# GigaScience

## Genome sequencing of deep-sea hydrothermal vent snails reveals adaptions to extreme environments
### --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-20-00187 |
| Full Title: | Genome sequencing of deep-sea hydrothermal vent snails reveals adaptions to extreme environments |
| Article Type: | Data Note |
| Funding Information: | National Key R&D Programme of China (No. 2018YFC0310702) — Dr. Xiang Zeng |

**Abstract:**

Background

The scaly-foot snail ( Chrysomallon squamiferum ) is highly adapted to deep-sea hydrothermal vents and drew people's interest once it was found. However, the limited information on its genome inmedes related research and understanding of its adaptation to deep-sea hydrothermal vents.

Findings

Here, we report the whole-genome sequencing and assembly of the scaly-foot snail and another snail ( Gigantopelta aegi ), which inhabits similar environments. Using ONT, 10X genomic, and Hi-C technologies, we obtained a chromosome-level genome of C. squamiferum with an N50 size of 20.71 Mb. By constructing a phylogenetic tree, we found that these two deep-sea snails were independent of other snails, and their divergence from each other occurred approximately 66.3 million years ago. Comparative genomic analysis showed that different snails have diverse genome sizes and repeat contents. Deep-sea snails have more DNA transposons and LTRs, but fewer LINEs, than other snails. Gene family analysis revealed that deep-sea snails experienced stronger selective pressures than shallow-water snails, and the nervous system, immune system, metabolism, DNA stability, antioxidation and biomineralization-related gene families were significantly expanded in scaly-foot snails. We also found 251 class II histocompatibility antigen H2-Aal, which uniquely exist in the Gigantopelta aegi genome, which is important for investigating the evolution of MHC genes.

Conclusion

Our study provides new insights into deep-sea snail genomes and valuable resources for further studies.

| | |
|---|---|
| Corresponding Author: | Yaolei Zhang<br>BGI<br>Qingdao, CHINA |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | BGI |
| Corresponding Author's Secondary Institution: | |
| First Author: | Yaolei Zhang |
| First Author Secondary Information: | |
| Order of Authors: | Yaolei Zhang |
| | Xiang Zeng |
| | Lingfeng Meng |
| | Guangyi Fan |

| | Jie Bai |
| --- | --- |
| | Jianwei Chen |
| | Yue Song |
| | Inge Seim |
| | Congyan Wang |
| | Zenghua Shao |
| | Nanxi Liu |
| | Haorong Lu |
| | Xiaoteng Fu |
| | Liping Wang |
| | Xin Liu |
| | Shanshan Liu |
| | Zongze Shao |

| 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**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible. | Yes |

| | |
|---|---|
| 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 |

main text

1    **Genome sequencing of deep-sea hydrothermal vent snails reveals adaptions to**

2    **extreme environments**

3    Xiang Zeng[1†], Yaolei Zhang[2,3,4†], Lingfeng Meng[2†], Guangyi Fan[2,3,6], Jie Bai[3], Jianwei

4    Chen[2], Yue Song[2] , Inge Seim[7,8], Congyan Wang[2], Zenghua Shao[2], Nanxi Liu[3], Haorong

5    Lu[3], Xiaoteng Fu[1], Liping Wang[1],Xin Liu[2,3,5], Shanshan Liu[2*],  Zongze Shao[1*]

6

7    [1]Key Laboratory of Marine Biogenetic Resources, Third Institute of Oceanography,

8    Ministry of Natural Resources, No.178 Daxue Road, Xiamen 361005, China.

9    [2]BGI-Qingdao, BGI-Shenzhen, Qingdao 266555, China

10   [3]BGI-Shenzhen, Shenzhen, 518083, China

11   [4]Department of Biotechnology and Biomedicine, Technical University of Denmark,

12   Lyngby, 2800, Denmark

13   [5]China National GeneBank, BGI-Shenzhen, Shenzhen 518120, China

14   [6]State Key Laboratory of Agricultural Genomics, BGI-Shenzhen, Shenzhen 518083,

15   China.

16   [7]Integrative Biology Laboratory, College of Life Sciences, Nanjing Normal University, Nanjing,

17   210046, China

18   [8]Comparative and Endocrine Biology Laboratory, Translational Research Institute-Institute of

19   Health and Biomedical Innovation, School of Biomedical Sciences, Queensland University of

20   Technology, Woolloongabba, 4102, Australia

21

22   [†] These authors contributed equally.

23   *Correspondence: shaozz@163.com (Z.S.); liushanshan@genomics.cn (S.L.)

24

25

26

27

28

29

30

31

32

33

34 **Abstract**

35 **Background**

36 The scaly-foot snail (*Chrysomallon squamiferum*) is highly adapted to deep-sea

37 hydrothermal vents and drew people's interest once it was found. However, the limited

38 information on its genome inmedes related research and understanding of its adaptation

39 to deep-sea hydrothermal vents.

40 **Findings**

41 Here, we report the whole-genome sequencing and assembly of the scaly-foot snail and

42 another snail (*Gigantopelta aegi*), which inhabits similar environments. Using ONT, 10X

43 genomic, and Hi-C technologies, we obtained a chromosome-level genome of *C.*

44 *squamiferum* with an N50 size of 20.71 Mb. By constructing a phylogenetic tree, we

45 found that these two deep-sea snails were independent of other snails, and their

46 divergence from each other occurred approximately 66.3 million years ago. Comparative

47 genomic analysis showed that different snails have diverse genome sizes and repeat

48 contents. Deep-sea snails have more DNA transposons and LTRs, but fewer LINEs, than

49 other snails. Gene family analysis revealed that deep-sea snails experienced stronger

50 selective pressures than shallow-water snails, and the nervous system, immune system,

51 metabolism, DNA stability, antioxidation and biomineralization-related gene families

52 were significantly expanded in scaly-foot snails. We also found 251 class II

53 histocompatibility antigen H2-Aal, which uniquely exist in the *Gigantopelta aegi* genome,

54 which is important for investigating the evolution of MHC genes.

55 **Conclusion**

56 Our study provides new insights into deep-sea snail genomes and valuable resources for

57 further studies.

58

59 **Keywords:** Deep-sea snails;  Genome assembly; Comparative genomics;

60 Biomineralization;

61

62

63

**Background**

64

65 The discovery of deep-sea hydrothermal vents in the late 1970s expanded our

66 knowledge of the extent of life on Earth [1]. Deep-sea macrobenthos, which are animals

67 that inhabit deep-sea hydrothermal vents, face high hydrostatic pressure, variable

68 temperatures and pH, and high levels of hydrogen sulphide, methane, and heavy metals

69 [2]. To date, the literature contains a limited number of studies on the genetics of

70 macrobenthos. A recent report on the genome of deep-sea hydrothermal vent/cold seep

71 mussels (*Bathymodiolus platifrons*) showed that, while most of the genes present in a

72 related shallow-water mussel (*Modiolus philippinarum*) have been retained, many gene

73 families have expanded in the *B. platifrons* genome. These include families that are

74 associated with stabilising protein structures, removing toxic substances from cells, and

75 the immune response to symbionts [3].

76 Gastropods represent the largest class of the phylum Mollusca, with different

77 estimates of diversity varying from 80,000 to 150,000 species [4]. More than 218

78 gastropod (i.e. snail and slugs) species have been described from chemosynthetic

79 ecosystems (i.e. solely rely on endosymbiotic bacterias for sustenance), of which more

80 than 138 are believed to be endemic to these ecosystems [5]. Gastropods are an important

81 component of the fauna in hydrothermal vents in terms of abundance and biomass [6].

82 Due to the lack of samples and fossil evidence, studies on the evolution and adaptation of

83 deep sea chemosynthetic gastropods are very limited. The scaly-foot snail (*Chrysomallon*

84 *squamiferum*) is only found in hydrothermal vents at a depth of ~3,000-metres in the

85 Indian Ocean. There are two types of varieties without genetic differences: black (due to

86 greigite, which is an iron sulphide mineral that covers its exterior) scaly-foot individuals

87 from the Kairei field on the central Indian ridge and Longqi field on the Southwest Indian

88 ridge, and white scaly-foot individuals from the Solitaire field on the Central Indian

89 Ridge and the Carlsberg ridge on the Northwest Indian Ridge [7]. *C. squamiferum* was

90 included in the International Union for Conservation of Nature (IUCN) Red List of

91 Endangered Species on July 18, 2019 [8]. *Gigantopelta* spp. is a major megafaunal

92 gastropod genus in some hydrothermal fields. The genus includes two species,

93 *Gigantopelta chessoia* sp. nov. from East Scotia Ridge and *Gigantopelta aegis* sp. nov.

