

## Genome of the ramshorn snail *Biomphalaria straminea* - an intermediate vector of schistosomiasis --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-21-00243	
<b>Full Title:</b>	Genome of the ramshorn snail <i>Biomphalaria straminea</i> - an intermediate vector of schistosomiasis	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	Hong Kong Research Grant Council (Collaborative Research Fund (C4015-20EF), General Research Fund (14100919))	Dr. Jerome Hui
<b>Abstract:</b>	<p><b>Background</b> Schistosomiasis or bilharzia is a parasitic disease caused by trematode flatworms of the genus <i>Schistosoma</i>. Infection of <i>Schistosoma mansoni</i> in humans results when cercariae emerge into water from freshwater snails in the genus <i>Biomphalaria</i>, and seek out and penetrate human skin. The snail <i>Biomphalaria straminea</i> was native to South America and is now also present in Central America and China, and represents a potential reservoir for spreading schistosomiasis. To date, genomic information for the genus is restricted to the neotropical species <i>Biomphalaria glabrata</i>. This hinders understanding of the biology and management of other schistosomiasis vectors, such as <i>B. straminea</i>.</p> <p><b>Findings</b> Using a combination of Illumina short-read, 10X Genomics linked-read, and Hi-C sequencing data, our 1.005 Gbp <i>B. straminea</i> genome assembly is of high contiguity, with a scaffold N50 of 25.3 Mbp. Developmental homeobox genes, hormonal genes, and stress-response genes were identified, and repeat content was annotated (40.68% of genomic content). Comparisons with other mollusc genomes revealed syntenic conservation, patterns of homeobox gene linkage indicative of evolutionary changes to gene clusters, expansion of heat shock protein genes, and the presence of sesquiterpenoid and cholesterol metabolic pathway genes in certain mollusc lineages.</p> <p><b>Conclusion</b> This study provides the first genome assembly for the snail <i>B. straminea</i> and offers an unprecedented opportunity to address a variety of biology related to schistosomiasis, as well as evolutionary and genomics questions related to molluscs more widely.</p>	
<b>Corresponding Author:</b>	Jerome Hui  HONG KONG	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Wenyan Nong	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Wenyan Nong	
	Yifei Yu	
	Madeleine E. Aase-Remedios	
	Yichun Xie	
	Wai Lok So	
	Yiqian Li	

	Cheuk Fung Wong
	Toby Baril
	Sean TS Law
	Sheung Yee Lai
	Jasmine Haimovitz
	Thomas Swale
	Shan-shan Chen
	Zhen-peng Kai
	Xi Sun
	Zhongdao Wu
	Alexander Hayward
	David Ferrier
	Jerome Hui
<b>Order of Authors Secondary Information:</b>	
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>  Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a> . Information essential to interpreting the data presented should be made available in the figure legends.  Have you included all the information requested in your manuscript?	Yes
<b>Resources</b>  A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.	Yes

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

1 **Genome of the ramshorn snail *Biomphalaria straminea* - an intermediate vector of**  
2 **schistosomiasis**

3 Wenyan Nong<sup>1,^</sup>, Yifei Yu<sup>1,^</sup>, Madeleine E. Aase-Remedios<sup>3,^</sup>, Yichun Xie<sup>1,^</sup>, Wai Lok So<sup>1,^</sup>,  
4 Yiqian Li<sup>1,^</sup>, Cheuk Fung Wong<sup>1,^</sup>, Toby Baril<sup>2,^</sup>, Sean T. S. Law<sup>1,^</sup>, Sheung Yee Lai<sup>1</sup>, Jasmine  
5 Haimovitz<sup>4</sup>, Thomas Swale<sup>4</sup>, Shan-shan Chen<sup>5</sup>, Zhen-peng Kai<sup>6</sup>, Sun Xi<sup>7</sup>, Zhongdao Wu<sup>7</sup>,  
6 Alexander Hayward<sup>2,\*</sup>, David E.K. Ferrier<sup>3,\*</sup>, Jerome H.L. Hui<sup>1,\*</sup>

7 1. School of Life Science, Simon F.S. Li Marine Science Laboratory, State Key Laboratory of  
8 Agrobiotechnology, The Chinese University of Hong Kong, Hong Kong, China

9 2. University of Exeter, United Kingdom

10 3. The Scottish Oceans Institute, Gatty Marine Laboratory, School of Biology, University of St.  
11 Andrews, United Kingdom

12 4. Dovetail Genomics, United States of America

13 5. Institute of Agro-food Standard and Testing Technology, Shanghai Academy of Agricultural  
14 Sciences, Shanghai, China

15 6. School of Chemical and Environmental Engineering, Shanghai Institute of Technology,  
16 Shanghai, China

17 7. Sun Yat-sen University, Guangdong, China

18 ^ = contributed equally;

19 \*=correspondence:

20 alex.hayward@exeter.ac.uk; dekf@st-andrews.ac.uk; jeromehui@cuhk.edu.hk

1 **Abstract**

2 **Background**

3 Schistosomiasis or bilharzia is a parasitic disease caused by trematode flatworms of the  
4 genus *Schistosoma*. Infection of *Schistosoma mansoni* in humans results when cercariae emerge  
5 into water from freshwater snails in the genus *Biomphalaria*, and seek out and penetrate human  
6 skin. The snail *Biomphalaria straminea* was native to South America and is now also present in  
7 Central America and China, and represents a potential reservoir for spreading schistosomiasis. To  
8 date, genomic information for the genus is restricted to the neotropical species *Biomphalaria*  
9 *glabrata*. This hinders understanding of the biology and management of other schistosomiasis  
10 vectors, such as *B. straminea*.

11 **Findings**

12 Using a combination of Illumina short- read, 10X Genomics linked- read, and Hi- C sequencing  
13 data, our 1.005 Gbp *B. straminea* genome assembly is of high contiguity, with a scaffold N50 of  
14 25.3 Mbp. Developmental homeobox genes, hormonal genes, and stress-response genes were  
15 identified, and repeat content was annotated (40.68% of genomic content). Comparisons with other  
16 mollusc genomes revealed syntenic conservation, patterns of homeobox gene linkage indicative of  
17 evolutionary changes to gene clusters, expansion of heat shock protein genes, and the presence of  
18 sesquiterpenoid and cholesterol metabolic pathway genes in certain mollusc lineages.

19 **Conclusion**

20 This study provides the first genome assembly for the snail *B. straminea* and offers an  
21 unprecedented opportunity to address a variety of biology related to schistosomiasis, as well as  
22 evolutionary and genomics questions related to molluscs more widely.

## 1 **Background**

2           With over 240 million people worldwide estimated to require preventive treatment, the  
3 World Health Organisation considers schistosomiasis to be the second most prevalent parasitic  
4 disease after malaria (<https://www.who.int/health-topics/schistosomiasis>). As such,  
5 schistosomiasis is a global health problem that causes considerable economic and social burdens.

6           Infection by *Schistosoma mansoni* in humans results when cercariae emerge into the water  
7 from their freshwater snail intermediate hosts in the genus *Biomphalaria*, and seek out and  
8 penetrate submerged body parts through the skin. Once inside the human body, adult worms  
9 migrate to the mesenteric venules of the bowel or rectum and lay thousands of eggs that circulate  
10 to the liver and leave the body via faeces. Miracidia larvae hatch from eggs that reach contaminated  
11 water, then seek out and penetrate a new snail intermediate host. Following this, sporocysts  
12 develop in the infected snails, and subsequently further free-living cercariae emerge from the snail  
13 into the water, completing the parasitic life cycle. Among the 34 described species  
14 of *Biomphalaria* snails, 18 species (including *B. straminea*) have been demonstrated to be infected  
15 by *S. mansoni*. Different geographical locations are dominated by different species  
16 of *Biomphalaria*.

17           The native range of *Biomphalaria* snails is South America and Africa (Campbell et al  
18 2000). However, several species have been introduced to other areas, presenting a risk of  
19 schistosomiasis infection. The occurrence of *B. straminea* in Asia was first reported at Lam Tsuen  
20 valley in Hong Kong during the 1970s (Meier-Brook 1974; Figure 1A), presumably having  
21 somehow spread from its native range in South America into Central America and southern China  
22 (Yang et al., 2018). *B. straminea* have since been identified at a number of locations in Hong Kong

1 and Guangdong Province (Attwood et al 2015; Dudgeon and Yipp 1983; Meier-Brook 1974;  
2 Woodruff et al 1985; Zeng et al 2017). While *S. mansoni* is not yet endemic in either Hong Kong  
3 or mainland China, cases of schistosomiasis caused by the parasite are currently increasing in  
4 China (Zhu and Xu 2014; Wang et al 2020).

5 Whole genome sequences are valuable resources for obtaining deeper understanding of the  
6 biology of any organism. In the case of *B. straminea*, such a resource will impact questions of how  
7 they may interact with *S. mansoni* and how similar the genetic mechanisms are between different  
8 *Biomphalaria* species, with possible implications for how treatments and management strategies  
9 might be transferable. To date, only the genome of *Biomphalaria glabrata* has been sequenced  
10 and analysed (Adema et al 2017; Tennessen et al 2020; Figure 1B), and a high-quality genome of  
11 *B. straminea* is lacking, hindering further understanding of the species. To address this issue, we  
12 provide and analyse a high-quality genome assembly for *B. straminea* together with accompanying  
13 transcriptomes.

14

15

16

17

18

19

## 1 **Results and Discussion**

### 2 ***Genome quality evaluation***

3           Genomic DNA was extracted from single individuals of *B. straminea* (Figure 1A). Genome  
4 sequences were first assembled using short-reads followed by scaffolding with Hi-C data. The  
5 genome assembly is 1.005 Gbp with a scaffold N50 of 25.3 Mbp (Figure 1B). This high physical  
6 contiguity is matched by high completeness, with an 87.0% complete BUSCO score (Simao et al  
7 2015)(Figure 1B). A total of 43,340 gene models, including 3,122 tRNA and 40,218 protein-  
8 coding genes, were generated by mapping transcriptome data to the genome assembly (S1.  
9 Sequencing data). The mean exon length is 262 bp, mean intron length is 1,603 bp, and mean  
10 deduced protein length is 377 aa. The genome quality generated in this study is comparable to the  
11 previously published genome assemblies of another schistosomiasis carrying vector snail, *B.*  
12 *glabrata* (Adema et al 2017; Tennessen et al 2020; Figure 1B).

13

### 14 ***Repeat element analysis***

15           We identified a total repeat content of 40.68% in the genome of *B. straminea* (Figure 1C),  
16 demonstrating that repeats make up a large proportion of total genome size in the species. A  
17 considerable proportion of repeats were unclassified (15.81%), suggesting that many of the  
18 annotated repeats represent new repeat families (Figure 1C), which is not unexpected given the  
19 relatively sparse attention given to the analysis of repeats in gastropod molluscs to date. Of the  
20 remaining repeats, LINE elements and DNA transposons are most abundant (LINEs: 10.48%,  
21 DNA transposons: 8.32%), whereas SINEs, LTR elements, and rolling-circle elements are present



1 only in low proportions (LTR elements: 2.7%, rolling-circle elements: 1.71%, SINEs: 1.31%)  
2 (Figure 1C). Consideration of a repeat landscape plot suggests that there has been a long-term  
3 ongoing expansion of repeats in *B. straminea*, with a recent spike in activity (Figure 1C). LINES  
4 and DNA transposons have expanded most significantly, however, there has also been a less  
5 considerable expansion of LTR and Rolling circle elements (Figure 1C).

6

### 7 ***Homeobox-containing gene content and linkage***

#### 8 a) Hox cluster genes

9 Homeobox genes are transcription factors involved in regulating animal development. Not  
10 only are they highly conserved between distantly related lineages, but also many of the genes are  
11 linked or clustered in genomes. Besides the most well-known clusters like the Hox and ParaHox  
12 clusters, many homeobox genes are linked including other ANTP class genes in NK and SuperHox  
13 clusters, and also amongst other classes of PRD, TALE, and SINE homeobox genes (Butts et al.,  
14 2008; Mazza et al., 2010; Ferrier, 2016). These clusters have been maintained or dispersed  
15 differently in different animal lineages. Changes to gene clustering may represent the breakdown  
16 of regulatory constraints which normally maintain clusters and are thought to be the mechanism  
17 holding together the tightly regulated Hox cluster, for instance. Genomic clustering also reflects  
18 the ancient origins of many of these homeobox genes by tandem duplication, e.g., the four ANTP  
19 clusters in the Bilaterian ancestor that arose via subsequent expansions from a single Proto-ANTP  
20 gene (Hui et al., 2012). Among molluscs, a diverse phylum to which gastropods belong, alongside  
21 other conchiferans (monoplacophorans, bivalves, scaphopods, and cephalopods), as well as  
22 aculiferans (aplacophorans and polyplacophorans), some of the diversity of body plans may be  
23 underpinned by changes to developmental genes like homeobox genes. Hox genes have been co-

1 opted to novel structures in cephalopods (Lee et al., 2003), and this corresponds to a breakdown  
2 of the Hox cluster across several chromosomes, and the loss of a few genes (Albertin et al., 2015).  
3 Other mollusc genomes show a breakdown of homeobox clustering overall, like the Pacific oyster  
4 (*Crassostrea gigas*; Paps et al., 2015), while a more recent chromosome-level assembly reveals  
5 large-scale patterns of linkage in another oyster (Li et al., 2020). This genome assembly of *B.*  
6 *straminea* improves our understanding of homeobox gene linkage in comparison to other molluscs  
7 and well-studied ecdysozoans like flies or vertebrates.

