK, Rundell RJ, Cowie RH. Food Preference and Reproductive Plasticity

659        in an Invasive Freshwater Snail. Biol Invasions. 2000;2:279-88.

660    12.    Mochida O. Spread of freshwater *Pomacea* snails (Pilidae, Mollusca) from Argentina to

661        Asia. Micronesica. 1991;3 51-62.

662    13.    Shan L, Zhang Y, Steinmann P, Zhou X. Emerging Angiostrongyliasis in Mainland China.

663        Emerg Infect Dis. 2008;14:161-4.

664    14.    Caldeira RL, Mendonca CL, Goveia CO, Lenzi HL, Graeff-TeixeiraC Lima WS, et al.

665        First record of molluscs naturally infected with *Angiostrongylus cantonensis* (Chen, 1935)

666        (Nematoda: Metastrongylidae) in Brazil. Memórias do Instituto Oswaldo Cruz.

667        2007;102:887-9.

668    15.    McMichael AJ, Beaglehole R. The changing global context of public health. Lancet

669        (London, England). 2000;356:495-9.

670    16.    Chapman A. Numbers of Living Species in Australia and the World. Australian Biological

671        Resources Study; 2009.

672    17.    Lindberg DR, Ponder WF, Haszprunar G. The Mollusca: relationships and patterns from

673        their first half-billion years. Oxford University Press, Oxford; 2004.

674    18.    Hayes KA, Cowie RH, Thiengo SC. A global phylogeny of apple snails: Gondwanan

675        origin, generic relationships, and the influence of outgroup choice (Caenogastropoda:

676        Ampullariidae). Biol J Linn Soc Lond. 2009;98:61-76.

677    19.    Matsukura K, Tsumuki H,Izumi Y, Wada T. Physiological response to low temperature in

678        the freshwater apple snail, *Pomacea canaliculata* (Gastropoda: Ampullariidae). J Exp

679      Biol. 2009;212:2558-63.

680   20.   Yusa Y, Wada T, Takahashi S. Effects of dormant duration, body size, self-burial and

681      water condition on the long-term survival of the apple snail, *Pomacea canaliculata*

682      (Gastropoda: Ampullariidae). Appl Entomol Zool. 2006;41:627-32.

683   21.   Seuffert ME, Burela S, Martín PR. Influence of water temperature on the activity of the

684      freshwater snail Pomacea canaliculata (Caenogastropoda: Ampullariidae) at its

685      southernmost limit (Southern Pampas, Argentina). Journal of Thermal Biology. 2010;

686      35:77-84.

687   22.   Kruatrachue M, Sumritdee C, Pokethitiyook P, Singhakaew S. Histopathological effects

688      of contaminated sediments on golden apple snail (*Pomacea canaliculata*, Lamarck 1822).

689      Bull Environ Contam Toxicol. 2011;86:610-4.

690   23.   Dreon MS, Frassa MV, Ceolín M, Ituarte S, Qiu JW, Sun J, et al. Novel animal defenses

691      against predation: a snail egg neurotoxin combining lectin and pore-forming chains that

692      resembles plant defense and bacteria attack toxins. PLoS One. 2013;8:e63782.

693      doi:10.1371/journal.pone.0063782.

694   24.   Ottaviani E, Caselgrandi E, Fontanili P, Franceschi C. Evolution, immune responses and

695      stress: studies on molluscan cells. Acta Biol Hung. 1992;43:293-8.

696   25.   Ottaviani E, Accorsi A, Rigillo G, Malagoli D, Blom JM, Tascedda F. Epigenetic

697      modification in neurons of the mollusc *Pomacea canaliculata* after immune challenge.

698      Brain Res. 2013;1537:18-26.

699   26.   Mercado Laczkó AC, Lopretto EC. Estudio cromosómico y cariotípico de p*omacea*

700      *canaliculata* (Lamarck, 1801) (Gastropoda, Ampullariidae). Revista del Museo Argentino

33

701    de Ciencias Naturales "Bernardino Rivadavia" Hidrobiología. 1998;8:15-20.

702  27.  Xu J, Han X, Li N, Yu J, Qian C, Bao Z. Analysis of genetic diversity of three geographic

703    populations of *Pomacea canaliculata* by AFLP. Acta Ecol Sin. 2009;29:4119- 26.

704  28.  Chen L, Xu H, Li H, Wu J, Ding H, Liu Y. Isolation and characterization of sixteen

705    polymorphic microsatellite loci in the golden apple snail *Pomacea canaliculata*. Int J Mol

706    Sci. 2011;12:5993-8.

707  29.  Mu X, Hou G, Song H, Xu P, Luo D, Gu D, et al. Transcriptome analysis between

708    invasive *Pomacea canaliculata* and indigenous *Cipangopaludina cahayensis* reveals

709    genomic divergence and diagnostic microsatellite/SSR markers. BMC Genet. 2015;16:12.