94 from the Southwest Indian Ridge [6]. Both *Chrysomallon* and *Gigantopelta* are members

95  of the family Peltospiridae. They live in high-density aggregations and share several

96  features, such as a large body size (up to > 45 mm, compared to typical sizes in other taxa

97  of 10-15 mm, a 10-50 fold increase in body volume) and an enlarged oesophageal gland

98  [9].

99  In this study, we sequenced and analysed chromosome-level genomes of the white

100  scaly-foot snail *Chrysomallon squamiferum* (*C. squamiferum*, **Figure 1a**) from the

101  Carlsberg ridge on the Northwest Indian Ridge and *Gigantopelta aegis* (*G. aegis*, **Figure**

102  **1a**) from the Southwest Indian Ridge. We gained insights into the evolution, gene family

103  expansions, and adaptations of these extremophile gastropods.

104

105  **Data Description**

106  **Genome assembly and annotation**

107  The *C. squamiferum* genome was sequenced using a combination of sequencing

108  libraries – 10X Genomics, Oxford Nanopore Technology (ONT), and Hi-C – to generate

109  ~369.03 Gb of raw data (**Table S1**). Due to the limited sample material, *G. aegis* was

110  sequenced from whole genome shotgun libraries (with 350 bp to 10 kb inserts on the

111  BGISEQ-500 platform) to generate 910.08 Gb of raw data (**Table S2**). The genome of *C.*

112  *squamiferum* was assembled with long ONT reads by using Canu [10] and WTDBG [11].

113  After polishing the genome with 10X Genomics sequencing data, a 454.58 Mb assembly

114  (a little smaller than the estimated genome size: 495 Mb, **Figure S1**) with 6,449 contigs

115  and an N50 of 541.32 kb was generated (**Table S3**). Next, Hi-C data were used to anchor

116  the assembly, yielding a 16-chromosome assembly (**Figure 1b**). This effort increased the

117  N50 size to ~20.71 Mb (**Table 1**). The 16 chromosomes cover ~80% of the whole

118  genome, and the average length, maximal length, and minimal length of the 16

119  chromosomes were 22.67, 46.78, and 10.64 Mb, respectively, (**Table S4**). A

120  Benchmarking Universal Single-Copy Orthologs (BUSCO) completeness score of 94.8%

121  for this genome suggests that it is of good quality (**Table S5**). This is the first

122  chromosome-level deep-sea snail genome assembly to date. An approximately 1.29 Gb (a

123  little smaller than the estimated genome size:1.50 Gb, **Figure S1**) genome assembly of *G.*

124  *aegis* with a scaffold N50 of 120.96 kb (**Table S6**) and a BUSCO completeness score of

125  88.4% (**Table S7**) was obtained using Platanus [12]. After masking repeat elements, we

126    employed homologous and *de novo* prediction methods to construct gene models for the

127    two genomes, obtaining 28,781 *C. squamiferum* genes and 25,601 *G. aegis* genes (**Tables**

128    **S8** and **S9**). The gene sets were functionally annotated using KEGG, Swiss-Prot, InterPro,

129    and TrEMBL (**Tables S10** and **S11**).

130

131    **Genome sizes and repeat contents.**

132    The genome assembly sizes of *C. squamiferum* (~455.36 Mb) and *G. aegis* (~1.29 Gb)

133    add to previous studies on freshwater snails (~916 Mb (*Biomphalaria glabrata*) [13] and

134    ~440 Mb (*Pomacea canaliculate*) [14]), which suggests that there is significant genome

135    size diversity within snails (**Figure 1c**). In the absence of ploidy effects [15, 16],

136    differences in genome size often stem from the accumulation of various repetitive

137    elements. A comparison of repeat elements (**Figure 1c** and **Table S12**) supports this

138    contention. The genomes of *C. squamiferum* and *P. canaliculate* (smaller genome sizes)

139    contained fewer repeats than *B. glabrata* and *G. aegis*, whereas *G. aegis* had more repeats

140    than *B. glabrata* (**Figure 1d**). This suggests that snail genome sizes correlate with repeat

141    content. Despite the similar genome sizes of *C. squamiferum* and *P. canaliculata,* their

142    genome landscapes are distinct. For example, ~10.17% of the *C. squamiferum* genome

143    consists of  tandem repeats compared to ~2.89% in *P. canaliculata* (**Table S12**). DNA

144    transposons and LTRs comprise ~17.73% and ~5.99% of the *C. squamiferum* genome,

145    respectively, but only ~6.84% and ~3.53% in *P. canaliculata*. LINEs make up ~8.63% of

146    the *P. canaliculata* genome compared to ~5.65% in *C. squamiferum*. Similarly, although

147    the larger *G. aegis* and *B. glabrata* genomes have similar proportions of tandem repeats,

148    *G. aegis* has a higher percentage of DNA transposons (~32.15% versus ~20.20%) and

149    LTRs (~13.32% versus ~3.75%). LINEs make up ~23.93% of the *B. glabrata* genome

150    compared to ~11.51% in *G. aegis*. Taken together, these data suggest that deep-sea

151    hydrothermal vent snail genomes have more DNA transposons and LTRs and fewer

152    LINEs than their freshwater counterparts. In particular, DNA/CMC-EnSpm,

153    DNA/TcMar−Tc1, and DNA/DNA were the main factors that caused the differences in

154    DNA transposon content in the four snail genomes (**Figure 1d**). We found that LINE/L2,

155    LINE/RTE-BovB, LINE/LINE, and LINE/CR1 were much higher in fresh-water snail

156    genomes than in deep-sea snails. Although most of the precise functions of these repeats

157  have not been studied in-depth, repeats have been thought to have a regulatory function

158  in related genes that play an important role in the life cycle. Thus, we might infer that the

159  expansion of DNA transposons and LTRs, as well as the abandonment of some LINEs,

160  may be closely associated with adaptation to extreme environments for deep-sea snails.

161

162  **Construction of phylogenetic relationships for deep-sea snails**

163  To determine the phylogenetic relationships between deep-sea snails and other

164  molluscs, we compared two mussels, two freshwater snails, and two shallow-water snails.

165  The California two-spot octopus and the freshwater leech *Helobdella robusta* were used

166  as the outgroup. We identified 26,668 gene families in the ten species examined (**Table**

167  **S13**). Phylogenetic trees were constructed from 406 shared single-copy orthologs. Both

168  ML and Bayesian methods revealed the same topology (**Figure 2a** and **Figure S2**), which

169  is consistent with a recent study [14]. In the tree, mussels and snails are clearly separated

170  and the two deep-sea snails are located on the same branch and are independent of other

171  snails (although their genome sizes are quite different). We estimated that *C.*

172  *squamiferum* and *G. aegis* diverged from a common ancestor approximately 66.3 million

173  years ago (MYA). This time is consistent with the most recent 'mass extinction', at the

174  end of the Cretaceous geological period ~66 MYA, where ~76% of species became

175  extinct [17].

176

177  **Demographic histories of the deep-sea snails**

178  As the speciation of the two deep-sea snails may be related to geological events (*see*

179  *above*), we estimated their historical effective population size ($N_e$) using whole-genome

180  genetic variation. We identified ~3.51 and ~3.19 million heterozygous SNPs with

181  nucleotide diversities of 0.0077 and 0.0025 for *C. squamiferum* and *G. aegis*, respectively.

182  We estimated $N_e$ changes using the pairwise sequential Markovian coalescent (PSMC)

183  method, which can infer demography from approximately 20,000 to 1 million years ago

184  (MYA) [18]. The effective population sizes of *C. squamiferum* and *G. aegis* – species

185  derived from different geographical locations in the Indian Ocean – are distinct (**Figure**

186  **2b**). The demographic history of *G. aegis* decreased until ~250 KYA (thousand years

187  ago), followed by an $N_e$ increase, from ~50,000 to 450,000 individuals, 20,000 years ago.