8         We found 114 homeobox genes in the genome of *B. straminea*, belonging to eleven  
9 recognized classes and one lophotrochozoan-specific gene, *Lopx* (Supplementary information S2a;  
10 Barton-Owen et al., 2018). Many of these genes are clustered or linked in the genome (Figure 2).  
11 Nine of the eleven Hox genes are found on scaffold 32695, in an arrangement that suggests several  
12 intrachromosomal rearrangements. In an ordered cluster as seen in *L. gigantea*, for instance, the  
13 Hox genes are situated in the genome in the ancestral order from anterior-acting *Hox1* to posterior-  
14 acting *Post1*, and no other non-Hox genes are found amongst the Hox genes (Simakov et al., 2013).  
15 Here, however, we find that *Hox2*, *Hox3*, and *Hox4* are downstream of *Hox5*. In addition, *Hox2*-  
16 *Hox5* are downstream of the posterior half of the cluster, including *Lox5*, *Hox7*, *Lox4*, *Lox2*, and  
17 *Post1*. *Hox1* is found on another scaffold, while the sequence for *Post2* is not in the genomic  
18 assembly, though its sequence is found in our transcriptome data. The Hox arrangement in *B.*  
19 *straminea* provides more linkage information than the *B. glabrata* assembly, where the short  
20 scaffolds corroborate only fragments of the Hox cluster like the linkage of *Hox4*, *Hox3*, and *Hox2*,  
21 but do not confirm the rearrangements in *B. straminea*, such as the linkage of *Hox5* to *Hox2*  
22 (Supplementary information S2b). We do see a difference in the arrangement of the posterior half  
23 of the Hox cluster, however, where in *B. glabrata*, *Lox4*, *Lox2*, *Post2*, and *Post1* are linked in that

1 order on scaffold 139, with *Lox4* and *Lox2* in the negative strand and *Post2* and *Post1* on the  
2 positive, which is slightly different from many other molluscs in which only *Post1* differs in  
3 orientation relative to the remainder of the posterior end of the Hox cluster genes (Simakov et al.,  
4 2013; Li et al., 2020). In *B. straminea*, there has been a rearrangement separating *Post1*, placing it  
5 with *Lox5* and *Hox7* and in the same orientation as *Lox4* and *Lox2* (Figure 2). Thus, the Hox genes  
6 of *Biomphalaria* seem highly rearranged relative to the ancestral order and each other. Clearly  
7 then, there are no (or minimal) long-range regulatory mechanisms operating across these genes  
8 that could have constrained their organization and prevented rearrangement. At most, there may  
9 be remains of some form of sub-cluster mechanisms, such as enhancer sharing, operating over the  
10 small regions (i.e. *Hox2-4* and *Lox2-4*) whose similar arrangement may be indicative of constraints  
11 conserved across *Biomphalaria* species. Future expression and regulatory element analyses may  
12 help resolve this possibility.

### 13 b) ParaHox cluster genes

14 The ParaHox cluster is the evolutionary sister to the Hox cluster (Brooke et al., 1998). The  
15 homeodomains of the three ParaHox genes (*Gsx*, *Xlox* and *Cdx*) are found on three separate  
16 scaffolds in *B. straminea* (Figure 2), however, three upstream exons of *Cdx* are on scaffold 5393,  
17 which also has the *Xlox* gene (Supplementary information S2a). This is in contrast to the genome  
18 of *B. glabrata*, where *Gsx* and *Xlox* are linked on scaffold 3 (Supplementary information S2a-b).  
19 Perhaps this pattern reflects maintained linkage between all three ParaHox genes in *Biomphalaria*  
20 species and only because of the draft level of all the assemblies this is not evident. However, if  
21 this is the case, the ParaHox genes are separated by large amounts of sequence and have not  
22 retained the ancestral order of *Gsx-Xlox-Cdx*. *B. glabrata Xlox* is nearly 4 Mb from the start of its

1 scaffold, while in *B. straminea*, *Xlox* is at a location with another homeobox-containing gene  
2 (*Phox*) 15 Mb away on one side and the first three *Cdx* exons are almost 5 Mb away on the other  
3 side of *Xlox*. Thus, although the *Biomphalaria* ParaHox genes may be linked, they cannot be  
4 considered to be clustered. This dispersal of ParaHox genes is typical for molluscs in general, with  
5 several species also showing loose linkage of some of the genes (Li et al., 2020), which contrasts  
6 with the relatively tight clustering of these genes in many deuterostomes (Osborne et al., 2009;  
7 Ikuta et al., 2013; Zhang et al., 2017) and the likely pan-cluster regulation that may operate in these  
8 deuterostomes.

9 c) ANTP-class homeobox genes

10 Beyond Hox and ParaHox, there are other linkages among and between the classes of  
11 homeobox genes that hint at their ancient evolutionary origins and genomic arrangement in clusters.  
12 Despite the many rearrangements to the Hox cluster, many genes linked to Hox clusters in other  
13 species are also found on the same scaffold in *B. straminea*, including *Mnx*, *Gbx-a* and *Gbx-b*, *En-*  
14 *a*, *Evx-a* and *Evx-b*, and *Dlx* (Castro and Holland, 2003; Chourrout et al., 2006; Butts et al., 2008;  
15 Hui et al., 2012; Li et al., 2020). These linkages give further support for the hypothesized Super-  
16 Hox cluster of non-Hox ANTP-class genes linked to the Hox genes in bilaterians (Butts et al.,  
17 2008).

18 d) SINE homeobox genes

19 Another highly conserved cluster besides Hox and ParaHox is the SINE-class cluster,  
20 typically composed of the *Six3/6*, *1/2*, and *4/5* genes or their protostome orthologues (Ferrier,  
21 2016). In *B. straminea*, *Six4/5* and *Six1/2* are on the same scaffold, but with a number of genes

1 between them, and *Six3/6* is on a distinct scaffold (Figure 2). In *B. glabrata*, *Six3/6* is linked to  
2 *Hlx* (Figure S2b), the last homeobox gene at the end of the *Six4/5-Six1/2* scaffold in *B. straminea*  
3 (Figure 2). Thus, there is clearly not a SINE-class gene cluster conserved in *B. straminea*, but the  
4 linkage of at least some of these genes indicates that the dispersal of this cluster has not yet  
5 proceeded to the extent of these genes being separated onto different chromosomes. Also, the  
6 location of the *Hlx* gene relative to different *Six* genes indicates a certain degree of genomic  
7 rearrangement between the two *Biomphalaria* species (i.e. conserved macrosynteny, but divergent  
8 microsynteny).

9 e) IRX homeobox genes

10 Homeobox genes in the IRX family within the TALE class, are also observed to be clustered in  
11 several lineages, for instance the three-gene (*ara*, *caup*, and *mirr*) cluster in *Drosophila*, two three-  
12 gene clusters in vertebrates, and four genes in the limpet *L. gigantea* (*irx4*, *irx2*, *irx1*, and *irx3*)  
13 (Irimia et al. 2008; Takatori et al., 2008; Kerner et al. 2009). These clusters are thought likely to  
14 have arisen convergently by independent tandem duplications in the arthropod, vertebrate, and  
15 mollusc lineages (Irimia et al., 2008; Takatori et al., 2008; Kerner et al., 2009; Chipman et al.,  
16 2014). Both *Biomphalaria* species have five IRX-family genes, one pair of which appears to be a  
17 product of a more recent, possibly *Biomphalaria*-specific, duplication (*Irx1-a* and *Irx1-b*). Perhaps  
18 surprisingly, none of the *Biomphalaria* *Irx* genes, *Irx1* (*a* and *b*), *Irx2*, *Irx3*, and *Irx4*, show clear  
19 orthology to specific limpet or oyster genes in a phylogenetic tree (Supplementary information  
20 S2c). A paucity of phylogenetically-informative amino-acid changes is the most likely explanation  
21 for this lack of resolution. Despite this lack of resolution of *Irx* orthology across species the *B.*  
22 *straminea* genome assembly does provide a new example of *Irx* gene clustering. *Irx3*, *Irx2*, and

1 *Irx4* are closely clustered in the genome, while *Irx1-b* is 7 Mb away on the same scaffold, also  
2 with *Zhx*, a ZF-class gene another 6 Mb further. The two *Irx1* paralogues, however, are on separate  
3 scaffolds, which may represent either a rearrangement following their duplication, convergence of  
4 the sequence of the homeodomain, or thirdly, an assembly artefact. In *B. glabrata*, only the linkage  
5 of *Irx4* with *Irx2* is corroborated due to the shorter scaffold lengths of that assembly. Further work,  
6 perhaps using other conserved domains from these genes and with a wider breadth of  
7 lophotrochozoan species could potentially determine whether in fact the four *Irx* gene types in  
8 *Biomphalaria* species are orthologous to genes in other species' *Irx* clusters. A multi-gene IRX-  
9 family cluster in *Biomphalaria* species with evidence of at least one independent expansion (*Irx1-*  
10 *a* and *Irx1-b*) provides an interesting addition to our understanding of IRX-family clusters, and the  
11 mechanisms behind gene expansions and subsequent maintenance of clustering in general.

12 f) PRD- and LIM- class homeobox genes

13 We also observe linkages amongst PRD-class genes, with clusters on scaffolds 13536,  
14 2216, 46009, and 563 (Figure 2). The widely found PRD-class cluster is the so-called HRO cluster,  
15 composed of the genes *Otp*, *Rx/Rax* and *Hbn/Arx-like* (Mazza et al., 2010; Ferrier, 2016), which  
16 ancestrally was likely embedded within a more extensive PRD/LIM-class mega-cluster, including  
17 the PRD-class genes *Gsc* and *Otx* and the LIM-class gene *Isl* (Ferrier, 2016). In *B. straminea* there  
18 is a remnant of the HRO cluster, with *Otp* clustered with *Hbn*, internally on a large scaffold (563)  
19 and flanked by other homeobox genes (Figure 2) including another PRD-class gene (*Arx-a*) now  
20 in this *Biomphalaria* PRD-class cluster, but the *Rax* genes are on other scaffolds. Interestingly, the  
21 *Isl* gene is also on this large 563 scaffold in *B. straminea*, consistent with descent from the  
22 hypothesized PRD/LIM-class mega-cluster (Ferrier, 2016). *B. glabrata* provides an interesting

1 contrast as the HRO cluster is now complete (with *Otp*, *Hbn* and *Rax-b*) in contrast to *B. straminea*,  
2 and again *Arx-a* is also in the *Biomphalaria* cluster (Figure 2; Supplementary Figure S2b). Why  
3 the PRD-class HRO cluster would remain intact in one species of *Biomphalaria* but not the other  
4 remains to be resolved. Also, whether the inclusion of the *Arx-a* gene in this cluster in these snails  
5 is found elsewhere in the animal kingdom and is of any functional significance also remains a topic  
6 for future work. Overall, the PRD-class gene clustering provides a mixed signal, of both  
7 conservation of remnants of ancient clustering alongside rearrangements between closely related,  
8 con-generic species.

9 g) Duplicated homeobox genes

10 There are several duplications shared between the two species, which we infer to be at least  
11 ancestral to the genus. These include paralogues of *Arx*, *Pax4/6*, *Irx1*, *En*, *Evx*, *Abox*, *Barhl*, *Pbx*,  
12 and *Tlx*, as well as three paralogues of *Vsx* and *Cers*. Notably, the three paralogues each of *Vsx*  
13 and *Cers* genes remain clustered in the genome, reflecting their likely origin by tandem duplication.  
14 This is also seen for *En*, *Tlx*, *Evx*, and *Abox*. *B. straminea* is the only species of the two with two  
15 paralogues of *Gbx*, though one has an apparently odd arrangement that would mean it is unlikely  
16 to be a functional gene, if this arrangement were real. The homeodomain is split across two exons,  
17 the first of which is in one orientation, while there are two copies of the second exon in the opposite  
18 orientation, indicating the second *Gbx* gene may be a pseudogene or an assembly artefact  
19 (Supplementary information S2a).

20 h) Giga-cluster homeobox genes

1 An overarching framework for understanding the genomic organization of homeobox-  
2 containing genes comes from hypotheses about their ancient linkage patterns following their  
3 presumed origins largely via tandem duplications. This ancestral clustering goes beyond the class-  
4 specific clusters already described above and is captured by the Giga-cluster hypothesis (Ferrier,  
5 2016). High-quality genome assemblies, such as the one described here for *B. straminea*, are key  
6 resources for testing this hypothesis and potentially expanding it. Several instances of linkage of  
7 different classes of homeobox gene are present in the *B. straminea* assembly, most notably on  
8 scaffolds 563, 8789, 2216 and 24987 (Figure 2). Scaffold 2216 is interesting for the linkage of the  
9 SINE-class genes *Six4/5* and *Six1/2* with some of the members of the ancestral PRD/LIM-class  
10 Mega-cluster (i.e. the PRD-class genes *Gsc* and *Otx*) that has undergone some dispersal in  
11 the *Biomphalaria* lineage (as described above). Also, some of the other members of this dispersed  
12 PRD/LIM Mega-cluster (*Isl*, *Otp*, *Hbn*) are on scaffold 563, which are now linked with many  
13 members of the dispersed NK-cluster (e.g. *NK5*, *NK4*, *Msx*, *Tlx-a* and *-b*, and *NK3*) as well as a  
14 member of the ancestral SuperHox cluster (i.e. *Hhex*) (Butts et al., 2008; Ferrier, 2016). Other  
15 members of the SuperHox cluster are still linked with the true Hox genes (EuHox genes) on  
16 scaffold 32695. These linkages of genes from different homeobox classes along with the further  
17 new instances of inter-class linkage on scaffolds 8798 (Figure 2) are all consistent with the Giga-  
18 cluster hypothesis (Ferrier, 2016). However, how much of all of these linkages represent ancestral  
19 associations (i.e. descended from primary clustering) versus instances of coming together in the  
20 genome convergently in evolution (i.e. secondary clustering) should be resolvable with  
21 comparisons to further high-quality genome sequences as well as a better understanding of the  
22 dynamics of genome evolution and rearrangements (reviewed in Ferrier, 2016).



