

1 **The genome and transcriptome of the snail *Biomphalaria sudanica* s.l.: Immune gene**
2 **diversification and highly polymorphic genomic regions in an important African vector of**
3 ***Schistosoma mansoni***

4
5 Tom Pennance¹, Javier Calvelo², Jacob A. Tennessen³, Ryan Burd¹, Jared Cayton¹, Stephanie R.
6 Bollmann⁴, Michael S. Blouin⁴, Johann M. Spaan¹, Federico G Hoffmann⁵, George Ogara⁶,
7 Fredrick Rawago⁶, Kennedy Andiego⁶, Boaz Mulonga⁶, Meredith Odhiambo⁶, Eric S. Loker⁷,
8 Martina R. Laidemitt⁷, Lijun Lu⁷, Andrés Iriarte², Maurice Odiere⁵, and Michelle L. Steinauer¹

9
10 ¹ College of Osteopathic Medicine of the Pacific – Northwest, Western University of Health
11 Sciences, Lebanon OR, USA

12 ² Laboratorio Biología Computacional, Departamento de Desarrollo Biotecnológico, Instituto de
13 Higiene, Facultad de Medicina, Universidad de la República, Montevideo 11600, Uruguay

14 ³ Harvard T.H. Chan School of Public Health, Boston, MA USA

15 ⁴ Oregon State University, Corvallis, OR USA

16 ⁵ Department of Biochemistry, Molecular Biology, Entomology, and Plant Pathology, Mississippi
17 State University, Starkville, MS USA

18 ⁶ Centre for Global Health Research, Kenya Medical Research Institute (KEMRI), P. O. Box 1578-
19 40100, Kisumu, Kenya

20 ⁷ Department of Biology, Center for Evolutionary and Theoretical Immunology, Parasite Division
21 Museum of Southwestern Biology, University of New Mexico, Albuquerque, New Mexico 87131,
22 U.S.A.

23

24 **Corresponding authors:** Tom Pennance (tpennance@westernu.edu) and Michelle L. Steinauer
25 (msteinauer@westernu.edu)

26

27

28

29 **Abstract**

30 **Background:** Control and elimination of schistosomiasis is an arduous task, with current strategies
31 proving inadequate to break transmission. Exploration of genetic approaches to interrupt
32 *Schistosoma mansoni* transmission, the causative agent for human intestinal schistosomiasis in
33 sub-Saharan Africa and South America, has led to genomic research of the snail vector hosts of
34 the genus *Biomphalaria*. Few complete genomic resources exist, with African *Biomphalaria*
35 species being particularly underrepresented despite this being where the majority of *S. mansoni*
36 infections occur. Here we generate and annotate the first genome assembly of *Biomphalaria*
37 *sudanica* sensu lato, a species responsible for *S. mansoni* transmission in lake and marsh habitats
38 of the African Rift Valley. Supported by whole-genome diversity data among five inbred lines, we
39 describe orthologs of immune-relevant gene regions in the South American vector *B. glabrata* and
40 present a bioinformatic pipeline to identify candidate novel pathogen recognition receptors
41 (PRRs).

42 **Results:** *De novo* genome and transcriptome assembly of inbred *B. sudanica* originating from the
43 shoreline of Lake Victoria (Kisumu, Kenya) resulted in a haploid genome size of ~944.2 Mb (6732
44 fragments, N50=1.067 Mb), comprising 23,598 genes (BUSCO=93.6% complete). The *B.*
45 *sudanica* genome contains orthologues to all described immune genes/regions tied to protection
46 against *S. mansoni* in *B. glabrata*. The *B. sudanica* *PTC2* candidate immune genomic region
47 contained many PRR-like genes across a much wider genomic region than has been shown in *B.*
48 *glabrata*, as well as a large inversion between species. High levels of intra-species nucleotide
49 diversity were seen in *PTC2*, as well as in regions linked to *PTC1* and *RADres* orthologues.
50 Immune related and putative PRR gene families were significantly over-represented in the sub-set
51 of *B. sudanica* genes determined as hyperdiverse, including high extracellular diversity in
52 transmembrane genes, which could be under pathogen-mediated balancing selection. However, no
53 overall expansion in immunity related genes were seen in African compared to South American
54 lineages.

55 **Conclusions:** The *B. sudanica* genome and analyses presented here will facilitate future research
56 in vector immune defense mechanisms against pathogens. This genomic/transcriptomic resource
57 provides necessary data for the future development of molecular snail vector control/surveillance
58 tools, facilitating schistosome transmission interruption mechanisms in Africa.