188 Several cycles of $N_e$ increase and decrease have been observed for *C. squamiferum*, with
189 the effective population size recovering and stabilising at 35,000 individuals from 70
190 KYA onwards. Thus, although deep-sea habitats are inhabited, deep-sea snail populations
191 are sensitive to habitat disturbances, such as major geological events. Unfortunately, the
192 *C. squamiferum* population size has dramatically decreased recently due to deep-sea
193 mining [8], which has made this species endangered.
194

**Evolution of single-copy orthologous genes**

196 The evolution and expression of single-copy orthologous genes are unique features of
197 organisms. To explore the evolutionary rate of single-copy orthologous genes, we
198 calculated the synonymous substitution rate (*Ka*) and nonsynonymous substitution rate
199 (*Ks*) values of 1,324 single-copy orthologous genes shared by the two deep-sea snails,
200 one shallow-water snail (*L. gigantea*), and two freshwater snails (*B. glabrata* and *P.*
201 *canaliculate*) (**Figure 2c**, **Figure S3,** and **Table S15**). We found that the *Ka* values of the
202 two deep-sea snails (average: 0.37 and 0.41) were higher than that of the shallow-water
203 snail (0.35) but similar to those of two freshwater snails (0.39 and 0.41), which suggests
204 that the genes of deep-sea and freshwater snails both evolved faster after their divergence
205 from shallow-water snails. The *Ks* values of the deep-sea (3.34 and 3.09) and freshwater
206 (3.19 and 3.24) snails were also similar and lower than those of the shallow-water snails
207 (3.72). Additionally, the *Ka*/*Ks* values of the deep-sea snails (average: 0.13 and 0.15)
208 were approximately ~20% and ~40% higher than those of the shallow-water snails (0.11);
209 from this we could infer that deep-sea snails have experienced stronger selective
210 pressures, possibly to allow adaptation to life in hydrothermal vents.
211

**Expanded gene families in deep-sea snail genomes**
213 *Nervous system*
214 Using CAFÉ [19] (*see details in methods*), we identified two significantly (*p*-value <
215 0.01) expanded gene families in the two deep-sea snail genomes compared to the
216 freshwater and shallow-water snails. BTB/POZ domain-containing protein 6 (*BTBD6*)
217 had 56 copies in *C. squamiferum* and 35 copies in *G. aegis*, while fewer than 5 copies
218 were found in the four other snail species examined (**Figure 3a**). We found 17 *BTBD6*

219   genes on chromosome 16 of *C. squamiferum,* and these genes showed traces of tandem

220   duplications (**Figure 3b**). In *G. aegis*, we also found several tandem gene clusters

221   (**Figure 3b**). *HTR4* (5-hydroxytryptamine receptor 4) had 12 copies in *C. squamiferum*

222   and 18 copies in *G. aegis*, while only one copy was found in the other snail species

223   (**Figure 3c**). The expansions of these gene families also displayed tandem duplications

224   (**Figure S4**). Both of these genes have roles in neuroregulation; *BTBD6*, which is an

225   adaptor of the Cul3 ubiquitin ligase complex, is essential for neural differentiation [20],

226   while *HTR4* modulates the release of various neurotransmitters[21]. A previous study

227   revealed that a large unganglionated nervous system exists in *C. squamiferum* [7] (**Figure**

228   **3d**). We speculate that the expansions of *BTBD6* and *HTR4* contribute to this system by

229   sustaining life in a deep-sea environment.

230

231   *Metabolism related genes*

232     *C. squamiferum* houses abundant endosymbionts in its greatly enlarged oesophageal

233   gland, and these endosymbionts supply nutrition for its host. KEGG enrichment analysis

234   on the 183 expanded gene families of *C. squamiferum* revealed significant enrichment

235   for metabolic pathways ($q$-value < 0.0001, **Table S16**). Among these genes, nine gene

236   families encoded enzymes in the glycolysis pathway and citrate cycle (TCA cycle). For

237   example, isocitrate dehydrogenase (IDH), which catalyses the oxidative decarboxylation

238   of isocitrate to produce α-ketoglutarate and $CO_2$, expanded significantly ($p < 0.01$). The

239   α-ketoglutarate dehydrogenase complex (OGDC) consists of three components:

240   oxoglutarate dehydrogenase (OGDH), dihydrolipoyl succinyltransferase (DLST), and

241   dihydrolipoyl dehydrogenase (DLD), among which OGDH expanded ($p < 0.01$, **Figure**

242   **4a**). IDH and OGDC are two rate-limiting enzymes in the TCA cycle, and related

243   biochemical reactions are irreversible (**Figure 4b**).

244

245   *Defence mechanisms*

246     Endosymbiotic bacteria are critical for snail life in deep-sea hydrothermal vent

247   ecosystems [22]. These bacterial taxa are largely restricted to chemosynthetic

248   environments, with some being exclusive to vents [23]. The different genome

249    evolutionary processes of *C. squamiferum* and  *G. aegis* may generate diverse defence
250    mechanisms that are used to adapt to different gene evolutions.

251      A total of 183 expanded gene families were identified in the *C. squamiferum* genome.
252    As expected, many of these have roles in the immune system. However, unlike the
253    freshwater snail *B. glabrata* [13] and deep-sea mussels [3], we did not detect an
254    expansion of the Toll-like receptor 13 (*TLR13*) gene family, but identified other
255    expanded gene families (**Figure 4a**). For example, increased  expression of thioredoxin 1
256    (*Txn1*; 22 copies in *C. squamiferum*) plays a pivotal role in T-cell activation in mice [24].
257    Although T-cell related adaptive immunity only appears in vertebrates, the existence and
258    expansion of this gene may assist the innate immune system of *C. squamiferum.*
259    Glutamine-fructose-6-phosphate transaminase (*GAFT*; 21 copies in *C. squamiferum*)
260    promotes the biosynthesis of chitin [25, 26], which is one of the stable components of the
261    crustacean shell, and provides protection against predation and infection.

262      We identified expanded gene families that maintain the stability of nucleic acids and
263    proteins, such as heat shock protein 90 (Hsp90; 13 copies in *C. squamiferum***, Figure 4a**),
264    which protects proteins against heat stress [27]; the single-stranded DNA-binding
265    proteins, encoded by SSB genes (19 copies in *C. squamiferum,* and 1 copy in other
266    species, **Figure 4a)**, which are required for DNA replication, recombination, and repair
267    processes [28]; and catalase (*CAT*, 6 copies *C. squamiferum*; **Figure 4c**), which is critical
268    in the response against oxidative stress [29]. The elevated levels of heavy metals and
269    sulphide, and high temperatures in hydrothermal vents are likely to greatly increase the
270    risk of DNA damage and misfolded proteins. Thus, these expanded gene families may
271    help these snails resist environmental stress.

272      We also found a special gene family, deleted in malignant brain tumours 1 (*DMBT1*),
273    expanded (70 copies, **Figure 4a**) in the *C. squamiferum* genome. *DMBT1* can encode
274    three glycoproteins (DMBT1 (deleted in malignant brain tumours 1 protein), SAG
275    (salivary agglutinin), and GP340 (lung glycoprotein-340)) and belongs to the scavenger
276    receptor cysteine-rich (SRCR) protein superfamily of the immune system [30]. This gene
277    consists of the SRCR, CUB, and zona pellucida domains, and all 70 copies of this gene in
278    *C. squamiferum* contain the SRCR domain, which can bind a broad range of pathogens,
279    including cariogenic *streptococci*, *Helicobacter pylori,* and HIV [31]. However, previous

280  studies have shown that SRCR domains that contain proteins are commonly expressed in
281  the shell martrix[32] and have been proven to be potentially linked to
282  biomineralisation[33], which would be associated with the foot scales of *C. squamiferum*.
283  Nonetheless, the expansion of this gene family will either strengthen the immune ability
284  or help construct the scale armour of these snails.

285  Correspondingly, we identified the expansion of 198 gene families (containing 4,515
286  genes) in the *G. aegis* genome. These families were enriched in 58 KEGG pathways (*q*-
287  value < 0.05) (**Table S17**). The majority of these pathways were associated with the
288  immune and disease response, and included terms such as 'infection', 'NOD-like receptor
289  signalling', 'Tumour necrosis factor (TNF) signalling pathway', and 'Antigen processing
290  and presentation' (**Figure S5**). Surprisingly, we found 251 copies of the H-2 class II
291  histocompatibility antigen, A-U alpha chain-like (H2-Aal) genes, which is one of the
292  major histocompatibility complex (MHC) genes in vertebrates [34]. The existence and
293  super expansion of this gene family in an invertebrate positions *G. aegis* as useful for the
294  study of immune system evolution.

295

296  **Discussion**

297  Molluscs are a highly diverse group, and their high biodiversity makes them an
298  excellent model to address topics such as biogeography, adaptability, and evolutionary
299  processes [35]. Members of the family Peltospiridae in the gastropod clade Neomphalina
300  are restricted to chemosynthetic ecosystems and, so far, are only known from hot vents
301  [6]. Based on the chromosome-scale genome assembly analyses of the scaly-foot snail (*C.*
302  *squamiferum*) and deep-sea snail (*G. aegi*), which both belong to the Peltospiridae family
303  from chemosynthetic ecosystems, our results provide insight into the possible evolution
304  and adaptation mechanisms of hydrothermal vent animals.