## 1 *Synteny analysis of B. straminea with other molluscs*

2           The homeobox analyses described above provide instances of linkages that indicate varied  
3 synteny conservation across various mollusc and animal clades, even between the two  
4 *Biomphalaria* species now sequenced. The *B. straminea* genome shows considerable conserved  
5 linkage within and between classes of homeobox, and the maintenance of certain conserved  
6 clusters or linkages observed throughout wider lineages (i.e. instances of remnants of the Hox,  
7 ParaHox, SuperHox, and Giga-clusters (Ferrier, 2016)). In comparison to *B. glabrata*, in which  
8 less linkage can be observed because of shorter scaffold lengths, there is some conserved synteny.  
9 A few differences between the species may be due to species-specific genomic rearrangements  
10 resulting in the disruption of gene order, but the alternative possibility of assembly artefacts cannot  
11 be excluded entirely at present without further work. Of particular interest for further study is the  
12 major rearrangement of the Hox cluster in *B. straminea*. Perhaps more thorough sequencing of *B.*  
13 *glabrata* could determine if this is shared in the genus, or if it is a novelty of *B. straminea*.  
14 Regardless of this, the impact of this rearrangement on Hox gene expression and function is of  
15 interest. Similarly, the impacts of the dispersal of the ParaHox cluster on gene expression will be  
16 interesting to resolve. Homeobox genes are good markers for genome organization, and these  
17 results show that key differences between the species may represent higher levels of genomic  
18 divergence than expected for these two snails. Here we observe specific cases of differences  
19 between our new *B. straminea* genome and that of *B. glabrata* within the context of ancestral  
20 linkages, and this pattern may be a good indicator of wider differences between the genetics and  
21 molecular processes operating in the two species.

22           To examine the syntenic relationships more generally between *Biomphalaria* and mollusc  
23 genomes, we constructed Oxford dot-plots, comparing the chromosomal positions of orthologous

1 genes between mollusc genomes. As shown in Figure 4, the relationship of pseudo-chromosomes  
2 and scaffolds between *B. straminea* and molluscs from other genus was conserved in most cases.  
3 Previous phylogenetic tree constructions for different *Biomphalaria* species suggested a  
4 monophyletic clade of African species with the remaining lineages being neotropical species.  
5 Based on this phylogenetic relationship, our data show that the neotropical species have not  
6 undergone any significant inter-chromosomal rearrangements from their last common ancestor  
7 after separation to different geographical regions. One-to-one synteny block could be identified  
8 between *B. straminea* and *Achatina immaculata*. However, in the comparison of *B. straminea* to  
9 the more evolutionary distant species, a few one-to-many blocks were found. These patterns  
10 indicated that some chromosome duplication and alteration occurred from the most recent common  
11 ancestor of *B. straminea*, *B. glabrata* and *A. immaculata*. Further, species with closer evolutionary  
12 distance shared more similar synteny patterns against *B. straminea* (for example, *Pomacea*  
13 *canaliculata* and *Marisa cornuarietis*, *Crassostrea gigas* and *Magallana hongkongensis*),  
14 suggesting the dynamic changes of chromosomes arrangement in different molluscs. In *Octopus*  
15 *sinensis*, the gene order and synteny blocks to *B. straminea* were largely lost suggesting more  
16 duplication, translocation and rearrangement events occurred since the divergence of *O. sinensis*  
17 (Cephalopoda) and the common ancestor of Gastropoda and Bivalvia.

18

### 19 *Ecdysteroid genes*

20 Ecdysteroids play important roles in regulating growth (in particular molting and metamorphosis)  
21 and sexual maturation of insects and other arthropods (Cheong et al 2015; Qu et al 2015). Although it has  
22 long been known that gastropods contain ecdysteroids, and that beta-ecdysone could stimulate host location  
23 activities in *S. mansoni* miracidia and enhance growth and egg production in *B. glabrata* (Bayne 1972; Shiff

1 and Dossaji 1991), the biosynthetic pathway genes for ecdysteroids have not been systemically studied in  
2 mollusc genomes to date. As shown in Figure 3A-B, typical genes involved in this pathway including  
3 *CYP307A1*, *CYP306A1*, *CYP302A1*, *CYP315A1*, *CYP314A1*, and *CYP18A1* are all absent from the *B.*  
4 *straminea* genome assembly and transcriptome data. Nevertheless, the receptors including EcR, RXR/USP  
5 and oxygenase-like protein Nvd that are essential regulators of cholesterol metabolism are revealed in *B.*  
6 *straminea* and other mollusc genomes (Figure 3A-B; Supplementary information S3). We thus treated *B.*  
7 *straminea* with  $10^{-6}$  M ecdysteroid 20-hydroxyecdysone for 24 hours but did not observe any significant  
8 expression changes in the downstream genes *E74*, *FOXO*, and *Nvd* (Figure 3C). It is unclear whether only  
9 certain forms of ecdysteroids may induce endogenous ecdysteroid pathway genes under particular  
10 conditions and this warrants further investigation. This is the first systematic analyses of ecdysteroid  
11 pathway genes in a mollusc genome, thus providing the foundations for future work to determine how  
12 ecdysteroids have their effect in these animals.

13

#### 14 ***Insulin signaling pathway genes***

15 Peptide hormones involved in growth and reproduction have been suggested as candidates  
16 for the development of novel methods of schistosomiasis control via manipulation of snail numbers  
17 (Acker et al 2019). Insulin is another understudied hormonal pathway in molluscs despite its  
18 potential functional roles. For instance, in the pond snail *Lymnaea stagnalis*, a decrease of insulin  
19 in the central nervous system correlated with high memory scores (Totani et al 2019), while  
20 insulin-related peptides with potential roles in sexual reproduction have been identified in the  
21 oyster *Crassostrea gigas* (Cherif-Feidel et al 2019). In both *B. straminea* and *B. glabrata* genomes,  
22 we were able to identify all key signalling pathway genes (Figure 3D-E, Supplementary S4). This  
23 establishes a foundation on which to further explore the functions of these hormones in molluscs.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

***Widespread gene turnover between Biomphalaria snails and other molluscs***

a) Gene gains and losses in mollusc genomes

A phylogenomic tree was constructed using 2,047 orthogroups with at least 12 out of 13 mollusc genomes having single-copy genes in each orthogroup (Supplementary information 6). Gene family analysis among these genomes revealed the expansion of 1,869 orthogroups and contraction of 623 orthogroups in *B. straminea* (Figure 5). This data highlights the importance of having the *B. straminea* genomic resource, and potentially suggested that specific control strategies might be needed for *B. straminea* rather than treating it as identical to *B. glabrata*.

b) Expansion of heat shock protein family in certain mollusc lineages

Heat shock proteins are important stress-responsive candidates involved in protein folding for molluscs, activated in response to such things as changing pH, oxygen level, and temperature. In some mollusc genomes, such as that of the Pacific oyster *Crassostrea gigas*, an expansion of heat shock protein 70 (HSP70) has been observed in the genome and hypothesized to be important to its adaptation (Zhang et al 2012). We thus identified the heat shock protein family genes in *Biomphalaria* and compared these to other lophotrochozoans to understand their evolution in different lineages (Figure 6). Among the different heat shock protein families in the investigated set of gastropods, bivalves, cephalopods, annelids, and platyhelminthes, a dramatic expansion is seen specifically in the HSP70 family in the bivalve molluscs (Figure 6; Supplementary information S7). Our data and analyses agree with previous studies (e.g. Zhang et al 2012), suggesting that the expansion of HSP70 is linked to the life history of molluscs having a sessile

1 stage. This survey also provides the foundation for future work on the expression and function of  
2 particular HSP genes/proteins and their activity in these parasite vectors, which may contribute to  
3 their adaptive ability as invasive species, and possibly contributing to the recent range expansion  
4 of *B. straminea*.

5

6 c) Differential sesquiterpenoid and cholesterol genes in certain mollusc lineages

7 Sesquiterpenoid hormones were once considered specific to insects and crustaceans where  
8 they control development and reproduction (Cheong et al 2015; Qu et al 2018; Tsang et al 2020).  
9 However, recent analyses have shown that the sesquiterpenoid system is also present in myriapods,  
10 annelids, and cnidarians (Chipman et al 2014; Qu et al 2015; Schenk et al 2016; Nong et al 2020).  
11 Conversely, vertebrates can only produce cholesterol but not sesquiterpenoids (Tobe and Bendena  
12 1999; Hui et al 2013), and a recent study revealed the canonical cholesterol biosynthesis pathway  
13 in sponges, placozoans and deuterostomes, suggesting cnidarians and protostomes experienced  
14 massive losses of these genes (Zhang et al 2019; Figure 7A). Treatment of  $10^{-6}$  M simvastatin and  
15 methyl farnesoate on the snail *B. straminea* can change the expression of sesquiterpenoid pathway  
16 genes HMGCR and FPPS, suggesting a sesquiterpenoid responsive system (Figure 7B-C).  
17 Comparison of sesquiterpenoid pathway genes in mollusc genomes further identified differential  
18 utilization of biogenesis pathways in bivalves and gastropods, where only gastropods but not the  
19 bivalves are able to produce cholesterol similar to vertebrates (Figure 7D-F). This is the first  
20 systematic study showing the differential sesquiterpenoid and cholesterol pathways taken by  
21 different mollusc lineages.

22

## 1 **Conclusion**

2 This study presents the first high quality genome assembly for a schistosomiasis-transmitting snail  
3 in China and Asia. The snail *Biomphalaria straminea* is important scientifically as well as holding  
4 considerable medical relevance. Our work also provides the dynamics of homeobox, ecdysteroid,  
5 insulin, heat shock protein, and sesquiterpenoid pathway genes, suggesting extensive molecular  
6 differences between *B. straminea* and *B. glabrata* as well as between molluscs. More generally,  
7 our high-quality *B. straminea* genome provides a useful reference point for further understanding  
8 molluscs biology, ecology and evolution.

9

10

## 11 **Methods**

### 12 *Sample collection and genome sequencing*

13 Specimens of the ramshorn snail (*B. straminea*) were collected from the New Territories,  
14 Hong Kong, and samples for genome sequencing originate from a single individual (Figure 1A).  
15 Genomic DNA (gDNA) was extracted using the PureLink Genomic DNA Mini Kit (Invitrogen)  
16 following the manufacturer's protocol. Extracted gDNA was subjected to quality control using a  
17 Nanodrop spectrophotometer (Thermo Scientific) and gel electrophoresis. Qualifying samples  
18 were sent to Novogene, and Dovetail Genomics for library preparation and sequencing. The  
19 resulting library was sequenced on an Illumina HiSeq X platform to produce  $2 \times 150$  paired-end  
20 sequences. The length-weighted mean molecule length is 22.2 kb, and the raw data can be found  
21 at NCBI's Small Read Archive (SRR12963913).

1

## 2 ***Dovetail Omni-C library preparation and sequencing***

3           For each Dovetail Omni-C library, chromatin was fixed with formaldehyde and extracted.  
4 Fixed chromatin was digested with DNase I, and chromatin ends were repaired and ligated to a  
5 biotinylated bridge adapter followed by proximity ligation of adapter containing ends. After  
6 proximity ligation, crosslinks were reversed and the DNA was purified. Purified DNA was treated  
7 to remove biotin that was not internal to ligated fragments. Sequencing libraries were generated  
8 using NEBNext Ultra enzymes and Illumina-compatible adapters. Biotin-containing fragments  
9 were isolated using streptavidin beads before PCR enrichment of each library. The library was  
10 sequenced on an Illumina HiSeqX platform to produce 128 million 150 bp read pairs, and the raw  
11 data can be found at NCBI's Small Read Archive (SRR12963914).

12

## 13 ***Transcriptome sequencing***

14           Total RNA from different tissues were isolated using a combination method of  
15 cetyltrimethylammonium bromide (CTAB) pre-treatment (Jordon-Thaden et. al. 2015) and  
16 mirVana™ miRNA Isolation Kit (Ambion) following the manufacturer's protocol. The extracted  
17 total RNA was subjected to quality control using a Nanodrop spectrophotometer (Thermo  
18 Scientific), gel electrophoresis, and an Agilent 2100 Bioanalyzer (Agilent RNA 6000 Nano Kit).  
19 Qualifying samples underwent library construction and sequencing at Novogene; polyA-selected  
20 RNA-Sequencing libraries were prepared using the TruSeq RNA Sample Prep Kit v2. Insert sizes  
21 and library concentrations of final libraries were determined using an Agilent 2100 bioanalyzer

1 instrument (Agilent DNA 1000 Reagents) and real-time quantitative PCR (TaqMan Probe)  
2 respectively. Details of the sequencing data can be found in Supplementary information S1.