59

60

61 **Keywords**

62 *Biomphalaria sudanica*, *Biomphalaria choanomphala*, schistosomiasis, snail vector, *de novo*
63 genome assembly, polymorphism, immunogenetics, gene family evolution, balancing selection,
64 pathogen recognition;

65

66

67 **Background**

68

69 Freshwater pulmonate snails of the genus *Biomphalaria* are intermediate hosts for a
70 diversity of trematode parasites but are most notorious for their role in the transmission of
71 *Schistosoma mansoni*, to humans. Schistosomiasis is a chronic and inflammatory disease with
72 devastating impacts to human health, which likely are underestimated due to schistosomiasis
73 associated morbidities and mortalities being attributed to non-communicable diseases for which
74 the symptoms are similar (1). Despite the widely recognized toll of schistosomiasis on human
75 health, there are few effective and implementable options for controlling transmission of the
76 parasite (2).

77 The need for novel interventions to interrupt *S. mansoni* transmission has spurred on
78 investigatory genomic research of the *Biomphalaria* snail vector hosts (3). Genomic resources for
79 *Biomphalaria* will facilitate the discovery of genetic resistance mechanisms to schistosome
80 infection, which could be manipulated to block transmission to snails, and thus humans (4).
81 Currently, three *Biomphalaria* species represent the extent of genome resources, two South
82 American species: *Biomphalaria glabrata* (5,6); *Biomphalaria straminea* (7); and one sub-
83 Saharan African species *Biomphalaria pfeifferi* (8). African *Biomphalaria* species are therefore
84 underrepresented in terms of complete genomic information, even though African species
85 contribute to the vast majority of global *S. mansoni* transmission since approximately 90% of
86 human infections occur in Africa. Thus, genome-wide analysis of African *Biomphalaria* species
87 would facilitate the development of genetic based vector control in areas where it is highly relevant
88 to transmission. As more *Biomphalaria* genomes become available, evolutionary analysis of
89 immunity of these major vectors will be possible, whilst putting it into the context of species
90 divergence across the African continent after being introduced from South America sometime
91 between 1.8 and 5 MYA (8–11).

92 The African species *B. sudanica* was originally described in 1870 from Djur and Rek
93 tributaries of the White Nile in the Bahr el Ghazal region of Southern Sudan (12). It is distributed
94 throughout the Nile Basin in marsh and lacustrine habitats in Uganda, Kenya, Sudan, Tanzania
95 and Ethiopia (13–18). *Biomphalaria sudanica* distribution in East Africa corresponds to the
96 geographic region where the genetic diversity of *S. mansoni* is the greatest (19,20). Most research
97 regarding *B. sudanica* has been focused on populations from Lake Victoria, where *S. mansoni*

98 remains highly endemic even following repeated and widespread mass drug administration of
99 schistosomiasis preventative chemotherapy (21–23). Although *B. sudanica* inhabits the marshy
100 fringes and nearshore shallow waters of Lake Victoria where human-freshwater contact takes place
101 (24), another snail vector of *S. mansoni*, described as *B. choanomphala* (25), occurs in deep water
102 habitats of Lake Victoria (13,26,27). The taxonomic status of these two species of Lake Victoria
103 snails is in question, as DNA divergence of mitochondrial genes (and the few nuclear genes that
104 have been sequenced) between these species suggests they may represent ecomorphs of a single
105 species (28–30). However, distinct morphologies, habitats, and schistosome susceptibility profiles
106 (31), make the distinction of these two forms critical in the context of a genome report. Thus, we
107 follow conventional use of the species name, or *Biomphalaria sudanica* sensu lato. Genomic data
108 of these species will facilitate future population genomic analyses aimed at better understanding
109 the relationship between these taxa. Snail based schistosomiasis control cannot even be imagined
110 without understanding these basics.

111 Experimental infections have shown that *B. sudanica* displays the greatest natural
112 resistance to schistosome infection relative to *B. choanomphala*, and another closely related
113 species, *B. pfeifferi* (27,31), and thus offers an excellent target for the discovery of immune
114 relevant genes in an African *Biomphalaria* species. While some immune genes can be
115 characterized by conserved domains as a result of positive selection (32), others such as those
116 involved in host-pathogen interactions are rapidly evolving under balancing selection due to the
117 simultaneous arms races occurring between the host and its pathogens (33). Indeed, immune loci
118 are among the most diverse in many genomes, including the classic example of the vertebrate
119 major histocompatibility complex (MHC) (34,35); human innate immunity genes (36); R genes in
120 plants (37,38) and more recently shown in immune genes of invertebrate organisms such as
121 *Caenorhabditis elegans* (39). Virtually nothing is known regarding the *B. sudanica* immune
122 defense except for what can be inferred from orthologous gene searching strategies related to
123 experimental work with the South American congener, *B. glabrata* (3). The identification of *B.*
124 *glabrata* loci associated with resistance to schistosomes is an active research field with approaches
125 such as Quantitative Trait Locus (QTL) analysis providing valuable new insights (6).