305  By constructing a phylogenetic tree, we found that snails diverged from other molluscs
306  approximately 555.2 MYA (**Figure 2a**). These two deep-sea snails were found to be
307  independent of other shallow-water snails around 536.6 MYA and diverged from each
308  other approximately 66.3 MYA. This evolutionary time frame implies that the last
309  common ancestor of all molluscs (LCAM) already lived before the infamous Cambrian
310  Explosion (530-540 MYA), which was speculated by the palaeobiological hypothesis

311 [36]. It also elucidated that deep-sea gastropod lineages originated at least around 540

312 MYA and diverged from other gastropods in the same age of the oldest molluscs taxons,

313 Aculifera and Conchifera [37, 38]. The deep sea gastropod lineages were also confirmed

314 by the phylogenetic analysis of mitogenomes [39]. Further conceived by the evolutionary

315 rate of single-copy orthologous genes, deep sea gastropod lineages have experienced

316 stronger selective pressures than shallow-water snails (**Figure 2c**). At the end of the

317 Cretaceous geological period, ~66 MYA, *C. squamiferum* and *G. aegis* diverged from

318 each other and had different historical effective population sizes (*Ne*) later (**Figure 2b**).

319 This indicated that they faced different environmental factors and selected pressures.

320 The transposable elements (TEs) play multiple roles in driving genome evolution in

321 eukaryotes[40]. The genome sizes of four representative snails were quite divergent (440

322 Mb-1.29 Gb). The deep sea snail *G.aegi* had the largest genome (1.29Gb), with the

323 highest percentage of DNA transposons (32.15%). Deep sea snails (*C squamiferum* and

324 *G.aegi*) had more DNA transposons and LTRs than other snails, but fewer LINEs. LTR

325 class has been identified as the main contributor to open chromatin regions and

326 transcription factor binding sites[41] [42]. LINEs may be associated with the

327 duplicability of genomic regions, which are always shared between related lineages[43].

328 It also indicated that the evolution of deep sea snail linages depends more on adaptive

329 needs than on a region-specific feature shared between lineages.

330 Specifically, we analysed expanded gene families in deep-sea snail genomes

331 (**Figure 4a**). They both significantly expanded the nervous system, especially BTB/POZ

332 domain-containing protein 6 (*BTBD6*) and 5-hydroxytryptamine receptor 4 (*HTR4*),

333 which are involved in the neuroregulation of activities, such as movement, predation, and

334 resistance to environmental change. As for the chemosynthetic snails, they both expanded

335 immune system-related genes. In the *C. squamiferum* genome, the expansions of

336 thioredoxin 1 (Txn1) and Glutamine-fructose-6-phosphate transaminase (GAFT) were

337 found. In the *G. aegi* genome, different immune and disease response genes were

338 expanded; for example, the major histocompatibility complex (MHC) genes, H-2 class II

339 histocompatibility antigen-like (H2-Aal). These expanded gene families were different

340 from fresh water snails and deep sea mussels.

341      Interestingly, in the scaly-foot snail (*Chrysomallon squamiferum*) genome, it
342    significantly enriched the main metabolic pathways, including the glycolysis pathway
343    and citrate cycle (TCA cycle), expanded the single-stranded DNA-binding protein (*SSB)*
344    family to stabilise ssDNA, heat shock protein 90 (*Hsp90*) to keep proteins folded
345    properly, and catalase (*CAT*) to avoid free radical generation by the peroxide. These gene
346    expansions might have provided deep sea snails with better immune reactions with
347    symbionts, rapid nerve signal conduction, stronger metabolism, and effective resistance
348    for adaptation to their hydrothermal vent habitat.

349    In particular, we found that *DMBT1* gene families that encode multiple SRCR domains
350    expanded significantly. These genes play important roles in immune response and
351    biomineralisation, both of which are vital for deep sea snails.

352    In conclusion, the genome analysis of  deep-sea snails (*Chrysomallon squamiferum*
353    and *Gigantopelta aegi*) from hydrothermal vents revealed their evolution and different
354    molecular adaptation to extreme environments and will be a valuable resource for
355    studying the evolution of inveterbrates.

356

**Materials and Methods**

**Sample collection and DNA isolation**

*Chrysomallon squamiferum* samples were obtained from the Daxi hydrothermal field (60.5°W 6.4°N, 2919m depth) on the Carlsberg Ridge, northwest Indian Ocean, in March 2017 during the Chinese DY38[th] cruise. *Gigantopelta aegis* samples were obtained from the Longqi vent field (37.8°S, 49.9°E, 2,780 m) on the southwest Indian ridge in March 2015 during the Chinese DY35[th] cruise. DNA was extracted from muscle samples using the cetyl trimethylammonium bromide (CTAB) method and a DNeasy blood & tissue kit (QIAGEN). DNA quality and quantity were checked using pulsed field gel electrophoresis and a Qubit Fluorometer (Thermo Scientific).

**Libraries preparation and sequencing**

***Whole Genome Shotgun Sequencing***

Four WGS libraries were prepared for sequencing: one short insert size library (350 bp) and three mate-pair large insert size libraries (2 kb, 5 kb, and 10 kb). Libraries were constructed using an MGI Easy FS DNA Library Prep Set kit (MGI, China). Paired-end reads (100 bp) and mate-pair reads (50 bp) were obtained from the BGISEQ-500 platform.

***10X Genomics sequencing***

To prepare the Chromium library, 1 ng of high quality DNA was denatured, spiked into reaction mix, and mixed with gel beads and emulsification oil to generate droplets within a Chromium Genome chip. Then, the rest of the steps were completed following the standard protocols for performing PCR. After PCR, the standard circularisation step for BGISEQ-500 was carried out, and DNA nanoballs (DNBs) were prepared [44]. Paired-end reads with a length of 150 bp were generated on the BGISEQ-500 platform.

***Oxford Nanopore Technologies***

DNA for long-read sequencing was isolated from the muscle tissues of our samples. Using 5 flow cells of the ONT chemistry for the GridION X5 sequencer

387   following manufacturer's protocols, we generated 39.61 Gbp of raw genome sequencing

388   data.

389

### *Hi-C* library and sequencing

391   The Hi-C library was prepared following the standard *in situ* Hi-C [45] protocol for

392   muscle samples, using *DpnII* (NEB, Ipswich, America) as the restriction enzyme. After

393   that, a standard circularisation step was carried out, followed by DNA nanoballs (DNB)

394   preparation following the standard protocols of the BGISEQ-500 sequencing platform as

395   previously described [44]. Paired-end reads with a length of 100 bp were generated on

396   the BGISEQ-500 platform.

397

### Genome assembly

399   For the genome assembly of *Chrysomallon squamiferum*, Canu (v1.7) was first used to

400   perform corrections of ONT reads with the parameters "correctedErrorRate=0.105

401   corMinCoverage=0   minReadLength=1000   minOverlapLength=800".   Then,   wtdbg

402   (v1.2.8) was used to assemble the genome with the parameters "--tidy-reads 3000 -k 0 -p

403   21 -S 4 --rescue-low-cov-edges" using corrected reads generated by Canu. Next, we

404   made use of the sequencing reads from the 10X genomic library to carry out genome

405   polishing using Pilon (version 1.22) with its default parameters. Quality control of Hi-C

406   sequencing reads was first performed using the HiC-Pro pipeline [46] with the parameters

407   "[BOWTIE2_GLOBAL_OPTIONS = --very-sensitive -L 30 --score-min L,-0.6,-0.2 --

408   end-to-end –reorder;BOWTIE2_LOCAL_OPTIONS = --very-sensitive -L 20 --score-min

409   L,-0.6,-0.2 --end-to-end –reorder; IGATION_SITE = GATC; MIN_FRAG_SIZE = 100;

410   MAX_FRAG_SIZE = 100000; MIN_INSERT_SIZE = 50; MAX_INSERT_SIZE =

411   1500]". In total, 23,646,810 pairs of valid reads were obtained. Next, the valid Hi-C data

412   was used to anchor the nanopore contigs onto chromosomes separately by applying the

413   3D-DNA [47] pipeline. The contact maps were then generated by the Juicer pipeline[48],

414   and the boundaries for each chromosome were manually rectified by visualising the

415   inter.hic file in Juicebox [49]. 16 chromosomes were identified by combining the linkage

416   information from the agp file.