3

#### 4 ***Genome assembly***

5 Chromium WGS reads were used to construct a *de novo* assembly using Supernova (v 2.1.1)  
6 with default parameters (raw coverage = 68.32x). The Supernova output pseudohap assembly and  
7 Dovetail OmniC library reads were used as input data for HiRise, a software pipeline designed  
8 specifically for using proximity ligation data to scaffold genome assemblies (Putnam et al, 2016).  
9 Dovetail OmniC library sequences were aligned to the draft input assembly using bwa  
10 (<https://github.com/lh3/bwa>). The separations of Dovetail OmniC read pairs mapped within draft  
11 scaffolds were analyzed by HiRise to produce a likelihood model for genomic distance between  
12 read pairs, and the model was used to identify and break putative misjoins, to score prospective  
13 joins, and make joins above a threshold.

14

#### 15 ***Gene model prediction***

16 Gene models were predicted as described in the Hong Kong oyster (*Magallana*  
17 *hongkongensis*) genome (Li et al. 2020). Briefly, the gene models were trained and predicted using  
18 funannotate (v1.7.4,<https://github.com/nextgenusfs/funannotate>) (Palmer & Stajich, 2020) with  
19 the following parameters: “--repeats2evm --protein\_evidence uniprot\_sprot.fasta --  
20 genemark\_mode ET --busco\_seed\_species metazoa --optimize\_augustus --busco\_db metazoa --  
21 organism other --max\_intronlen 350000”. The gene models from several prediction sources  
22 including GeneMark, high-quality Augustus predictions (HiQ), pasa, Augustus, GlimmerHM and



1 snap were passed to Evidence Modeler and generated the gene model annotation files, followed  
2 by PASA to update the EVM consensus predictions, and add UTR annotations and models for  
3 alternatively spliced isoforms. Protein-coding genes were searched with BLASTp against the nr  
4 and swissprot databases by diamond (v0.9.24) (Buchfink et al., 2014) with parameters “--more-  
5 sensitive --evaluate 1e-3”, and mapped by HISAT2 (version 2.1.0) with transcriptome reads. Gene  
6 models with no similarity to any known protein and no mRNA support were removed from the  
7 final version.

8

### 9 ***Repetitive elements annotation***

10         Repetitive elements were identified using an in-house pipeline as follows. Firstly, elements  
11 were identified using RepeatMasker v.4.1 (Smit et al., 2013), using a sensitive (-s) search and  
12 ignoring low-complexity repeats (-nolow). Subsequently, a *de novo* repeat library was constructed  
13 using RepeatModeler v.1.0.11 (Smit et al., 2015), including RECON v.1.08 (Bao., et al 2002) and  
14 RepeatScout v.1.0.5 (Price et al., 2005). Identified novel repeats were analysed using a ‘BLAST,  
15 Extract, Extend’ process to characterise elements along their entire length (Platt et al., 2016)[23];  
16 Consensus sequences and classifications for each repeat family were generated, and the resulting  
17 *de novo* repeat library was utilised to identify repetitive elements in RepeatMasker. All plots were  
18 generated using Rstudio ver. 1.2.1335 with R ver. 3.5.1 (Team, 2013) and ggplot2 ver. 3.2.1  
19 (Wickham, 2016).

20

### 21 ***Gene family annotation and gene tree building***

1 Gene family sequences were first retrieved from the *B. straminea* genome using the  
2 tBLASTn algorithm on a local server. The identity of each retrieved gene was then checked by  
3 reciprocal searches against the Genbank nr database at NCBI with BLASTx. For phylogenetic  
4 analyses of gene families, DNA sequences were first translated into amino-acid sequences and  
5 aligned to other reference sequences (extracted from NCBI) using Clustal W. Gapped sites were  
6 removed from alignments using MEGA 7.0, and phylogenetic trees (neighbor-joining) were  
7 constructed using MEGA 7.0, where each phylogenetic node was analysed using 1000 bootstrap  
8 replicates. For homeobox-containing genes, homeodomains were annotated using tBLASTn  
9 searches with HomeoDB sequences, and sequences from representative lophotrochozoan families,  
10 including the expanded Spiralia TALEs (Barton-Owen et al., 2018). We also removed redundant  
11 hits based on their unique locations in the genome sequence, and manually detected any likely  
12 artefactual duplicates which were not carried forward into the protein sequences alignments  
13 (Supplementary Table S2). Alignments of each class were made using MUSCLE (Edgar, 2004),  
14 with homeodomain sequences from human (*Homo sapiens*), amphioxus (*Branchiostoma floridae*),  
15 fruitfly (*Drosophila melanogaster*), the red flour beetle (*Tribolium castaneum*), an oyster  
16 (*Crassostrea gigas*), a limpet (*Lottia gigantea*), a brachiopod (*Lingula anatina*), and the annelids  
17 *Platynereis dumerilii* and *Capitella teleta*, where available from other studies (Paps et al., 2015;  
18 Barton-Owen et al., 2018) and HomeoDB (Ying-Fu et al., 2011; 2008). The best substitution  
19 models were tested with ModelFinder, and Maximum Likelihood phylogenies were constructed  
20 with IQ-TREE with 1000 bootstrap replicates (Nguyen et al., 2015).

21

22 ***Identification of orthologous genes and gene families***

1 Orthologues and orthogroups in *B. straminea* and 12 other animal proteomes were inferred  
2 using OrthoFinder v. 2.5.2 [28] with default values and ‘-M msa’ activated. To cover the gene  
3 families, the longest protein of each gene was taken as the representative in OrthoFinder analysis.  
4 Gene duplication events were then identified. Duplication ratios per node/tip were calculated by  
5 dividing the number of duplications observed in each node/tip by the total number of gene trees  
6 containing that node. CAFE5 was used to infer gene gain and loss rates [29]. Orthogroups from  
7 output of OrthoFinder were regarded as gene families and fed to CAFE5. A divergence tree was  
8 inferred using r8s [30] from the species tree generated by OrthoFinder. We tested several gamma  
9 rate categories (-k) and k=1 showed the best likelihood.

10

### 11 ***Functional terms enrichment analysis***

12 Orthogroups were assigned Gene Ontology (GO), EuKaryotic Orthologous Groups  
13 (KOG), Kyoto Encyclopedia of Genes and Genomes (KEGG), and KEGG Orthology (KO) terms  
14 by inheriting the terms from genes found within the groups. The functional term annotations were  
15 performed using eggNOG [31]. Functional enrichment was tested for using function  
16 ‘compareCluster()’ in R package ‘clusterProfiler’ v.3.16.1 [32] under the environment of R 4.0.4  
17 [33]. Significantly enriched terms were determined with pvalueCutoff = 0.05, pAdjustMethod =  
18 "BH", and qvalueCutoff = 0.2. Data was visualised using R packages ‘ggplot2’ [25], ‘ggtree’ [34]  
19 and ‘pathview’ [35].

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

***Macrosynteny analysis***

Single-copy orthologues anchored by mutual best Diamond blastp v0.9.14.115[17] hits (evaluate 0.001) between *B. straminea* and 12 other animals with chromosome-level or near chromosome-level assemblies were used in macrosynteny analysis. Oxford synteny plots were generated following previously described methods [36] using R packages ‘ggplot2’ [25].

***Drug and hormone treatment and RT-qPCR***

Adult animals from culture were rinsed in double-distilled water to remove any contaminants. Three individuals per set were placed in a glass container, with a well of 3.5cm in radius and 0.8cm in depth, filled with 2ml of double-distilled water with either  $10^{-6}$ M or  $10^{-8}$ M of methyl farnesoate (MF) (Sigma),  $6 \times 10^{-5}$ M of simvastatin (Sigma) or  $10^{-6}$ M of 20-hydroxyecdysone (AbcamBiochemicals) in separate setups. The chemicals were first dissolved in acetone and diluted to the target concentration in the treatment container. The control setup contained the same number of individuals and was treated with the same concentration of acetone in corresponding experiments. Each replicate of snails was exposed for 24 hours to these treatments without any feeding. Post-treated animals were rinsed with double-distilled water and shells were removed for whole body total RNA extraction. The RNA from each experiment was isolated using TRIzol reagent following the manufacturer's protocol. Purified RNA was dissolved in nuclease-free water. The cDNA synthesis was performed using the iScript gDNA Clear cDNA Synthesis Kit (BioRad) following the manufacturer's protocol. The cDNA was used in subsequent

1 quantitative real time PCR. The amplification conditions were as follows: initial denaturation at  
2 95 °C for 30 s, followed by 40 cycles of 95 °C denaturation for 15s, 57 °C primer annealing for  
3 15s and 72°C extension for 15s. Each sample was analyzed in replicates. The expression of each  
4 target gene transcript was normalized to the housekeeping gene, myoglobin (Myo), and fold  
5 induction analyses were calculated using the  $\Delta\Delta C_t$  method.

6

## 7 **Ethics Statement**

8 N/A

9

## 10 **Availability of Supporting Data and Materials**

11 The raw genome and RNA sequencing data have been deposited in the SRA under Bioproject  
12 number PRJNA673593. The final chromosome assembly was submitted to NCBI Assembly under  
13 accession number JADKLZ000000000 in NCBI. All data is available from the corresponding  
14 author upon reasonable request.

15

## 16 **Competing interests**

17 The authors declare no competing interests.

18

1 **Figure legends**

2 **Figure 1.** A) Life cycle of snail *Biomphalaria straminea*; B) Comparison of snail *Biomphalaria*  
3 genome assembly quality; C) Transposable elements in *Biomphalaria straminea*.

4 **Figure 2.** Distribution of Homeoboxes in the genome of *Biomphalaria straminea*. Class is denoted  
5 by colour, arrows show orientation on each scaffold, which are represented by black lines and are  
6 numbered underneath. *Post2* is not found in the genomic sequence but is found in the  
7 transcriptome, so is not shown on a scaffold. Grey gene names and box outlines denote partial  
8 homeodomain sequences.

9 **Figure 3.** A) Schematic diagram of biosynthetic pathway of ecdysteroids; B) Presence and absence  
10 of ecdysteroid pathway genes in *B. straminea*; C) Expression of genes upon  $10^{-6}$ M 20-  
11 hydroxyecdysone treatment for 24 hours (n=13-15); D) Schematic diagram of biosynthetic  
12 pathway of insulin; E) Number of gene copies of insulin pathway genes in *B. straminea*.

13 **Figure 4.** Synteny between *B. straminea* and other 12 mollusc genomes. The species tree is  
14 constructed using 2,047 orthogroups with at least 12 out of 13 mollusc genomes having single-  
15 copy genes in each orthogroup. In the Oxford dot plot, each dot represents a pair of orthologous  
16 genes between *B. straminea* and the specific mollusc. Horizontal and vertical dashed lines  
17 represent chromosome or scaffold boundaries. Orthologous genes are colored according to their  
18 position in *B. straminea* scaffolds. Significance of synteny blocks is computed using one-tailed  
19 Fisher's exact test, and the color of synteny blocks with Benjamini & Hochberg corrected p over  
20 0.05 are turned into grey.

1 **Figure 5.** Summaries of gene families in *B. straminea* and other 12 mollusc. A) Gene family  
2 clustering, only the longest isoform for each gene was used; B) Gene family expansion and  
3 contraction between mollusc genomes. Brown and green color indicate the number of significantly  
4 ( $p < 0.05$ ) expanded or contracted gene families at each node, respectively.

5 **Figure 6.** A) Schematic diagram showing the heat shock proteins actions; B) Number of gene  
6 copies of heat shock proteins in different mollusc genomes. The purple box highlights the  
7 expansion of HSP70 in certain mollusc lineages.

8 **Figure 7.** A) Schematic diagram showing the mevalonate pathway, and the downstream  
9 sesquiterpenoid and *de novo* cholesterol synthesis pathways. B) Expression of genes upon  $6 \times 10^{-5}$   
10 M simvastatin,  $10^{-6}$ M and  $10^{-8}$ M methyl farnesoate treatment for 24 hours; \* =  $p < 0.05$ . C)  
11 Heatmap of mevalonate pathway orthologues identified in gastropod and bivalve genomes. D)  
12 Heatmap of sesquiterpenoid synthesis pathway orthologues identified in gastropod and bivalve  
13 genomes. E) Heatmap of *de novo* cholesterol synthesis pathway orthologues identified in  
14 gastropod and bivalve genomes. F) Schematic diagram showing the evolution of sesquiterpenoid  
15 pathway genes in bilaterians.

16

17 **Additional Files.**

18 **Supplementary information S1.** Sequencing data.

19 **Supplementary information S2.** a) Tables of homeobox genes sequences in *B. straminea*, *B.*  
20 *glabrata*, a synteny comparison of homeobox genes, and comparison of ParaHox gene linkage. b)

1 Distribution of Homeoboxes in the genome of *Biomphalaria glabrata*. c) Alignments and  
2 phylogenies of each class of Homeobox sequences.

3 **Supplementary information S3.** Ecdysteroid genes.