710  30.  Sun J, Wang M, Wang H, Zhang H, Zhang X, Thiyagarajan V, et al. De novo assembly of

711    the transcriptome of an invasive snail and its multiple ecological applications. Mol Ecol

712    Resour. 2012;12:1133-44.

713  31.  Mu H, Sun J , Fang L, Luan T, Williams GA, Cheung SG, et al. Genetic Basis of

714    Differential Heat Resistance between Two Species of Congeneric Freshwater Snails:

715    Insights from Quantitative Proteomics and Base Substitution Rate Analysis. J Proteome

716    Res. 2015;14:4296-308.

717  32.  Yang L, Cheng TY, Zhao FY. Comparative profiling of hepatopancreas transcriptomes in

718    satiated and starving *Pomacea canaliculata*. BMC Genet. 2017;18:18.

719  33.  Xiong YM, Yan ZH, Zhang JE, Li HY. Analysis of albumen gland proteins suggests

720    survival strategies of developing embryos of *Pomacea canaliculata*. Molluscan Res.

721    2017:1-6.

722  34.  Sun J, Mu H , Zhang H, Chandramouli KH, Qian PY, Wong CK, et al. Understanding the

723       regulation of estivation in a freshwater snail through iTRAQ-based comparative

724       proteomics. J Proteome res. 2013;12:5271-80.

725   35.   Sun J, Zhang H, Wang H, Heras H, Dreon MS, Ituarte S, et al. First proteome of the egg

726       perivitelline fluid of a freshwater gastropod with aerial oviposition. J Proteome Res.

727       2012;11:4240-8.

728   36.   Aplysia Genome Project. Broad Institute. Vertebrate Biology Group. 2009.

729       https://www.broadinstitute.org/aplysia/aplysia-genome-project

730   37.   Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al. The oyster genome reveals stress

731       adaptation and complexity of shell formation. Nature. 2012;490:49-54.

732   38.   Du X, Fan G, Jiao Y, Zhang H, Guo X, Huang R, et al. The pearl oyster Pinctada fucata

733       martensii genome and multi-omic analyses provide insights into biomineralization.

734       Gigascience. 2017;6:1-12.

735   39.   Takeuchi T, Kawashima T, Koyanagi R, Gyoja F, Tanaka M, Ikuta T, et al. Draft genome

736       of the pearl oyster *Pinctada fucata*: a platform for understanding bivalve biology. DNA

737       Res. 2012;19:117-30.

738   40.   Simakov O, Marletaz F, Cho SJ, Edsinger-Gonzales E, Havlak P, Hellsten U, et al.

739       Insights into bilaterian evolution from three spiralian genomes. Nature. 2013;493:526-31.

740   41.   Albertin CB, Simakov O, Mitros T, Wang ZY, Pungor JR, Edsinger-Gonzales E, et al. The

741       octopus genome and the evolution of cephalopod neural and morphological novelties.

742       Nature. 2015;524:220-4.

743   42.   Uliano-Silva M, Dondero F, Dan Otto T, Costa I, Lima NCB, Americo JA, et al. A

744       hybrid-hierarchical genome assembly strategy to sequence the invasive golden mussel

745    Limnoperna fortunei. Gigascience. 2017. doi: 10.1093/gigascience/gix128

746    43.    Adema CM, Hillier LW, Jones CS, Loker ES, Knight M, Minx P, et al. Corrigendum:

747           Whole genome analysis of a schistosomiasis-transmitting freshwater snail. Nat Commun.

748           2017;8:16153.

749    44.    Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, et al. Estimation of genomic characteristics

750           by analyzing k-mer frequency in de novo genome projects. Quantitative Biology

751           2013:arXiv:1308.2012 [q-bio.GN].

752    45.    Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms

753           and methods to estimate maximum-likelihood phylogenies: assessing the performance of

754           PhyML 3.0. Syst Biol 2010;59:307-21. doi:10.1093/sysbio/syq010.

755    46.    Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol.

756           2007;24:1586-91. doi:10.1093/molbev/msm088.

757    47.    Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome

758           comparisons dramatically improves orthogroup inference accuracy. Genome Biol.

759           2015;16:157. doi:10.1186/s13059-015-0721-2.

760    48.    Feschotte C, Wessler SR. Mariner-like transposases are widespread and diverse in

761           flowering plants. Proc Natl Acad Sci U S A 2002;99:280-5.

762    49.    Hua-Van A, Le Rouzic A, Boutin TS, Filée J, Capy P. The struggle for life of the

763           genome's selfish architects. Biol Direct. 2011;6:19.

764    50.    Werren JH. Selfish genetic elements, genetic conflict, and evolutionary innovation. Proc

765           Natl Acad Sci U S A. 2011;108:10863-70.

766    51.    Chrousos GP. Stress and disorders of the stress system. Nat Rev Endocrinol.

767    2009;5:374-81.

768    52.    Vabulas RM, Raychaudhuri S, Hayer-Hartl M. Protein folding in the cytoplasm and the

769          heat shock response. Cold Spring Harbor perspectives in biology. 2010;2:a004390.

770    53.    Chen B，Retzlaff M，Roos T，Frydman J. Cellular Strategies of Protein Quality Control.