417    For the genome assembly of *Gigantopelta aegis*, we only have WGS sequencing reads

418    because of limited DNA and tissue samples. Platanus (v1.2.4)[12] was used to perform

419    genome assembly with WGS clean data with the parameters "assemble –k 29 –u 0.2,

420    scaffold -l 3 -u 0.2 -v 32 -s 32 and gap_close –s 34 –k 32 –d 5000". BUSCO (v2) were

421    used to evaluate genome assemblies with the metazoan_odb9 database.

422

423    **Genome annotation**

424    ***Repeat annotation***

425    Homolog-based and *de novo* prediction methods were used to detect repeat contents. In

426    particular, RepeatMasker (v4.0.5) [50] and RepeatProteinMasker (v 4.0.5) were used to

427    detect transposable elements against the Repbase database[51] at the nuclear and protein

428    levels, respectively. RepeatMasker was used again to detect species-specific transposable

429    elements against databases generated by RepeatModeler (v1.0.8) and LTR-FINDER

430    (v1.0.6)[52]. Moreover, Tandem Repeat Finder (v4.0.7)[53] was utilised to predict

431    tandem repeats.

432

433    ***Gene annotation***

434    We combined homology-based and *de novo* evidence to predict protein-coding genes in

435    two genomes. For the homology-based method, we used six relative gene sets of *Aplysia*

436    *californica*, *Bathymodiolus platifrons*, *Biomphalaria glabrata*, *Lottiu gigantea*, *Modiolus*

437    *philippinarum*, and *Pomacea canaliculata*. First, these homologous protein sequences

438    were aligned onto each assembled genome using TBLASTN (RRID:SCR 011822), with

439    an *E*-value cut-off of $1 \times 10^{-5}$, and the alignment hits were linked to candidate gene loci

440    by GenBlastA [54]. Second, we extracted genomic sequences of candidate gene regions,

441    including 2 kb flanking sequences, and then used GeneWise (v2.2.0)[55] to determine

442    gene models.

443

444    In the *de novo* method, we used Augustus (Augustus, RRID:SCR 008417)[56] to predict

445    the gene models on repeat-masked genome sequences. We selected high-quality genes

446    with intact open reading frames (ORFs) and the highest GeneWise [55] score from a

447    homology-based gene set to train Augustus with default parameters before prediction.

448  Gene models with incomplete ORFs and small genes with protein-coding lengths less

449  than 150 bp were filtered out. Finally, a BLASTP (BLASTP, RRID:SCR 001010) search

450  of predicted genes was performed against the Swiss-Prot database (UniProt, RRID:SCR

451  002380) [57]. Genes with matches to Swiss-Prot proteins containing any one of the

452  following keywords were filtered: transpose, transposon, retrotransposon, retrovirus,

453  retrotransposon, reverse transcriptase, transposase, and retroviral. Finally, the results of

454  the homology- and *de novo*-based gene sets were merged using GLEAN to yield a

455  nonredundant reference gene set.

456

457  ***Gene function annotation***

458  We annotated the protein-coding genes by searching against the following public

459  databases: Swiss-Prot[58], the Kyoto Encyclopedia of Genes and Genomes (KEGG)[59],

460  InterPro[60], and TrEMBL[58].

461

462  **Phylogenetic tree reconstruction and divergence time estimation**

463  The TreeFam tool [61] was used to identify gene families as follows: first, all the protein

464  sequences from selected 10 representative species (*Aplysia californica*, *Octopus*

465  *bimaculoides*, *Biomphalaria glabrata*, *Crassostrea gigas*, *Lottia gigantea*, *Pomacea*

466  *canaliculata*, *Pinctada fucata*, *Chrysomallon squamiferum*, *Gigantopelta aegis,* and

467  *Helobdella robusta*) were compared using blastp with the *E*-value threshold set as 1e-7.

468  Then, alignment segments of each protein pair were concatenated using the in-house

469  software Solar. H-scores were computed based on Bit-scores and were used to evaluate

470  the similarity among proteins. Finally, gene families were obtained by clustering

471  homologous gene sequences using Hcluster sg (v0.5.0).

472

473  We obtained 406 one-to-one single-copy orthology gene families based on gene family

474  classification. Then, these gene families were extracted and aligned using guidance from

475  amino-acid alignments created using the default parameters of the MUSCLE [62]

476  programme. All sequence alignments were then concatenated to construct 1 super-matrix

477  and then a phylogenetic tree was constructed under a GTR+gamma model for nucleotide

478  sequences using ML and Bayesian methods. The same set of codon sequences were used

479  for phylogenetic tree construction and estimation of divergence time. The PAML

480  mcmctree programme [63, 64] was used to determine divergence times with the

481  approximate likelihood calculation method,  and the correlated molecular clock and REV

482  substitution model. The concatenated CDS of one-to-one orthologous genes and the

483  phylogenomics topology were used as inputs. We used five calibration time points based

484  on fossil records: *A. californica - C. gigas* (~516.3 - 558.3 million years ago (Mya)), *A.*

485  *californica - P. canaliculata* (~310 – 496 Mya), *A. californica - Octopus bimaculoides*

486  (~551 – 628 Mya), *C. gigas - H. robusta* (~585 – 790 Mya), and *C. gigas– P. fucata*

487  (394 Mya) (http://www.timetree.org), were used as constraints in the MCMCTree

488  estimation.

489

490  **Expansion and contraction of gene families**

491  We used CAFE (Computational Analysis of gene Family Evolution) v2.1[19] to analyse

492  gene family expansion and contraction under the maximum likelihood framework. The

493  gene family results from the TreeFam pipeline and the estimated divergence time

494  between species were used as inputs. We used the parameters "-p 0.01, -r 10000, -s" to

495  search for the birth and death parameter ($\lambda$) of gene families, calculated the probability of

496  each gene family with observed sizes using 10,000 Monte Carlo random samplings, and

497  reported birth and death parameters in gene families with probabilities less than 0.01.

498

499

**Figure legends**

501  **Figure 1. Genome characteristics of *C. squamiferum* and *G. aegis*. a)** Photos of two

502  species. Left: *C. squamiferum*; right: *G. aegis.* **b**) Heat map of chromatin interaction

503  relationships at a 125 kb resolution of 16 chromosomes. **c**) Genome sizes and

504  transposable elements in *C. squamiferum, G. aegis,* and two representative freshwater

505  snail genomes. **d**) Distribution of repeat sub-types of four species.

506

507  **Figure 2. Phylogenetic tree, estimated $N_e$, and evolution of single copy orthologous**

508  **genes of deep-sea snails**. **a**) Phylogenetic tree of ten representative molluscs. Expanded

509  and contracted gene families were identified using CAFE. Divergence time was estimated

510  using mcmctree. Species names in red represent two deep-sea snails. The timescale refers

511  to the TimeTree database. **b**) Estimated demographic histories of two deep-sea snails.

512  The generation time set to "3" refers to the land snail [65]. The $\mu$ values are calculated

513  in **Table S15**. **c**). Box plot of $K$a/$K$s values for five snails.

514

515  **Figure 3. Expansion of nervous system-related genes a**) Phylogentic tree of *BTBD6*

516  genes in the examined species. The grey ellipses mark different clusters of genes. **b**)

517  Expansion pattern of *BTBD6* genes in two deep-sea snails. Grey lines represent scaffold

518  sequences. Coloured rectangles represent *BTBD6* genes. Symbols "//" represent other

519  genes along the scaffolds. The blue numbers: "1" represent only one gene between the

520  tandem duplicated genes. **c**) Expansion of *HTR4* genes. The species legend in the middle

521  was used for **a** and **c**. Gene trees of **a** and **c** were constructed using MUSCLE (v3.8.31)

522  and FastTree (v2.1.10). **d**) Sketch map of the large unganglionated nervous system of *C.*

523  *squamiferum.* The prunosus represents the nervous system. The right amplifying

524  represents one example of neurotransmitter release.

525

526  **Figure 4**. **Expansion of immune, metabolism, DNA stability, and antioxidation genes**.

527  **a**) Gene numbers of four defence-related genes (*DMBT1*, *GAFT*, *Hsp90,* and *Txn1*), three

528  metabolism-related genes (*OGDHE1*, *OGDHE2,* and *IDH*), and the *SSB* gene. **b**) TCA

529  cycle signal pathway. The brown ellipses represent important enzymes and the expansion

530    of these genes (*OGDHE1*, *OGDHE2,* and *IDH*). **c**) Expansion of the catalase (CAT) gene

531    in selected species.