4 **Supplementary information S4.** Insulin pathway genes.

5 **Supplementary information S5.** Synteny information

6 **Supplementary information S6.** Gene expansion and contraction.

7 **Supplementary information S7.** Heat shock protein family genes.

8 **Supplementary information S8.** Cholesterol genes and primers.

9 **Supplementary information S9.** Phylogenetic trees.

10 **Supplementary information S10.** Tables.

11 **Abbreviations.**

12 BLAST: Basic Local Alignment Search Tool; BUSCO: Benchmarking Universal Single-Copy  
13 Orthologs; kb: kilobase pairs; Mb: megabase pairs; NCBI: National Center for Biotechnology  
14 Information; TE: transposable element

15

16 **Competing Interests**

17 The authors declare that they have no competing interests.



1

## 2 **Funding**

3 This work was supported by the Hong Kong Research Grant Council Collaborative Research Fund  
4 (C4015-20EF), General Research Fund (14100919), and The Chinese University of Hong Kong  
5 Direct Grant (133134084). YY, WLS, CFW, STSL, and YL were supported by the PhD  
6 studentships of The Chinese University of Hong Kong. AH is supported by a Biotechnology and  
7 Biological Sciences Research Council (BBSRC) David Phillips Fellowship (BB/N020146/1). TB  
8 is supported by a studentship from the Biotechnology and Biological Sciences Research Council-  
9 funded South West Biosciences Doctoral Training Partnership (BB/M009122/1). MEAR is  
10 supported by a PhD studentship from the School of Biology and St Andrews University.

11

## 12 **Authors' Contributions**

13 JHLH, DEKF, AH, ZW, SX, ZPK, SSC conceived the study. JHLH, DEKF, AH supervised the  
14 study. WN, JH, TS assembled the genome. WN carried out the gene model prediction and  
15 comparison. YY carried out the heat shock proteins analyses. YX carried out the gene gain and  
16 loss and synteny analyses. WLS and CFW carried out the sesquiterpenoid analyses. YY, WLS and  
17 SYL carried out the ecdysteroid analyses. MEAR and YL carried out the homeobox gene analyses.  
18 TB carried out the transposable element analyses. STSL carried out the insulin analyses. WN,  
19 YY, YX, WLS, MEAR, TB, AH, DEKF, JHLH wrote the first draft of manuscript. All authors  
20 approved the final version of the manuscript.

21

1 **Acknowledgements**

2 The authors would like to thank Elaine Huang and Ho Yin Yip for collection and maintenance of  
3 snails. We thank Thomas Barton-Owen for help and advice on homeobox searches.

4

5

6

7

8

9

10

11 **References**

12 1. Campbell G, Jones CS, Lockyer AE, Hughes S, Brown D, Noble LR, et al.. Molecular evidence  
13 supports an African affinity of the Neotropical freshwater gastropod, *Biomphalaria glabrata*, Say  
14 1818, an intermediate host for *Schistosoma mansoni*. *Proc R Soc B Biol Sci*. Royal Society; 2000;  
15 doi: 10.1098/rspb.2000.1291.

16 2. Meier Brook C. A snail intermediate host of *Schistosoma mansoni* introduced into Hong Kong.  
17 *Bull World Health Organ*. World Health Organization; 51:6611975;

18 3. Attwood SW, Huo GN, Qiu JW. Update on the distribution and phylogenetics of *Biomphalaria*  
19 (*Gastropoda: Planorbidae*) populations in Guangdong Province, China. *Acta Trop*. Elsevier; 2015;

- 1 doi: 10.1016/j.actatropica.2014.04.032.
- 2 4. Dudgeon D, Yipp MW. A report on the gastropod fauna of aquarium fish farms in Hong Kong,  
3 with special reference to an introduced human schistosome host species, *Biomphalaria straminea*  
4 (Pulmonata: Planorbidae). *Malacol Rev* 16. :93–4 1983;
- 5 5. Meier-Brook C. A snail intermediate host of *Schistosoma mansoni* introduced into Hong Kong.  
6 *Bull World Health Organ.* 51:6611974;
- 7 6. Woodruff DS, Mulvey M, Yipp MW. The continued introduction of intermediate host snails to  
8 *Schistosoma mansoni* into Hong Kong. *Bull World Heal Organ.* 631985;
- 9 7. Zeng X, Yiu WC, Cheung KH, Yip HY, Nong W, He P, et al.. Distribution and current infection  
10 status of *Biomphalaria straminea* in Hong Kong. *Parasites and Vectors.* BioMed Central Ltd.;;  
11 2017; doi: 10.1186/s13071-017-2285-3.
- 12 8. Zhu R, Xu J. Epidemic situation of imported schistosomiasis in China: prevention and control.  
13 *Chin J Schisto Control.* 262014;
- 14 9. Wang L, Wu X, Li X, Zheng X, Wang F, Qi Z, et al.. Imported Schistosomiasis: A New Public  
15 Health Challenge for China. *Front Med.* Frontiers Media S.A.; 2020; doi:  
16 10.3389/fmed.2020.553487.
- 17 10. Adema CM, Hillier LDW, Jones CS, Loker ES, Knight M, Minx P, et al.. Whole genome  
18 analysis of a schistosomiasis-transmitting freshwater snail. *Nat Commun.* Nature Publishing  
19 Group; 2017; doi: 10.1038/ncomms15451.
- 20 11. Tennesen JA, Bollmann SR, Peremyslova E, Kronmiller BA, Sergi C, Hamali B, et al..  
21 Clusters of polymorphic transmembrane genes control resistance to schistosomes in snail vectors.  
22 *Elife.* eLife Sciences Publications Ltd; 2020; doi: 10.7554/ELIFE.59395.

- 1 12. Jordon-Thaden IE, Chanderbali AS, Gitzendanner MA, Soltis DE. Modified CTAB and TRIzol  
2 Protocols Improve RNA Extraction from Chemically Complex Embryophyta. *Appl Plant Sci.*  
3 Wiley; 2015; doi: 10.3732/apps.1400105.
- 4 13. Putnam NH, Connell BO, Stites JC, Rice BJ, Hartley PD, Sugnet CW, et al.. Chromosome-  
5 scale shotgun assembly using an in vitro method for long-range linkage arXiv : 1502 . 05331v1 [  
6 q-bio . GN ] 18 Feb 2015. *Genome Res.* 2016; doi: 10.1101/gr.193474.115.Freely.
- 7 14. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
8 *Bioinformatics.* 2009; doi: 10.1093/bioinformatics/btp324.
- 9 15. Li Y, Nong W, Baril T, Yip HY, Swale T, Hayward A, et al.. Reconstruction of ancient  
10 homeobox gene linkages inferred from a new high-quality assembly of the Hong Kong oyster  
11 (*Magallana hongkongensis*) genome. *BMC Genomics.* 2020; doi: 10.1186/s12864-020-07027-6.
- 12 16. Palmer JM, Stajich J. nextgenusfs/funannotate: funannotate v1.7.4 (Version 1.7.4). Zenodo.
- 13 17. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat*  
14 *Methods* 2014 121. Nature Publishing Group; 2014; doi: 10.1038/nmeth.3176.
- 15 18. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and  
16 genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol.* Nature Publishing Group; 2019;  
17 doi: 10.1038/s41587-019-0201-4.
- 18 19. Smit AFA, Hubley RR, Green PR. RepeatMasker Open-4.0. <http://repeatmasker.org>.
- 19 20. Smit A, Hubley R. RepeatModeler Open-1.0. <http://repeatmasker.org>. 2015;
- 20 21. Bao Z, Eddy SR. Automated De Novo Identification of Repeat Sequence Families in  
21 Sequenced Genomes. doi: 10.1101/gr.88502.
- 22 22. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes.

- 1 *Bioinformatics*. 2005; doi: 10.1093/bioinformatics/bti1018.
- 2 23. Platt RN, Blanco-Berdugo L, Ray DA. Accurate transposable element annotation is vital when  
3 analyzing new genome assemblies. *Genome Biol Evol*. 2016; doi: 10.1093/gbe/evw009.
- 4 24. Team RC. R: A language and environment for statistical computing. Vienna, Austria; 2013;
- 5 25. Wickham H. ggplot2: Elegant graphics for data analysis. New York: Springer-Verlag;
- 6 26. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J*  
7 *Mol Biol*. 1990; doi: 10.1016/S0022-2836(05)80360-2.
- 8 27. Kumar S, Stecher G, Tamura K. MEGA7: Molecular Evolutionary Genetics Analysis Version  
9 7.0 for Bigger Datasets. *Mol Biol Evol*. 2016; doi: 10.1093/molbev/msw054.
- 10 28. Emms DM, Kelly S. OrthoFinder: Phylogenetic orthology inference for comparative genomics.  
11 *Genome Biol*. BioMed Central Ltd.; 2019; doi: 10.1186/s13059-019-1832-y.
- 12 29. Mendes FK, Vanderpool D, Fulton B, Hahn MW. CAFE 5 models variation in evolutionary  
13 rates among gene families. Robinson P, editor. *Bioinformatics*. Oxford University Press (OUP);  
14 2020; doi: 10.1093/bioinformatics/btaa1022.
- 15 30. Sanderson MJ. r8s: Inferring absolute rates of molecular evolution and divergence times in the  
16 absence of a molecular clock. *Bioinformatics*. Oxford Academic; 2003; doi:  
17 10.1093/bioinformatics/19.2.301.
- 18 31. Huerta-Cepas J, Szklarczyk D, Heller D, Hernández-Plaza A, Forslund SK, Cook H, et al..  
19 eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based  
20 on 5090 organisms and 2502 viruses. *Nucleic Acids Res*. 2018; doi: 10.1093/nar/gky1085.
- 21 32. Yu G, Wang L-G, Han Y, He Q-Y. clusterProfiler: an R package for comparing biological  
22 themes among gene clusters. *OMICS*. 2012; doi: 10.1089/omi.2011.0118.

- 1 33. R Core Team. R: a language and environment for statistical computing. Vienna, Austria;
- 2 34. Yu G, Smith DK, Zhu H, Guan Y, Lam TTY. ggtree: an r package for visualization and  
3 annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol*  
4 *Evol.* British Ecological Society; 2017; doi: 10.1111/2041-210X.12628.
- 5 35. Luo W, Brouwer C. Pathview: An R/Bioconductor package for pathway-based data integration  
6 and visualization. *Bioinformatics.* Oxford Academic; 2013; doi: 10.1093/bioinformatics/btt285.
- 7 36. Simakov O, Marlétaz F, Yue JX, O'Connell B, Jenkins J, Brandt A, et al.. Deeply conserved  
8 synteny resolves early events in vertebrate evolution. *Nat Ecol Evol.* Springer US; 2020; doi:  
9 10.1038/s41559-020-1156-z.
- 10 37. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. BUSCO: Assessing  
11 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.*  
12 Oxford University Press; 2015; doi: 10.1093/bioinformatics/btv351.
- 13 38. Cheong SPS, Huang J, Bendena WG, Tobe SS, Hui JHL. Evolution of ecdysis and  
14 metamorphosis in arthropods: The rise of regulation of juvenile hormone. *Integr Comp Biol.*  
15 Oxford University Press;
- 16 39. Qu Z, Kenny NJ, Lam HM, Chan TF, Chu KH, Bendena WG, et al.. How Did Arthropod  
17 Sesquiterpenoids and Ecdysteroids Arise? Comparison of Hormonal Pathway Genes in Noninsect  
18 Arthropod Genomes. *Genome Biol Evol.* Oxford University Press; 2015; doi: 10.1093/gbe/evv120.
- 19 40. Bayne CJ. On the reported occurrence of an ecdysone-like steroid in the freshwater snail,  
20 *Biomphalaria glabrata* (Pulmonata; Basommatophora), intermediate host of *Schistosoma mansoni*.  
21 *Parasitology* . 1972 Jun;64(3):501-9. doi: 10.1017/s003118200004556x.
- 22 41. Shiff CJ, Dossaji SF. Ecdysteroids as regulators of host and parasite interactions: a study of

- 1 interrelationships between *Schistosoma mansoni* and the host snail, *Biomphalaria glabrata*. *Trop*  
2 *Med Parasitol.* 1991 Mar;42(1):11-6.
- 3 42. Acker MJ, Habib MR, Beach GA, Doyle JM, Miller MW, Croll RP. An immunohistochemical  
4 analysis of peptidergic neurons apparently associated with reproduction and growth in  
5 *Biomphalaria alexandrina*. *Gen Comp Endocrinol.* Academic Press Inc.; 2019; doi:  
6 10.1016/j.ygcen.2019.03.017.
- 7 43. Totani Y, Aonuma H, Oike A, Watanabe T, Hatakeyama D, Sakakibara M, et al.. Monoamines,  
8 insulin and the roles they play in associative learning in pond snails. *Front. Behav. Neurosci.*  
9 *Frontiers Media S.A.*;
- 10 44. Cherif—Feildel M, Heude Berthelin C, Adeline B, Rivière G, Favrel P, Kellner K. Molecular  
11 evolution and functional characterisation of insulin related peptides in molluscs: Contributions of  
12 *Crassostrea gigas* genomic and transcriptomic-wide screening. *Gen Comp Endocrinol.* Academic  
13 Press Inc.; 2019; doi: 10.1016/j.ygcen.2018.10.019.
- 14 45. Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al.. The oyster genome reveals stress adaptation  
15 and complexity of shell formation. *nature.com*.
- 16 46. Qu Z, Bendena WG, Tobe SS, Hui JHL. Juvenile hormone and sesquiterpenoids in arthropods:  
17 Biosynthesis, signaling, and role of MicroRNA. *J. Steroid Biochem. Mol. Biol.* Elsevier Ltd;
- 18 47. Tsang SSK, Law STS, Li C, Qu Z, Bendena WG, Tobe SS, et al.. Diversity of Insect  
19 Sesquiterpenoid Regulation. *Front. Genet.* *Frontiers Media S.A.*;
- 20 48. Chipman AD, Ferrier DEK, Brena C, Qu J, Hughes DST, Schröder R, et al.. The First Myriapod  
21 Genome Sequence Reveals Conservative Arthropod Gene Content and Genome Organisation in  
22 the Centipede *Strigamia maritima*. *PLoS Biol.* Public Library of Science; 2014; doi:

- 1 10.1371/journal.pbio.1002005.
- 2 49. Schenk S, Krauditsch C, Frühauf P, Gerner C, Raible F. Discovery of methylfarnesoate as the  
3 annelid brain hormone reveals an ancient role of sesquiterpenoids in reproduction. *Elife*. eLife  
4 Sciences Publications Ltd; 2016; doi: 10.7554/eLife.17126.
- 5 50. Nong W, Cao J, Li Y, Qu Z, Sun J, Swale T, et al.. Jellyfish genomes reveal distinct homeobox  
6 gene clusters and conservation of small RNA processing. *Nat Commun*. Nature Research; 2020;  
7 doi: 10.1038/s41467-020-16801-9.
- 8 51. Tobe SS, Bendena WG. The regulation of juvenile hormone production in arthropods.  
9 Functional and evolutionary perspectives. *Ann N Y Acad Sci*. New York Academy of Sciences;
- 10 52. Hui JHL, Bendena WG, Tobe SS. Future perspectives for research on the biosynthesis of  
11 juvenile hormones and related sesquiterpenoids in arthropod endocrinology and ecotoxicology.  
12 *Juv Horm Juvenoids Model Biol Eff Environ Fate*. 2013; doi: 10.1201/b14899.
- 13 53. Zhang T, Yuan D, Xie J, Lei Y, Li J, Fang G, et al.. Evolution of the Cholesterol Biosynthesis  
14 Pathway in Animals. *Mol Biol Evol*. Oxford University Press; 2019; doi: 10.1093/molbev/msz167.
- 15 54. Yang, Y, Cheng, W, Wu, X, Huang, S, Deng, Z, Zeng, X, et al. Prediction of the potential  
16 global distribution for *Biomphalaria straminea*, an intermediate host for *Schistosoma mansoni*.  
17 (2018) PLOS Negl Trop Dis 12(5):e0006548
- 18 55. Barton-Owen, T. B., Szabó, R., Somorjai, I. M. L., and Ferrier, D. E. K. (2018). A revised  
19 spiralian homeobox gene classification incorporating new polychaete transcriptomes reveals a  
20 diverse TALE class and a divergent hox gene. *Genome Biol. Evol.* 10, 2151–2167.  
21 doi:10.1093/gbe/evy144.
- 22 56. Butts, T., Holland, P. W. H., and Ferrier, D. E. K. (2008). The Urbilaterian Super-Hox cluster.



- 1 Trends Genet. 24, 259–262. doi:10.1016/j.tig.2007.09.006.
- 2 57. Castro, L. F. C., and Holland, P. W. H. (2003). Chromosomal mapping of ANTP class  
3 homeobox genes in amphioxus: piecing together ancestral genomes. *Evol. Dev.* 5, 459–465.  
4 doi:10.1046/J.1525-142X.2003.03052.X.
- 5 58. Chourrout, D., Delsuc, F., Chourrout, P., Edvardsen, R. B., Rentzsch, F., Renfer, E., et al.  
6 (2006). Minimal ProtoHox cluster inferred from bilaterian and cnidarian Hox complements. *Nat.*  
7 2006 4427103 442, 684–687. doi:10.1038/nature04863.
- 8 59. Edgar, R. C. (2004). MUSCLE: Multiple sequence alignment with high accuracy and high  
9 throughput. *Nucleic Acids Res.* 32, 1792–1797. doi:10.1093/nar/gkh340.
- 10 60. Ferrier, D. E. K. (2016). Evolution of Homeobox Gene Clusters in Animals: The Giga-Cluster  
11 and Primary vs. Secondary Clustering. *Front. Ecol. Evol.* 4, 1–13. doi:10.3389/fevo.2016.00036.
- 12 61. Hui, J. H. L., McDougall, C., Monteiro, A. S., Holland, P. W. H., Arendt, D., Balavoine, G.,  
13 et al. (2012). Extensive chordate and annelid macrosynteny reveals ancestral homeobox gene  
14 organization. *Mol. Biol. Evol.* 29, 157–165. doi:10.1093/molbev/msr175.
- 15 62. Li, Y., Nong, W., Baril, T., Yip, H. Y., Swale, T., Hayward, A., et al. (2020). Reconstruction  
16 of ancient homeobox gene linkages inferred from a new high-quality assembly of the Hong Kong  
17 oyster (*Magallana hongkongensis*) genome. *BMC Genomics* 21, 713. doi:10.1186/s12864-020-  
18 07027-6.
- 19 63. Mazza, M. E., Pang, K., Reitzel, A. M., Martindale, M. Q., and Finnerty, J. R. (2010). A  
20 conserved cluster of three PRD-class homeobox genes (homeobrain, rx and orthopedia) in the  
21 Cnidaria and Protostomia. *Evodevo* 1, 3. doi:10.1186/2041-9139-1-3.
- 22 64. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. (2015). IQ-TREE: A Fast

1 and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol.*  
2 *Evol.* 32, 268–274. doi:10.1093/molbev/msu300.

3 65. Paps, J., Xu, F., Zhang, G., and Holland, P. W. H. (2015). Reinforcing the egg-timer:  
4 Recruitment of novel Lophotrochozoa homeobox genes to early and late development in the  
5 Pacific oyster. *Genome Biol. Evol.* 7, 677–688. doi:10.1093/gbe/evv018.

6 66. Simakov, O., Marlétaz, F., Cho, S.-J., Edsinger-Gonzales, E., Havlak, P., Hellsten, U., et al.  
7 (2013). Insights into bilaterian evolution from three spiralian genomes. *Nature* 493, 526–531.  
8 doi:10.1038/nature11696.

9 67. Zhong, Y.-F., Butts, T., and Holland, P. W. H. (2008). HomeoDB: a database of homeobox  
10 gene diversity. *Evol. Dev.* 10, 516–518. doi:10.1111/J.1525-142X.2008.00266.X.

11 68. Zhong, Y., and Holland, P. W. H. (2011). HomeoDB2: functional expansion of a comparative  
12 homeobox gene database for evolutionary developmental biology. *Evol. Dev.* 13, 567–568.  
13 doi:10.1111/J.1525-142X.2011.00513.X.

14 69. Brooke, N.M., Garcia-Fernández, J. & Holland, P.W.H. (1998) *Nature* 392, 920-922

15 70. Osborne, P.W., Benoit, G., Laudet, V., Schubert, M., & Ferrier, D.E.K. (2009). Differential  
16 regulation of ParaHox genes by retinoic acid in the invertebrate chordate amphioxus  
17 (*Branchiostoma floridae*). *Developmental Biology* 327, 252-262.

18 71. Ikuta, T., et al., (2013). Identification of an intact ParaHox cluster with temporal colinearity  
19 but altered spatial colinearity in the hemichordate *Ptychodera flava*. *BMC Evolutionary Biology*  
20 13:129.

21 72. Zhang, H., et al. (2017). Lampreys, the jawless vertebrates, contain only two ParaHox gene  
22 clusters. *PNAS* 114(34), 9146-9151.

- 1 73. Irimia, M., Maeso, I., and Garcia-Fernandez, J. (2008). Convergent Evolution of Clustering  
2 of Iroquois Homeobox Genes across Metazoans. *Mol. Biol. Evol.* 25, 1521–1525.  
3 doi:10.1093/molbev/msn109.
- 4 74. Kerner, P., Ikmi, A., Coen, D., and Vervoort, M. (2009). Evolutionary history of the  
5 iroquois/Irx genes in metazoans. *BMC Evol. Biol.* 9, 1–14. doi:10.1186/1471-2148-9-74.
- 6 75. Takatori, N., Butts, T., Candiani, S., Pestarino, M., Ferrier, D.E.K., Saiga, H., and Holland,  
7 P.W.H. (2008). Comprehensive survey and classification of homeobox genes in the genome of  
8 amphioxus, *Branchiostoma floridae*. *Dev. Genes Evol.* 218, 579-590. Doi: 10.1007/s00427-008-  
9 0245-9.