771          Cold Spring Harbor Perspectives in Biology. 2011;3:a004374.

772    54.    Korennykh A and Walter P. Structural basis of the unfolded protein response. Annu Rev

773          Cell Dev Biol. 2012;28:251-77.

774    55.    Chambers JE and Yarbrough JD. Xenobiotic biotransformation systems in fishes. Comp

775          Biochem Physiol C. 1976;55:77-84.

776    56.    Mello DF, de Oliveira ES, Vieira RC, Simoes E, Trevisan R, Dafre AL, et al. Cellular and

777          Transcriptional Responses of *Crassostrea gigas* Hemocytes Exposed in Vitro to

778          Brevetoxin (PbTx-2) Mar Drugs. 2012;10: 583-97.

779    57.    Boutet I, Tanguy A, Moraga D. Characterisation and expression of four mRNA sequences

780          encoding glutathione S-transferases pi, mu, omega and sigma classes in the Pacific oyster

781          *Crassostrea gigas* exposed to hydrocarbons and pesticides. Mar Biol 2004;146:53-64.

782    58.    Deeley RG, Westlake C, Cole SP. Transmembrane transport of endo- and xenobiotics by

783          mammalian ATP-binding cassette multidrug resistance proteins. Physiol Rev.

784          2006;86:849-99.

785    59.    Liu C, Zhang T, Wang L, Wang M, Wang W, Jia Z, et al. The modulation of extracellular

786          superoxide dismutase in the specifically enhanced cellular immune response against

787          secondary challenge of Vibrio splendidus in Pacific oyster (*Crassostrea gigas*). Dev
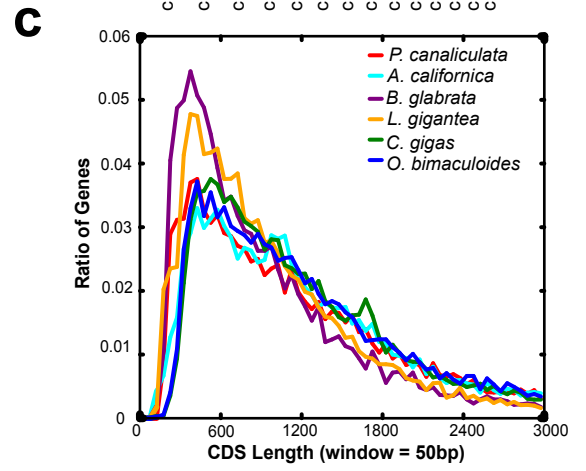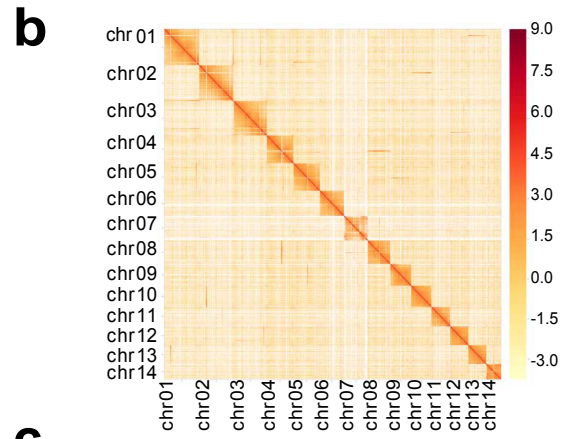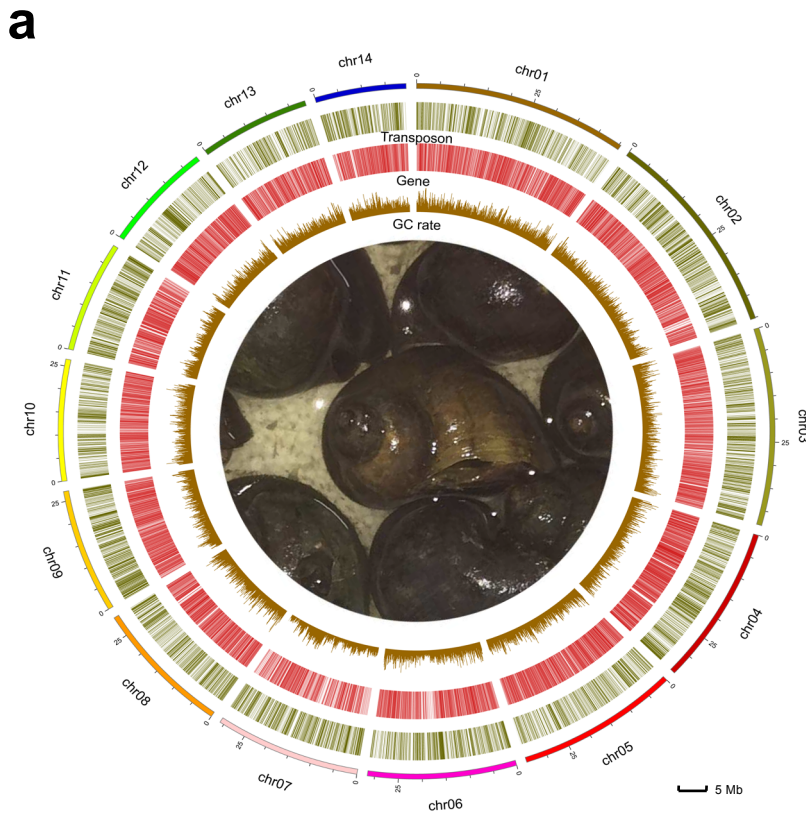
788          Comp Immunol. 2016;63:163-70.