532

533    **Table 1. Genome assembly and annotation of *Chrysomallon squamiferum* and**

534    ***Gigantopelta aegis*.**

| Species | *Chrysomallon squamiferum* | *Gigantopelta aegis* |
|---|---|---|
| Genome size | 455.36 Mb | 1.29 GB |
| Scaffold N50 | 20.7 Mb | 120.96 kb |
| Contig N50 | 541.32 kb | 6.96 kb |
| Number of genes | 28,781 | 25,601 |
| Repeat content | 30.56% | 64.17% |
| GC content | 34.48% | 37.45% |
| Complete BUSCO | 94.80% | 88.40% |

535

536    **Data and code availability**

537    The genome assemblies of these two genomes have been deposited in GenBank under the

538    accession number CNP0000854. The raw sequencing reads were also uploaded to the

539    SRA database under accession number  CNP0000854.

540    **Additional Files**

541    Additional File 1: Supplementary Figures and Tables.docx

542

543    **Author contributions**

544    Z.S., S.L., G.F., and X.L. conceived and managed this project and amended the

545    manuscript. X.Z., Y.Z., L.M., and I.S. performed the evolutionary analysis and wrote the

546    manuscript. L.M., J.C., and Y.S. performed genome assembly and annotation. J.B., S.L.,

547    X.F., C.W., Z.S., H.L., N.L., and L.W. were responsible for sample collection, DNA

548    extraction, and library construction.

549

550    **Acknowledgments**

## Supplemental information

Supplemental Information can be found online.

## Declaration of interests

The authors declare that they have no competing interests.

## References

1. Corliss JB, Dymond J, Gordon LI, Edmond JM, Von Herzen RP, Ballard RD, et al. Submarine Thermal Springs on the Galápagos Rift. Science. 1979;203 4385:1073-83.
2. VAN DOVER CL. The ecology of deep-sea hydrothermal vents. Princeton: Princeton University Press; 2000.
3. Sun J, Zhang Y, Xu T, Zhang Y, Mu H, Zhang Y, et al. Adaptation to deep-sea chemosynthetic environments as revealed by mussel genomes. Nature ecology & evolution. 2017;1 5:0121.
4. Parkhaev PY. The Cambrian 'basement' of gastropod evolution. Geological Society, London, Special Publications. 2007;286 1:415-21.
5. Sasaki T, Warén A, Kano Y, Okutani T and Fujikura K. Gastropods from Recent Hot Vents and Cold Seeps: Systematics, Diversity and Life Strategies. 2010.
6. Chong C, Linse K, Copley TJ and Rogers DA. The 'scaly-foot gastropod': a new genus and species of hydrothermal vent-endemic gastropod (Neomphalina: Peltospiridae) from the Indian Ocean. 2015;81 3:322-34.
7. Chen C, Copley JT, Linse K and Rogers AD. Low connectivity between 'scaly-foot gastropod' (Mollusca: Peltospiridae) populations at hydrothermal vents on the Southwest Indian Ridge and the Central Indian Ridge. Organisms Diversity & Evolution. 2015;15 4:663-70.
8. Sigwart JD, Chen C, Thomas EA, Allcock AL, Bohm M and Seddon M. Red Listing can protect deep-sea biodiversity. Nature Ecology and Evolution. 2019;3 8:1134-.
9. Chen C, Uematsu K, Linse K and Sigwart JD. By more ways than one: Rapid convergence at hydrothermal vents shown by 3D anatomical reconstruction of Gigantopelta (Mollusca: Neomphalina). BMC Evolutionary Biology. 2017;17 1:62.

589　10.　Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH and Phillippy AM. Canu:
590　　　　scalable and accurate long-read assembly via adaptive k-mer weighting and repeat
591　　　　separation. Genome research. 2017;27 5:722-36.
592　11.　Ruan J and Li H. Fast and accurate long-read assembly with wtdbg2. Nature
593　　　　Methods. 2019:1-4.
594　12.　Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al.
595　　　　Efficient de novo assembly of highly heterozygous genomes from whole-genome
596　　　　shotgun short reads. Genome Res. 2014;24 8:1384-95.
597　　　　doi:10.1101/gr.170720.113.
598　13.　Adema CM, Hillier LW, Jones CS, Loker ES, Knight M, Minx P, et al. Whole
599　　　　genome analysis of a schistosomiasis-transmitting freshwater snail. Nature
600　　　　communications. 2017;8:15451.
601　14.　Liu C, Zhang Y, Ren Y, Wang H, Li S, Jiang F, et al. The genome of the golden
602　　　　apple snail Pomacea canaliculata provides insight into stress tolerance and
603　　　　invasive adaptation. GigaScience. 2018;7 9:giy101.
604　15.　Biemont C. Genome size evolution: within-species variation in genome size.
605　　　　Nature Publishing Group, 2008.
606　16.　Dehal P and Boore JL. Two rounds of whole genome duplication in the ancestral
607　　　　vertebrate. PLoS biology. 2005;3 10.
608　17.　Barnosky AD, Matzke N, Tomiya S, Wogan GO, Swartz B, Quental TB, et al. Has
609　　　　the Earth's sixth mass extinction already arrived? Nature. 2011;471 7336:51-7.
610　18.　Li H and Durbin R. Inference of human population history from individual whole-
611　　　　genome sequences. Nature. 2011;475 7357:493-6.
612　19.　De Bie T, Cristianini N, Demuth JP and Hahn MW. CAFE: a computational tool
613　　　　for the study of gene family evolution. Bioinformatics. 2006;22 10:1269-71.
614　20.　Sobieszczuk DF, Poliakov A, Xu Q and Wilkinson DG. A feedback loop mediated
615　　　　by degradation of an inhibitor is required to initiate neuronal differentiation.
616　　　　Genes & development. 2010;24 2:206-18.
617　21.　Conductier G, Dusticier N, Lucas G, Côté F, Debonnel G, Daszuta A, et al.
618　　　　Adaptive changes in serotonin neurons of the raphe nuclei in 5-HT4 receptor
619　　　　knock-out mouse. European Journal of Neuroscience. 2006;24 4:1053-62.
620　22.　Goffredi SK, Warén A, Orphan VJ, Van Dover CL and Vrijenhoek RC. Novel
621　　　　forms of structural integration between microbes and a hydrothermal vent
622　　　　gastropod from the Indian Ocean. Appl Environ Microbiol. 2004;70 5:3082-90.
623　23.　Wolff T. Composition and endemism of the deep-sea hydrothermal vent fauna.
624　　　　CBM-Cahiers de Biologie Marine. 2005;46 2:97-104.
625　24.　Muri J, Heer S, Matsushita M, Pohlmeier L, Tortola L, Fuhrer T, et al. The
626　　　　thioredoxin-1 system is essential for fueling DNA synthesis during T-cell
627　　　　metabolic reprogramming and proliferation. Nature communications. 2018;9
628　　　　1:1851.
629　25.　Kato N, Dasgupta R, Smartt C and Christensen B. Glucosamine: fructose-6-
630　　　　phosphate aminotransferase: gene characterization, chitin biosynthesis and
631　　　　peritrophic matrix formation in Aedes aegypti. Insect molecular biology. 2002;11
632　　　　3:207-16.