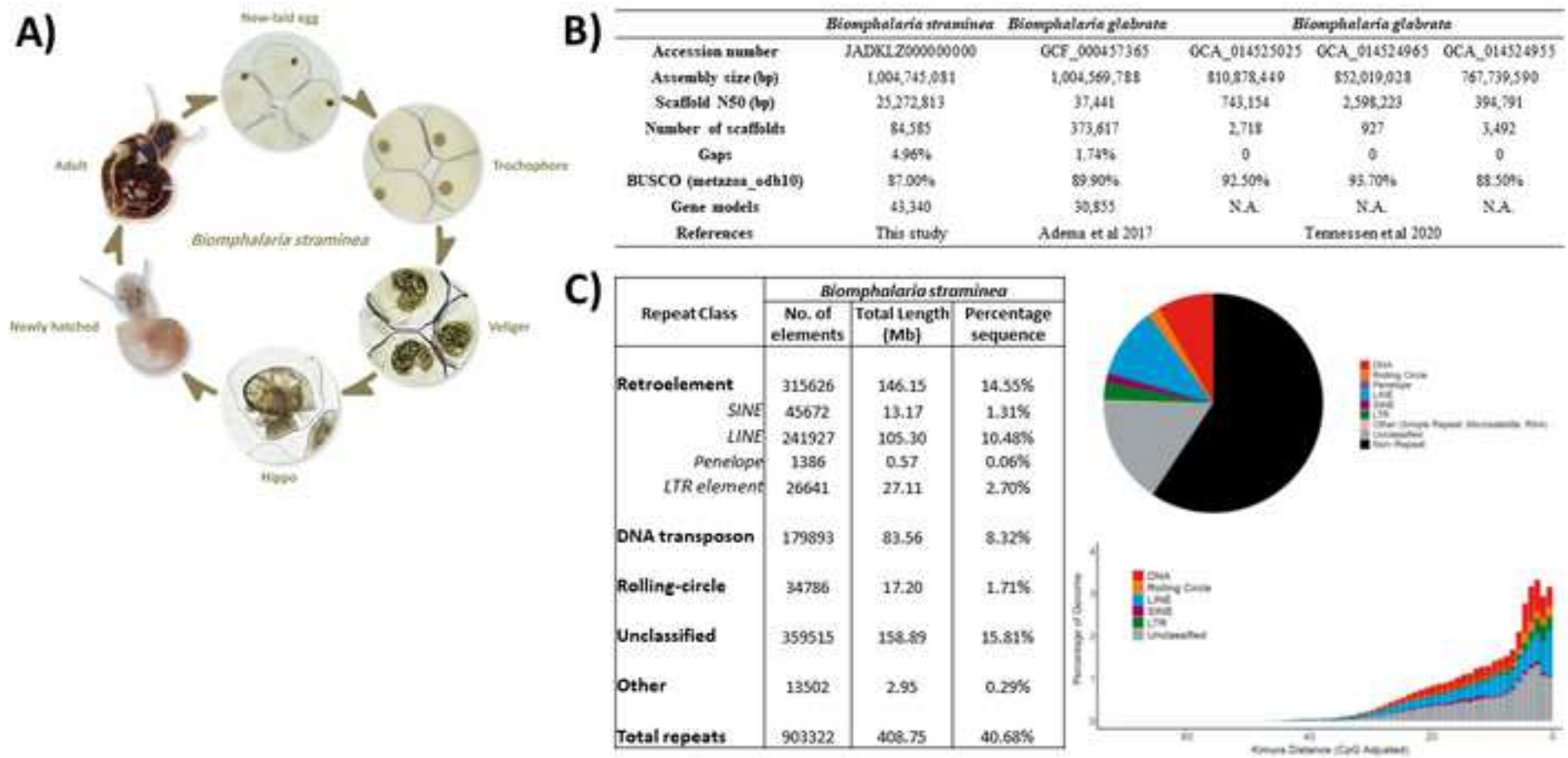


Fig 1

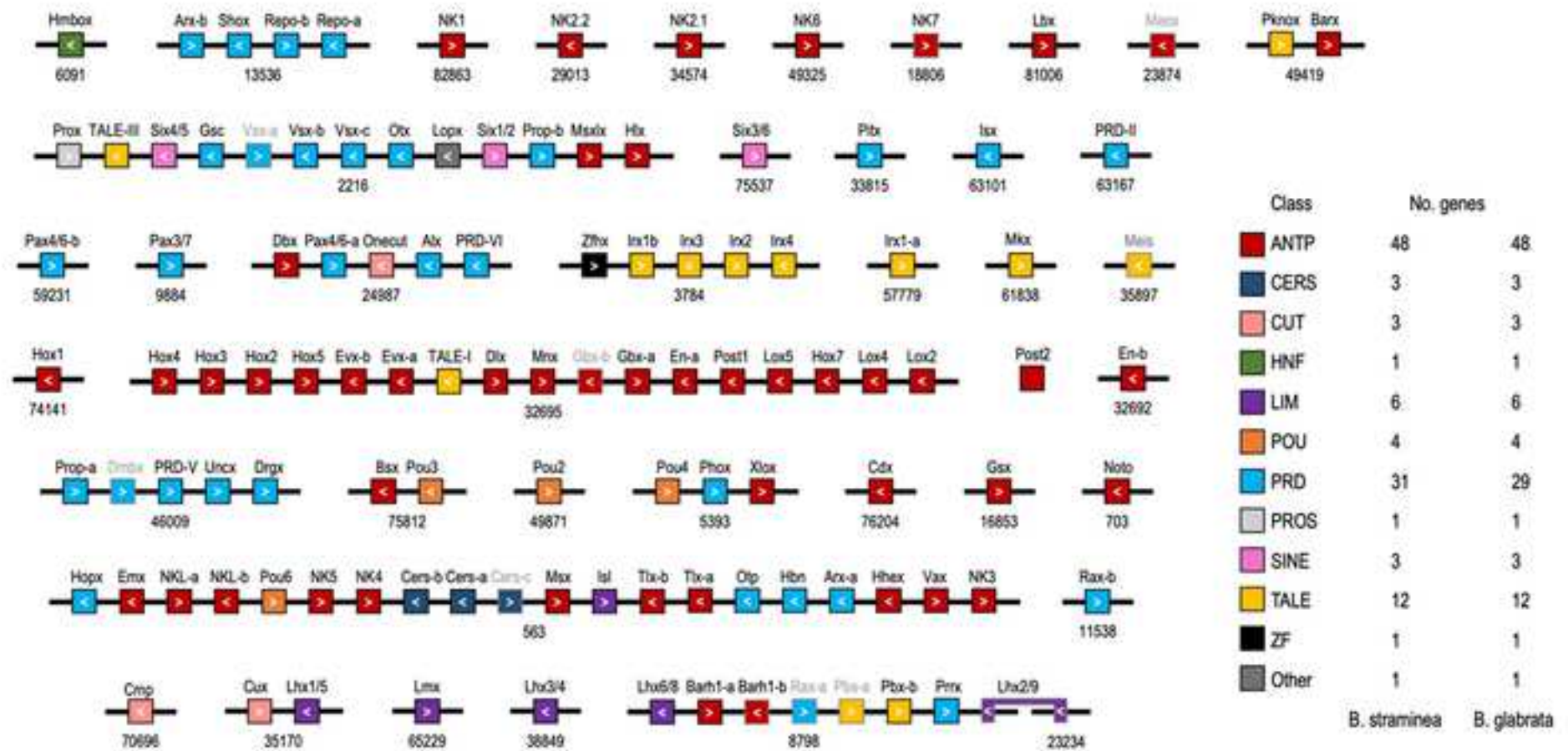


Fig 2

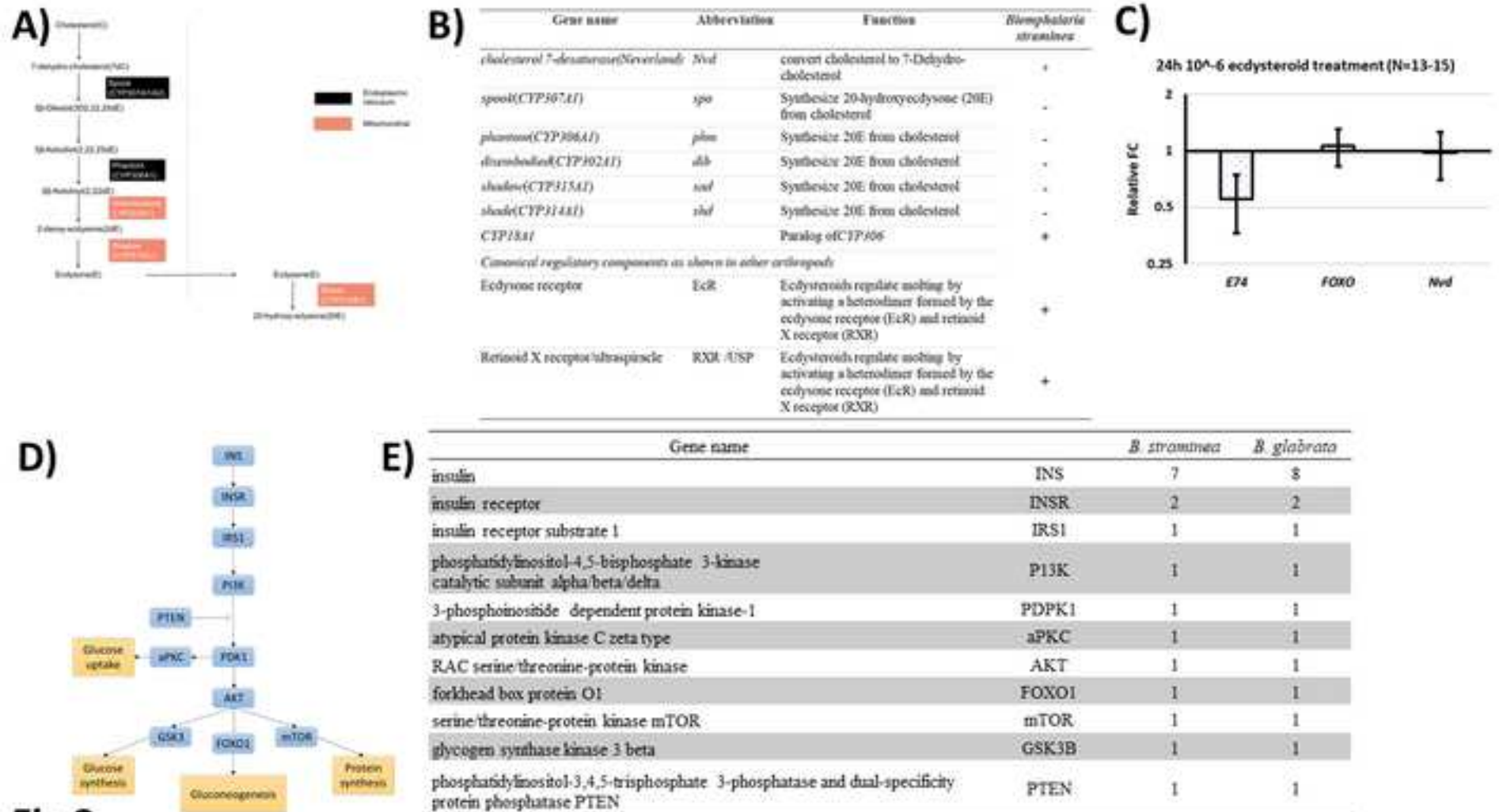
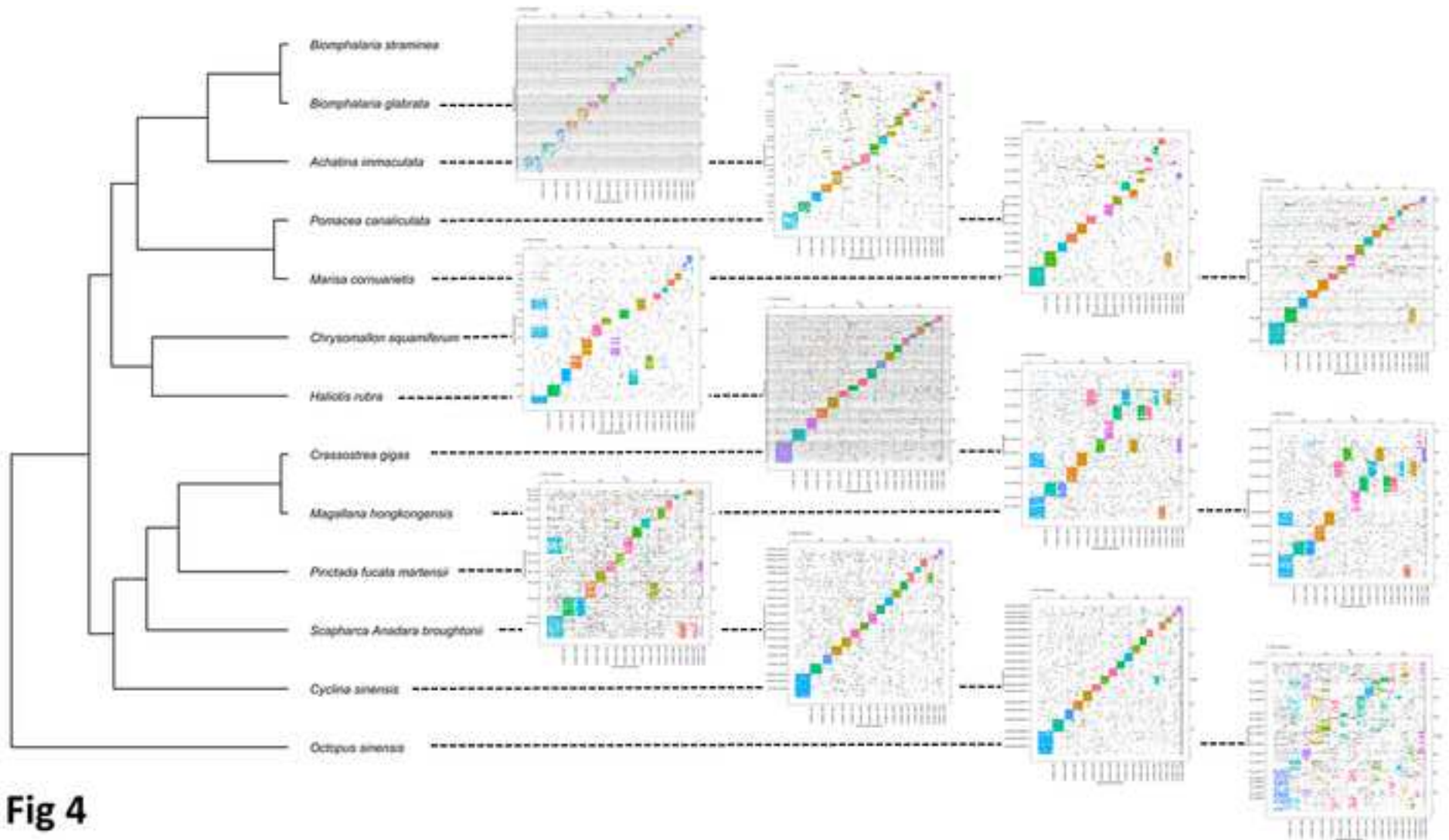
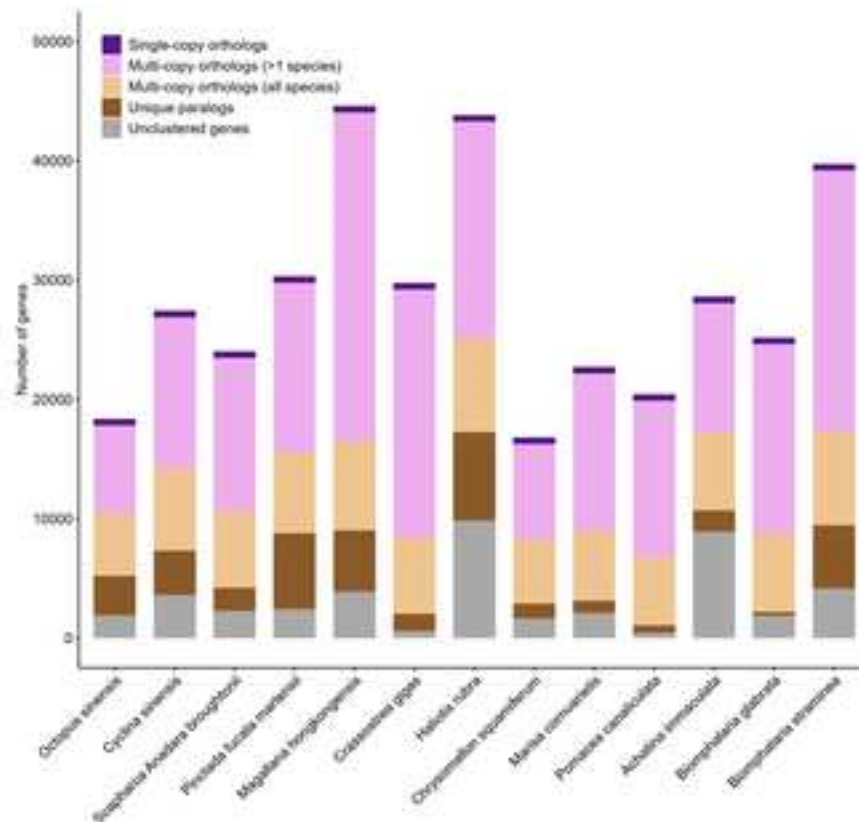


Fig 3

**Fig 4**

A)



B)