37

789    60.    Lamb DC, Lei L, Warrilow AG, Lepesheva GI, Mullins JG, Waterman MR, et al. The first

790         virally encoded cytochrome p450. J Virol. 2009;83:8266-9.

791    61.    Urlacher VB, Girhard M. Cytochrome P450 monooxygenases: an update on perspectives

792         for synthetic application. Trends Biotechnol. 2012;30:26-36.

793    62.    Sanderson T, van den Berg M. Topic 3.1: Interactions of xenobiotics with the steroid

794         hormone biosynthesis pathway. Pure Appl Chem. 2003;75:1957-71.

795    63.    Goldstone JV, McArthur AG, Kubota A, Zanette J, Parente T, Jönsson ME, et al.

796         Identification and developmental expression of the full complement of Cytochrome P450

797         genes in Zebrafish. BMC Genomics. 2010;11:643.

798    64.    Chuang SS, Helvig C, Taimi M, Ramshaw HA, Collop AH, Amad M, et al. CYP2U1, a

799         novel human thymus- and brain-specific cytochrome P450, catalyzes omega- and

800         (omega-1)-hydroxylation of fatty acids. J Biol Chem. 2004;279:6305-14.

801    65.    Fleming I. The pharmacology of the cytochrome P450 epoxygenase/soluble epoxide

802         hydrolase axis in the vasculature and cardiovascular disease. Pharmacol Rev.

803         2014;66:1106-40.

804    66.    Zhang G, Kodani S, Hammock BD. Stabilized epoxygenated fatty acids regulate

805         inflammation, pain, angiogenesis and cancer. Prog Lipid Res. 2014;53:108-23.

806    67.    de Jong-Brink M, Boer HH, Joosse J. Mollusca. In: Adiyodi, K.G., Adiyodi,

807         R.G. (Eds.), Reproductive Biology of invertebrates. Oogenesis oviposition and

808         oosorption, vol. 1. John Wiley & Sons Ltd., New York, 1983; pp. 297-355.

809    68.    Garin CF, Heras H, Pollero RJ. Lipoproteins of the egg perivitelline fluid of *Pomacea*

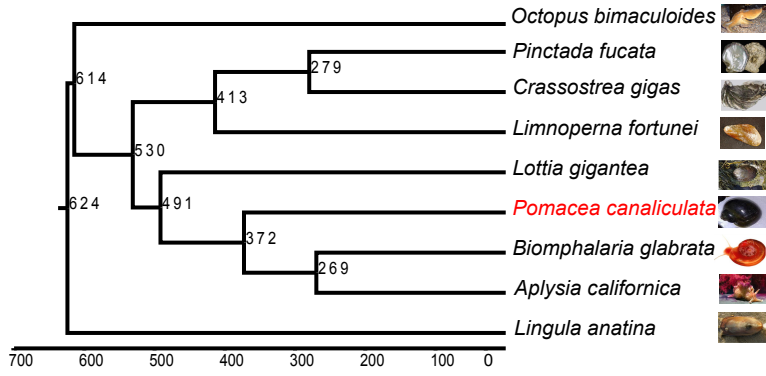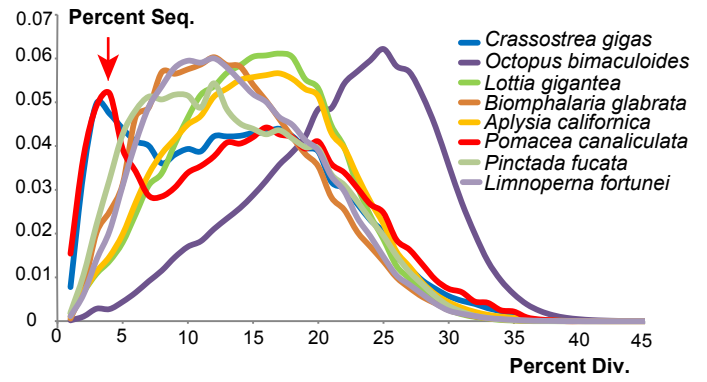810         *canaliculata* snails (Mollusca: Gastropoda). J Exp Zool. 1996;276:307-14.

811    69.    Dreon MS, Schinella G, Heras H, Pollero RJ. Antioxidant defense system in the apple

812           snail eggs, the role of ovorubin. Arch Biochem Biophys. 2004;422:1-8.

813    70.    Dreon MS, Ituarte S, Heras H. The role of the proteinase inhibitor ovorubin in apple snail

814           eggs resembles plant embryo defense against predation. PLoS One. 2010;5:e15059.

815           doi:10.1371/journal.pone.0015059.

816    71.    Cardoso AM, Cavalcante JJV, Vieira RP, Lima JL, Grieco MAB, Clementino MM, et al.

817           Gut Bacterial Communities in the Giant Land Snail *Achatina fulica* and Their

818           Modification by Sugarcane-Based Diet. Plos One. 2012;7 doi:ARTN

819           e3344010.1371/journal.pone.0033440.

820    72.    Cardoso AM, Cavalcante JJV, Cantão ME, Thompson CE, Flatschart RB, Glogauer A, et

821           al. Metagenomic Analysis of the Microbiota from the Crop of an Invasive Snail Reveals a

822           Rich Reservoir of Novel Genes. Plos One. 2012;7 doi:ARTN

823           e4850510.1371/journal.pone.0048505.