633 26. Lagorce A, Le Berre-Anton V, Aguilar-Uscanga B, Martin-Yken H,
634     Dagkessamanskaia A and François J. Involvement of GFA1, which encodes
635     glutamine–fructose-6-phosphate amidotransferase, in the activation of the chitin
636     synthesis pathway in response to cell-wall defects in Saccharomyces cerevisiae.
637     European journal of biochemistry. 2002;269 6:1697-707.
638 27. Csermely P, Schnaider T, So C, Prohászka Z and Nardai G. The 90-kDa molecular
639     chaperone family: structure, function, and clinical applications. A comprehensive
640     review. Pharmacology & therapeutics. 1998;79 2:129-68.
641 28. Marceau AH. Functions of single-strand DNA-binding proteins in DNA
642     replication, recombination, and repair. Single-Stranded DNA Binding Proteins.
643     Springer; 2012. p. 1-21.
644 29. Nazıroğlu M. Molecular role of catalase on oxidative stress-induced Ca2+
645     signaling and TRP cation channel activation in nervous system. Journal of
646     Receptors and Signal Transduction. 2012;32 3:134-41.
647 30. Hohenester E, Sasaki T and Timpl R. Crystal structure of a scavenger receptor
648     cysteine-rich domain sheds light on an ancient superfamily. Nature Structural &
649     Molecular Biology. 1999;6 3:228.
650 31. Ligtenberg AJ, Karlsson NG and Veerman EC. Deleted in malignant brain tumors-
651     1 protein (DMBT1): a pattern recognition receptor with multiple binding sites.
652     International journal of molecular sciences. 2010;11 12:5212-33.
653 32. Aguilera F, McDougall C and Degnan BM. Co-option and de novo gene evolution
654     underlie molluscan shell diversity. Molecular biology and evolution. 2017;34
655     4:779-92.
656 33. Mann K, Edsinger-Gonzales E and Mann M. In-depth proteomic analysis of a
657     mollusc shell: acid-soluble and acid-insoluble matrix of the limpet Lottia
658     gigantea. Proteome science. 2012;10 1:28.
659 34. Benoist CO, Mathis DJ, Kanter MR, Williams II VE and McDevitt HO. Regions
660     of allelic hypervariability in the murine Aα immune response gene. Cell. 1983;34
661     1:169-77.
662 35. Lindberg DR PW, Haszprunar G. The Mollusca: Relationships and Patterns
663     fromTheir First Half-Billion Years. Oxford: Oxford University Press; 2004.
664 36. Wanninger A and Wollesen T. The evolution of molluscs. Biological Reviews.
665     2019;94 1:102-15.
666 37. Vinther J. A molecular palaeobiological perspective on aculiferan evolution.
667     Journal of natural history. 2014;48 45-48:2805-23.
668 38. Vinther J. The origins of molluscs. Palaeontology. 2015;58 1:19-34.
669 39. Lee H, Chen W, Puillandre N, Aznarcormano L, Tsai M and Samadi S.
670     Incorporation of deep-sea and small-sized species provides new insights into
671     gastropods phylogeny. Molecular Phylogenetics and Evolution. 2019;135:136-47.
672 40. Kazazian HH. Mobile elements: drivers of genome evolution. science. 2004;303
673     5664:1626-32.
674 41. Jacques P-E, Jeyakani J and Bourque G. The majority of primate-specific
675     regulatory sequences are derived from transposable elements. PLoS genetics.
676     2013;9 5.
677 42. Sundaram V, Cheng Y, Ma Z, Li D, Xing X, Edge P, et al. Widespread

678  contribution of transposable elements to the innovation of gene regulatory
679  networks. Genome research. 2014;24 12:1963-76.

680  43.  Janoušek V, Laukaitis CM, Yanchukov A and Karn RC. The role of
681  retrotransposons in gene family expansions in the human and mouse genomes.
682  Genome biology and evolution. 2016;8 9:2632-50.

683  44.  Drmanac R, Sparks AB, Callow MJ, Halpern AL, Burns NL, Kermani BG, et al.
684  Human genome sequencing using unchained base reads on self-assembling DNA
685  nanoarrays. Science. 2010;327 5961:78-81. doi:10.1126/science.1181498.

686  45.  Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, et
687  al. A 3D map of the human genome at kilobase resolution reveals principles of
688  chromatin looping. Cell. 2014;159 7:1665-80. doi:10.1016/j.cell.2014.11.021.

689  46.  Servant N, Varoquaux N, Lajoie BR, Viara E, Chen CJ, Vert JP, et al. HiC-Pro: an
690  optimized and flexible pipeline for Hi-C data processing. Genome Biol.
691  2015;16:259. doi:10.1186/s13059-015-0831-x.

692  47.  Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al. De
693  novo assembly of the Aedes aegypti genome using Hi-C yields chromosome-
694  length scaffolds. Science. 2017;356 6333:92-5. doi:10.1126/science.aal3327.

695  48.  Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, et al.
696  Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C
697  Experiments. Cell Syst. 2016;3 1:95-8. doi:10.1016/j.cels.2016.07.002.

698  49.  Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, Lander ES, et al.
699  Juicebox Provides a Visualization System for Hi-C Contact Maps with Unlimited
700  Zoom. Cell Syst. 2016;3 1:99-101. doi:10.1016/j.cels.2015.07.012.

701  50.  Tarailo-Graovac M and Chen N. Using RepeatMasker to identify repetitive
702  elements in genomic sequences. Current Protocols in Bioinformatics. 2009:4.10.
703  1-4.. 4.

704  51.  Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O and Walichiewicz J.
705  Repbase Update, a database of eukaryotic repetitive elements. Cytogenetic and
706  genome research. 2005;110 1-4:462-7.

707  52.  Xu Z and Wang H. LTR_FINDER: an efficient tool for the prediction of full-
708  length LTR retrotransposons. Nucleic acids research. 2007;35 suppl 2:W265-W8.

709  53.  Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic
710  acids research. 1999;27 2:573.

711  54.  She R, Chu JS, Wang K, Pei J and Chen N. GenBlastA: enabling BLAST to
712  identify homologous gene sequences. Genome Res. 2009;19 1:143-9.
713  doi:10.1101/gr.082081.108.

714  55.  Birney E, Clamp M and Durbin R. GeneWise and Genomewise. Genome Res.
715  2004;14 5:988-95. doi:10.1101/gr.1865504.

716  56.  Keller O, Kollmar M, Stanke M and Waack S. A novel hybrid gene prediction
717  method employing protein multiple sequence alignments. Bioinformatics. 2011;27
718  6:757-63. doi:10.1093/bioinformatics/btr010.

719  57.  UniProt C. UniProt: a worldwide hub of protein knowledge. Nucleic Acids Res.
720  2019;47 D1:D506-D15. doi:10.1093/nar/gky1049.

721  58.  Boeckmann B, Bairoch A, Apweiler R, Blatter M-C, Estreicher A, Gasteiger E, et
722  al. The SWISS-PROT protein knowledgebase and its supplement TrEMBL in
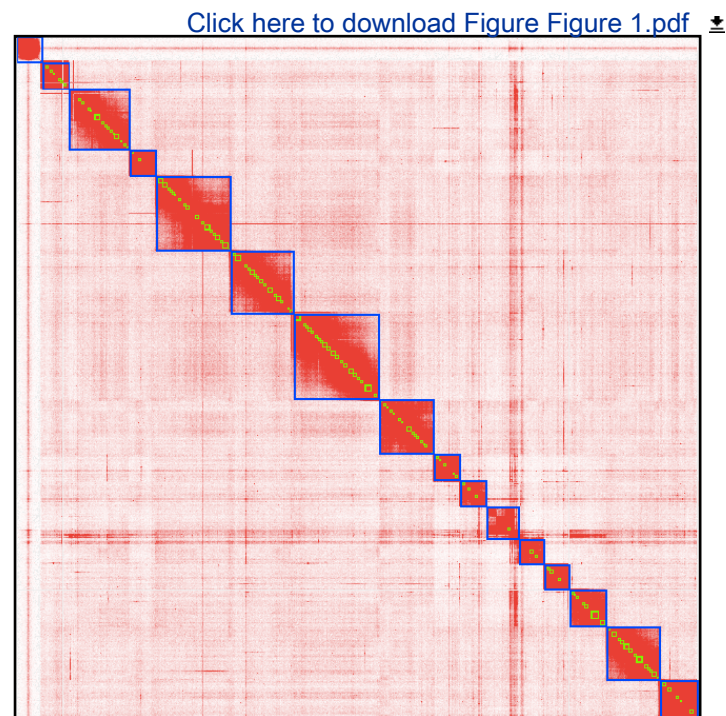723  2003. Nucleic acids research. 2003;31 1:365-70.