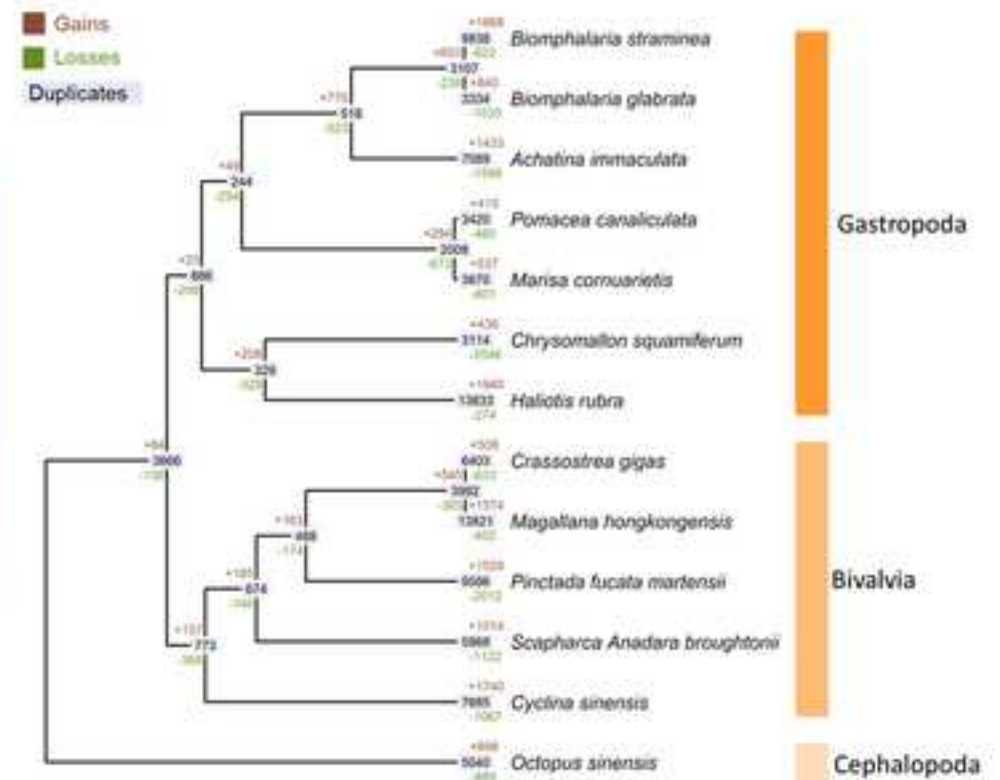
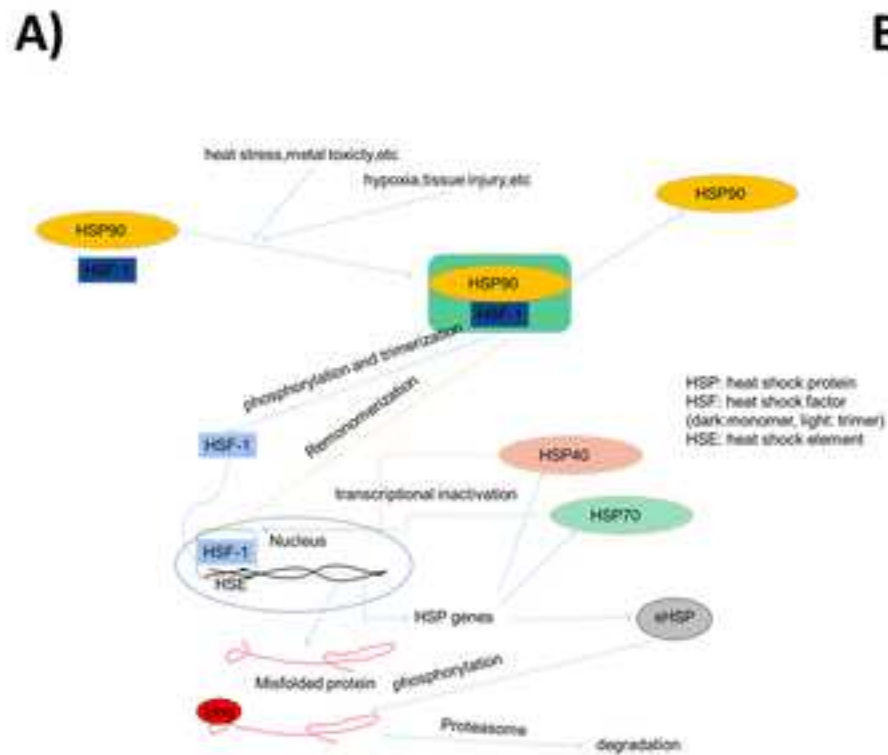


Fig 5





**B)**

Family	Species	HSF1*	HSF2	Dna/HSF40	HSF70	HSFPC/HSF90	HSF110
Gastropoda	Planorbidae <i>Biomphalaria straminea</i>	2	7	44	41	5	6
	Planorbidae <i>Biomphalaria glabrata</i>	2	5	32	12	2	1
	Hydrobia <i>Hydrobia ulvae</i>	3	8	51	24	5	2
	Ampullariidae <i>Lymnaea stagnalis</i>	2	11	36	15	4	2
	Ampullariidae <i>Lymnaea stagnalis</i>	2	9	37	13	3	1
	Ampullariidae <i>Melampus commersoni</i>	2	12	35	13	3	2
	Ampullariidae <i>Pomacea canaliculata</i>	2	13	36	12	3	2
	Ampullariidae <i>Pomacea maculata</i>	2	12	36	10	3	2
	Lymnaeidae <i>Radix auricularia</i>	2	9	28	24	3	2
	Pulmonata <i>Chrysomallon squaliformum</i>	2	10	36	16	3	1
Bivalvia	Achardidae <i>Achardoa fulva</i>	3	14	37	26	5	3
	Achardidae <i>Achardoa immaculata</i>	4	14	33	22	5	4
	Ostreida <i>Megallana hongkongensis</i>	2	20	52	123	5	2
	Ostreida <i>Crassostrea gigas</i>	2	12	56	141	5	2
	Ostreida <i>Crassostrea virginica</i>	3	18	56	141	4	2
	Ostreida <i>Saccostrea glomerata</i>	2	13	41	110	3	1
	Mytilidae <i>Bathymedulla planifrons</i>	2	7	43	103	3	2
	Mytilidae <i>Mytilus philippinensis</i>	2	14	40	87	3	2
	Pectinidae <i>Pinctada fucata</i>	2	12	35	89	3	3
	Pectinidae <i>P. f. martensii</i>	0	10	40	95	3	2
Pectinidae <i>Musculopecten yessoensis</i>	2	8	37	61	3	2	
Vanidae <i>Cyclina sinensis</i>	4	12	44	77	5	2	
Archoidea <i>Scapharca (Archidona) broughtonii</i>	2	11	42	81	4	2	
Cephalopoda	Octopodidae <i>Octopus bimaculoides</i>	2	8	32	9	4	3
	Architeuthidae <i>Architeuthis duxripud</i>	2	8	28	16	4	3
	Octopodidae <i>Octopus sinensis</i>	3	8	34	16	11	3
Annelida	Cirratia <i>Helobdella robusta</i>	2	13	37	10	3	2
	Polychaeta <i>Capitella teleta</i>	2	23	34	34	4	2
Pterobranchia	Cystoda <i>Leptococcus multiseptatus</i>	2	4	23	49	6	6
	Sessilata <i>Schistosoma mansoni</i>	1	11	23	6	3	1

Fig 6

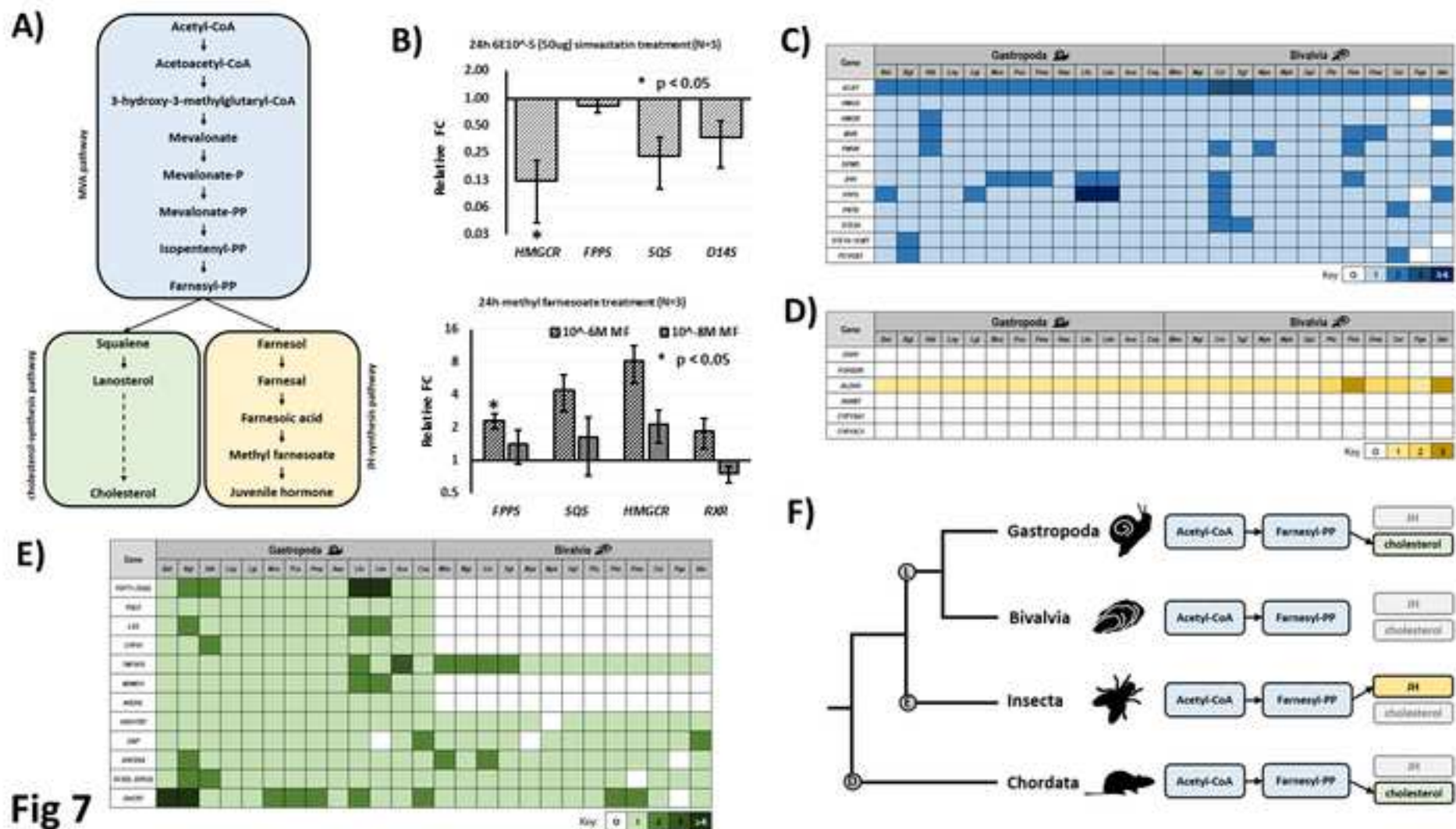

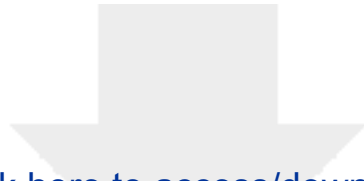


Fig 7



Click here to access/download  
**Supplementary Material**  
S1. Sequencing data.xls

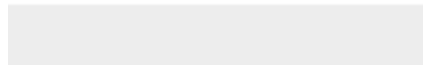





[Click here to access/download](#)


**Supplementary Material**

**S2a\_HboxSeqsSyntenyParaHox.xlsx**





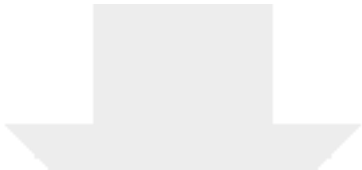
Click here to access/download  
**Supplementary Material**  
S2b\_B.glabrata\_Hboxes.jpg



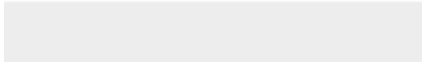



Click here to access/download  
**Supplementary Material**  
S3. Ecdysteroid.xlsx





Click here to access/download  
**Supplementary Material**  
S4. Insulin.xlsx





Click here to access/download  
**Supplementary Material**  
S5. Syntenypptx

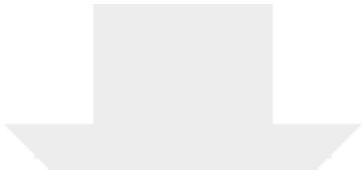




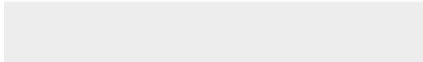



Click here to access/download  
**Supplementary Material**  
S6. Gene gain and loss.xlsx





Click here to access/download  
**Supplementary Material**  
S7. HSP.xlsx

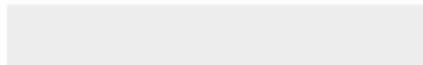




[Click here to access/download](#)

**Supplementary Material**

S8. Cholesterol and sesquiterpenoid.xlsx





Click here to access/download  
**Supplementary Material**  
S9. Phylogentic trees.pptx





Click here to access/download  
**Supplementary Material**  
S10. Tables.pptx