824    73.    Cabrera G, Pérez R, Gómez JM, Ábalos A, Cantero D. Toxic effects of dissolved heavy

825           metals on Desulfovibrio vulgaris and Desulfovibrio sp strains. J Hazard Mater

826           2006;135:40-6. doi:10.1016/j.jhazmat.2005.11.058.

827    74.    Finlay JA, Allan VJ, Conner A, Callow ME, Basnakova G, Macaskie LE. Phosphate

828           release and heavy metal accumulation by biofilm-immobilized and chemically-coupled

829           cells of a Citrobacter sp. pre-grown in continuous culture. Biotechnol Bioeng.

830           1999;63:87-97.

831    75.    Valls M, de Lorenzo V, Gonzalez-Duarte R, Atrian S. Engineering outer-membrane

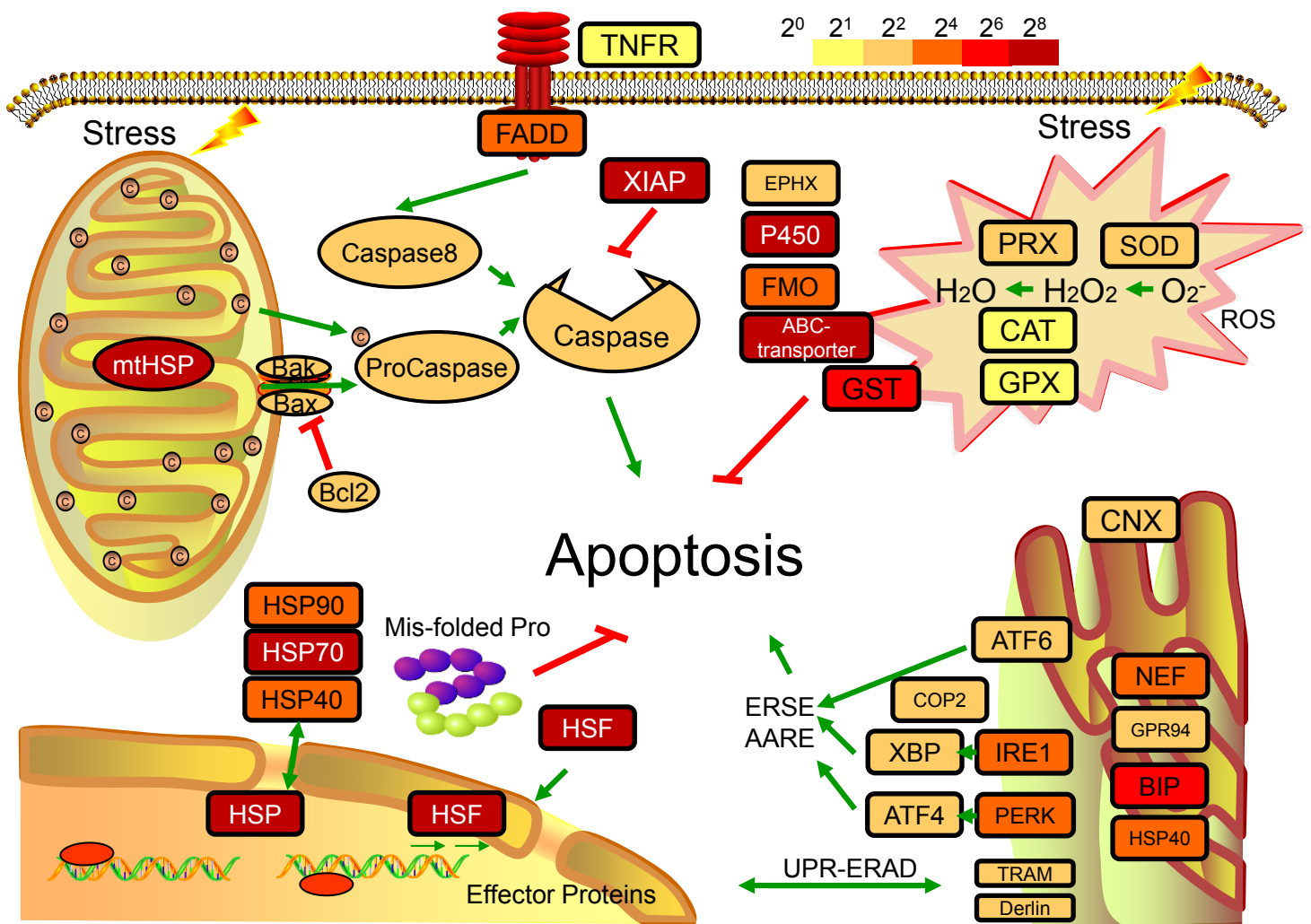832           proteins in *Pseudomonas putida* for enhanced heavy-metal bioadsorption. J Inorg

833    Biochem. 2000;79:219-23.

834    76.    Pinheiro GL, Correa RF, Cunha RS, Cardoso AM, Chaia C, Clementino MM, et al.

835    Isolation of aerobic cultivable cellulolytic bacteria from different regions of the

836    gastrointestinal tract of giant land snail *Achatina fulica*. Front Microbiol. 2015;6 doi:Artn

837    86010.3389/Fmicb.2015.00860.