724  59.  Kanehisa M and Goto S. KEGG: kyoto encyclopedia of genes and genomes.
725      Nucleic acids research. 2000;28 1:27-30.
726  60.  Apweiler R, Attwood TK, Bairoch A, Bateman A, Birney E, Biswas M, et al. The
727      InterPro database, an integrated documentation resource for protein families,
728      domains and functional sites. Nucleic acids research. 2001;29 1:37-40.
729  61.  Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, et al. TreeFam: a
730      curated database of phylogenetic trees of animal gene families. Nucleic Acids
731      Res. 2006;34 Database issue:D572-80. doi:10.1093/nar/gkj118.
732  62.  Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high
733      throughput. Nucleic Acids Res. 2004;32 5:1792-7. doi:10.1093/nar/gkh340.
734  63.  Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. Mol Biol Evol.
735      2007;24 8:1586-91. doi:10.1093/molbev/msm088.
736  64.  Yang Z and Rannala B. Bayesian estimation of species divergence times under a
737      molecular clock using multiple fossil calibrations with soft bounds. Mol Biol
738      Evol. 2006;23 1:212-26. doi:10.1093/molbev/msj024.
739  65.  Schilthuizen M. Rapid, habitat-related evolution of land snail colour morphs on
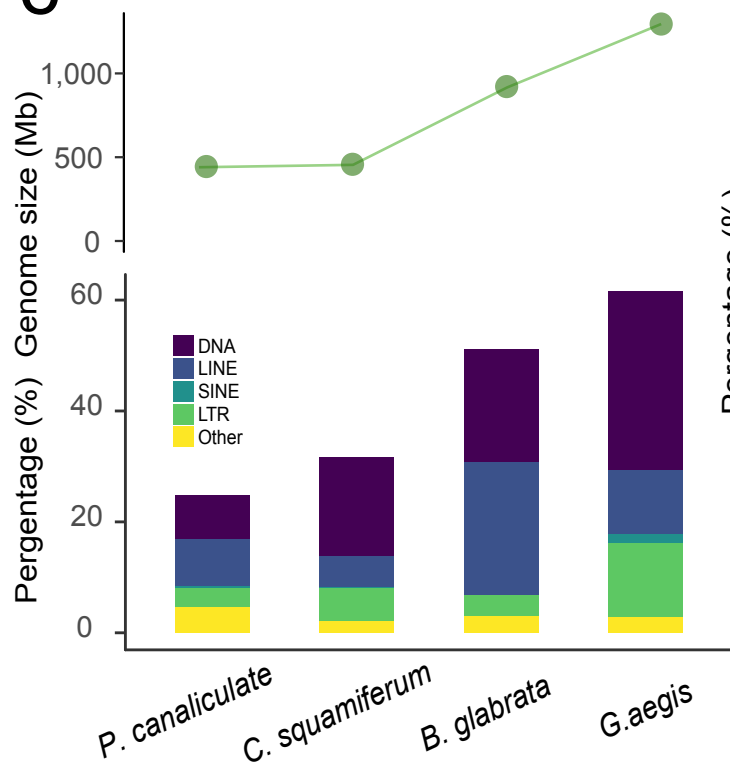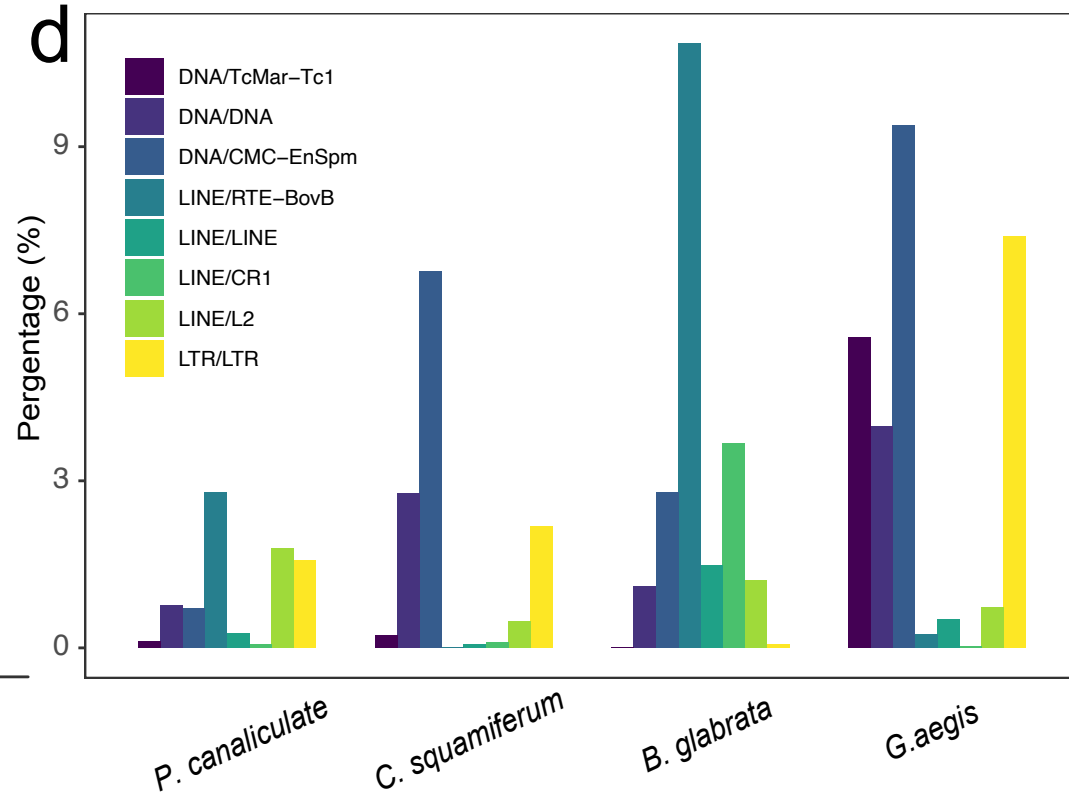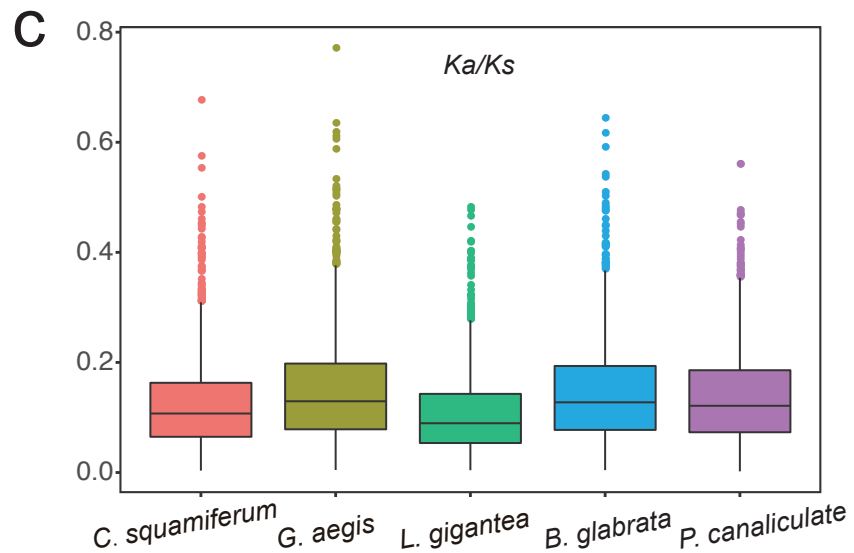740      reclaimed land. Heredity. 2013;110 3:247-52.
741

Figure 1

Figure 2

a

+1677 / -19359    *H. robusta*
+873 / -18124     *O. bimaculoides*
+2713 / -3182     *C. gigas*
+3777 / -3480     *P. fucata*
+2160 / -14397    *L. gigantea*
+2301 / -6786     *G. aegis*
66.3(42.4-100.0)
+3088 / -3857     *C. squamiferum*
536.6(525.1-543.3)
+1585 / -4825     *P. canaliculate*
+2383 / -1739     *B. glabrata*
+948 / -3786      *A. californica*

Divergence time
Expanded gene families
Contracted gene families

Neo-Proterozoic | Paleozoic | Mesozoic | Cenozoic
Proterozoic | Phanerozoic

688  600        400        200        0   MYA

b

C. squamiferum (μ=2.81 x10⁻⁹)
G. aegis (μ=4.19 x10⁻⁹)

Effective population size (x10⁴)

Time since present (years) (**g** = 3)

c

*Ka/Ks*

C. squamiferum, G. aegis, L. gigantea, B. glabrata, P. canaliculate

Figure 3

Click here to download Figure Figure 3.pdf ⬇



**a**

3

2

1

0.3

- *C. squamiferum*
- *G. aegis*
- *P. canaliculate*
- *B. glabrata*
- *O. bimaculoides*
- *L. gigantea*

**c**

1

2

0.3

**d**

BTBD6

HTR4

**b**

*C. squamiferum*

21542  22602  22601  22599  16079  16080  16082  16084  16096  16101  24917  02464  02463  02462  02461  07957  07956

1                    1      1

chr16: 0.76-10.35 Mb

*G. aegis*

10004001  10004002

scaffold1711 (1.58-9.00 kb)

10008840  10008842

1

scaffold31399 (70.10-102.49 kb)

10012143  10012145  10012146

1

scaffold48401 (128.41-159.78 kb)

Figure 4

Click here to access/download
**Supplementary Material**
Supplementary Figures and Tables.docx