838    77.    Zoetendal EG, Heilig HG, Klaassens ES, Booijink CC, Kleerebezem M, Smidt H, et al.

839    Isolation of DNA from bacterial samples of the human gastrointestinal tract. Nature

840    protocols 2006, 1(2): 870-873.

841    78.    Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids

842    Res. 1999;27:573-80.

843    79.    Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, et al. Automated eukaryotic

844    gene structure annotation using EVidenceModeler and the Program to Assemble Spliced

845    Alignments. Genome Biol. 2008;9:R7.

846    80.    Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, et al. InterProScan:

847    protein domains identifier. Nucleic Acids Res. 2005;33:W116-20.

848    81.    Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M. KEGG for integration and

849    interpretation of large-scale molecular data sets. Nucleic Acids Res. 2012;40:D109-D14.

850    82.    Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high

851    throughput. Nucleic Acids Res. 2004;32:1792-7.

852    83.    Darriba D, Taboada GL, Doallo R, Posada D. ProtTest 3: fast selection of best-fit models

853    of protein evolution. Bioinformatics. 2011;27:1164-5. doi:10.1093/bioinformatics/btr088.

854    84.    Sun J, Zhang Y, Xu T, Zhang Y, Mu H, Zhang Y, et al. Adaptation to deep-sea
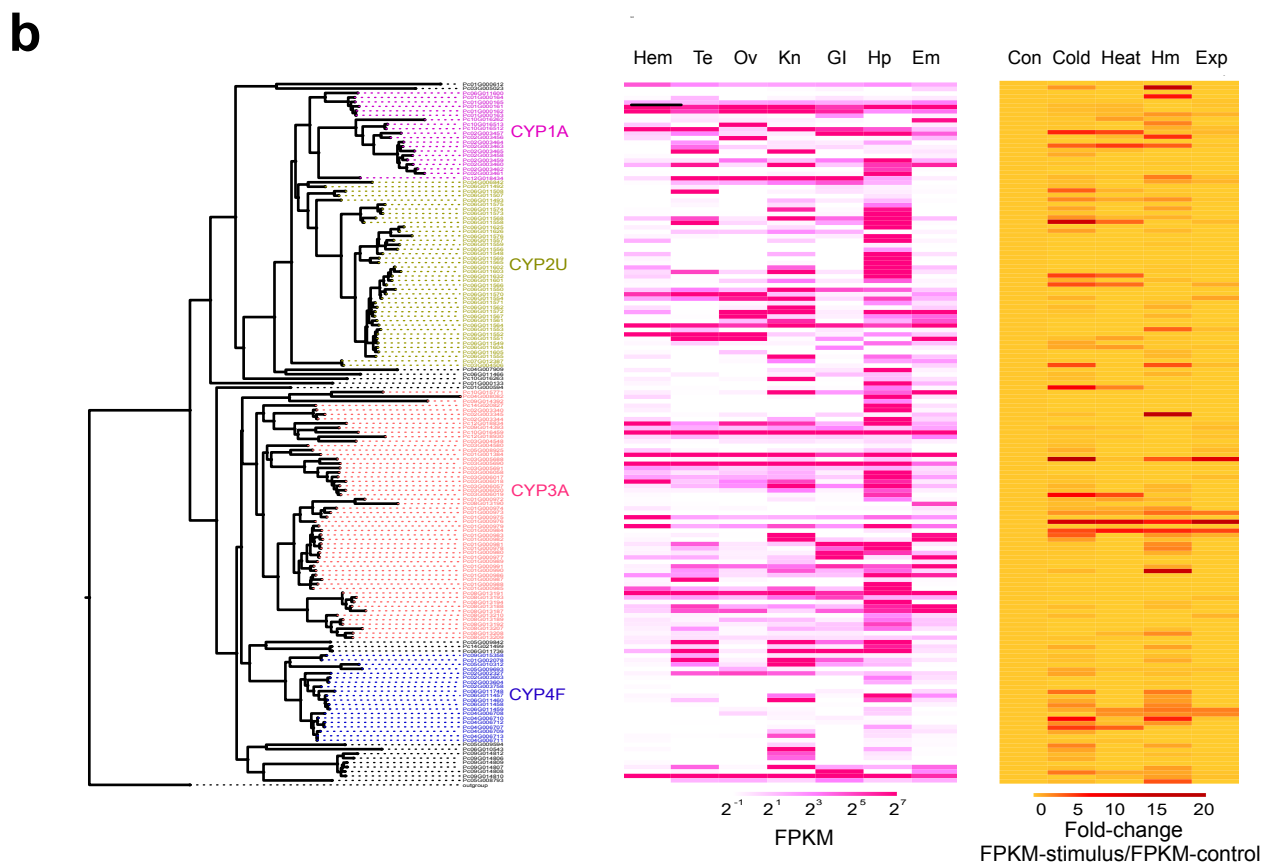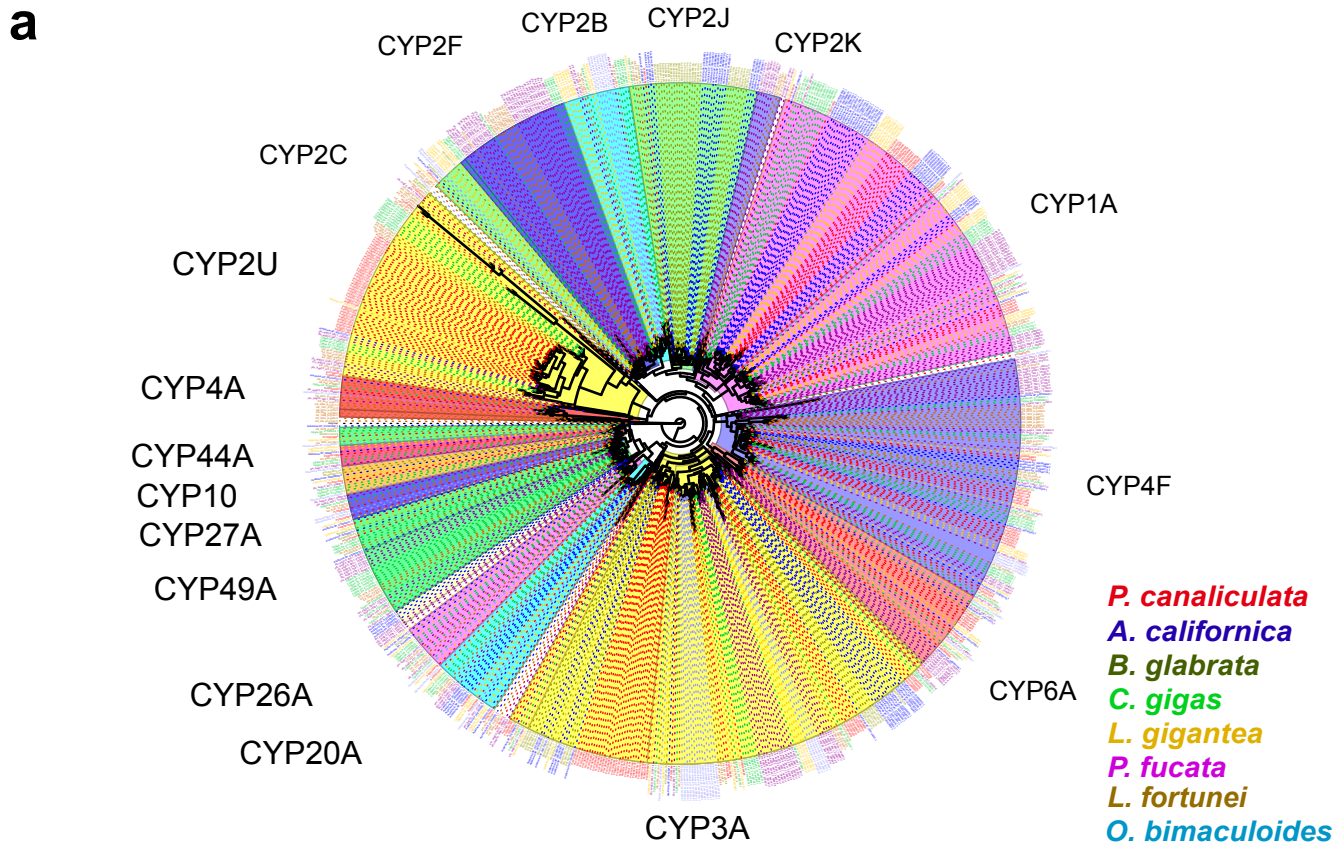
855    chemosynthetic environments as revealed by mussel genomes. Nature Ecology &

856    Evolution. 2017; 1: 121.

857  85.   Benton MJ, Donoghue PCJ, Asher RJ. in The Timetree of Life:Calibrating and

858    Constraining Molecular Clocks (eds Hedges, S. B. & Kumar, S.)35–86 (Oxford Univ.

859    Press, 2009.

860  86.   Zapata F, Wilson NG, Howison M, Andrade SC, Jörger KM, Schrödl M, et al.

861    Phylogenomic analyses of deep gastropod relationships reject Orthogastropoda. Proc Biol

862    Sci. 2014;281(1794):20141739. doi: 10.1098/rspb.2014.1739.

863  87.   Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler

864    transform. Bioinformatics. 2009;25:1754-60.

865  88.   Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. metaSPAdes: a new versatile

866    metagenomic assembler. Genome Res. 2017;27:824-34.

867  89.   Hyatt D, LoCascio PF, Hauser LJ, Uberbacher EC. Gene and translation initiation site

868    prediction in metagenomic sequences. Bioinformatics. 2012;28:2223-30.

869  90.   Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation

870    sequencing data. Bioinformatics. 2012;28:3150-2.

871  91.   Buchfink B, Chao X, Huson DH. Fast and sensitive protein alignment using DIAMOND.

872    Nat Methods. 2015;12:59-60.

873  92.   Gerlach W and Stoye J. Taxonomic classification of metagenomic shotgun sequences

874    with CARMA3. Nucleic Acids Res. 2011;39 doi:Artn E9110.1093/Nar/Gkr225.

875  93.   Kanehisa M, Goto S, Kawashima S, Okuno Y, Hattori M. The KEGG resource for

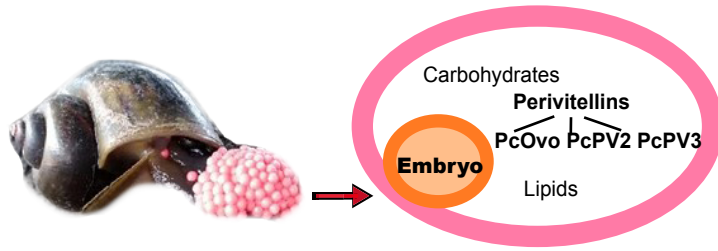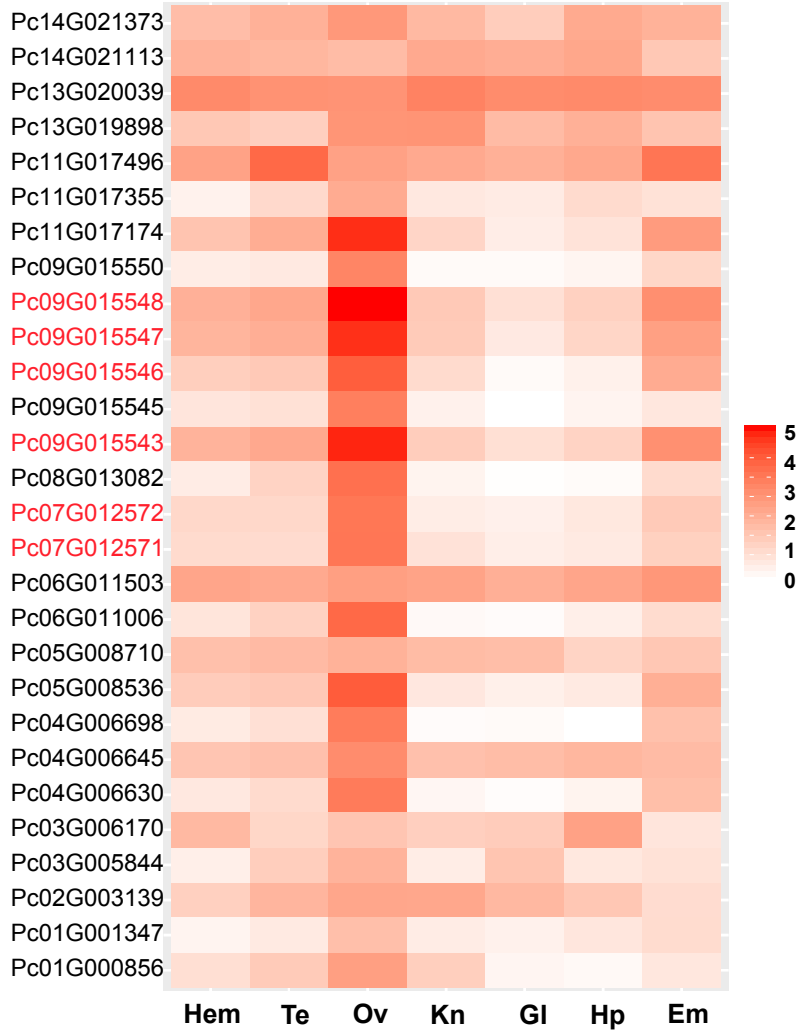876    deciphering the genome. Nucleic Acids Res. 2004;32:D277-80.

877    94.    Yin Y, Mao X, Yang J, Chen X, Mao F, Xu Y. dbCAN: a web resource for automated

878          carbohydrate-active enzyme annotation. Nucleic Acids Res. 2012;40:W445-51.

879    95.    Eddy SR. Accelerated Profile HMM Searches. Plos Comput Biol. 2011;7 doi:ARTN

880          e100219510.1371/journal.pcbi.1002195.

881    96.    Qin JJ, Li YR, Cai ZM, Li SH, Zhu JF, Zhang F, et al. A metagenome-wide association

882          study of gut microbiota in type 2 diabetes. Nature. 2012;490:55-60.

883    97.    Liu C, Zhang Y, Ren Y, Wang H, Li S, Jiang F, et al. Supporting data for "The genome of the

884          golden apple snail Pomacea canaliculata provides insight into stress tolerance and invasive

885          adaptation". GigaScience Database. 2018. http://dx.doi.org/10.5524/100485

**a**



**b**



**c**

# a



# b

**a**

CYP2F  CYP2B  CYP2J  CYP2K

CYP2C

CYP2U

CYP4A

CYP44A
CYP10
CYP27A

CYP49A

CYP26A

CYP20A

CYP3A

CYP1A

CYP4F

CYP6A

*P. canaliculata*
*A. californica*
*B. glabrata*
*C. gigas*
*L. gigantea*
*P. fucata*
*L. fortunei*
*O. bimaculoides*

**b**

Hem  Te  Ov  Kn  GI  Hp  Em

Con  Cold  Heat  Hm  Exp

CYP1A

CYP2U

CYP3A

CYP4F

$2^{-1}$  $2^1$  $2^3$  $2^5$  $2^7$
FPKM

0  5  10  15  20
Fold-change
FPKM-stimulus/FPKM-control

## a

**P. canaliculata**



## b

Click here to access/download

**Supplementary Material**

Supplemental_Information-final.doc