126 In this paper, we present the first description of the genome of *B. sudanica*. Our novel
127 annotated genome of *B. sudanica* 111 (Bs111), an inbred line maintained at Western University of
128 Health Sciences which originates from the Kisumu region (Kenya) of Lake Victoria, comprises

129 PacBio HiFi long-read DNA and RNA sequence data, as well as Illumina short-read RNA
130 sequence data. This combination of genomic and transcriptomic data provides a confident
131 annotation of functional gene boundaries, exon-intron structure, and isoforms for the
132 representative genome of this species. Here, we focus on identifying and describing gene regions
133 orthologous to those involved in immunity of the South American vector *B. glabrata* to *S. mansoni*,
134 such as the *PTC1* and *PTC2* genomic regions (40,41) and fibrinogen-related proteins (FREPs)
135 (42,43). We also used a new analysis pipeline to find novel pathogen recognition receptors (PRRs).
136 Following the hypothesis that PRRs are under balancing selection as with other highly
137 polymorphic immune loci, we searched in the most hyperdiverse genome regions through the
138 genomic comparison of five genetic lines of *B. sudanica* for signatures of candidate PRRs. This
139 allows us to identify immune related genes that do not maintain detectable sequence similarity
140 with known gene families and are not only specific to schistosome immunity. The description of
141 key features of the *B. sudanica* genome provides multiple exciting avenues for future research into
142 this important vector of *S. mansoni*.

143

144 **Results**

145

146 *Biomphalaria sudanica* 111 line genome assembly and nuclear genome annotation

147

148 The PacBio assembled *B. sudanica* (Bs111) haploid genome size is ~944.2 Mb, comprising
149 6732 contigs and scaffolds with an N50 of 1.067 Mb, and a mean sequencing coverage of ~23x
150 (Supplementary Table 1). The estimated size of the *B. sudanica* genome is somewhat larger than
151 those of *B. pfeifferi* (~771.8 Mb (8)) and *B. glabrata* (iM line: ~871.0 Mb (6)), but smaller than
152 that of *B. straminea* (~1,004.7 Mb (7)).

153

154 PacBio and Illumina RNA sequence data were obtained from Bs111 snails to aid in
155 annotation of the assembled *B. sudanica* genome. Pooled RNA was processed following a standard
156 PacBio IsoSeq procedure, which yielded 335.0 Gbases (N50 of ~115.5 Kbases) of long-read
157 transcript data. Of the original 6,945,781 circular consensus sequencing (ccs) reads, 3,798,283
158 (54.68%) passed the Q20 quality threshold determined in the longQC software (44), of which
159 1,708,667 (44.99%) were identified as potentially complete isoforms (i.e. they bear both 5' and 3'
adapter sequences) using Lima (github.com/PacificBiosciences/barcoding) (Supplementary Table

160 1). To supplement the RNA transcript long-read data, Illumina paired-end 150 short-read sequence
161 data yielded 45,125,478 paired reads of which 44,008,723 (97.53%) passed the trimming process
162 conducted in Trimomatic (45) (Supplementary Table 1). Overall mapping rate of long-read
163 transcripts using minimap2 (46,47) and Illumina RNA short-reads using STAR (48) to the
164 assembled genome was close to 100%. Transcript characterization using StringTie2 v2.2.1 (49)
165 identified 25,847 individual genes, of which 23,598 in TransDecoder.Predict v5.5 (50) had an
166 assigned open reading frame (ORF) (Supplementary File 1 and Supplementary File 2).
167 InterProScan v5.56-89.0 (51) identified at least one protein domain signature on 19,945 of the
168 23,598 genes (Supplementary File 3). BUSCO (52) completeness analysis shows that the latter set
169 (23,598 genes) represents a close to complete genome annotation relative to mollusca_odb10 (of
170 5295 BUSCO groups, 93.6% complete, 1.4% fragments and 5.0% missing).

171 The *B. sudanica* genome contains orthologues to at least 18 candidate immune loci of *B.*
172 *glabrata* that function in protection against *S. mansoni* (Supplementary Table 2). These include
173 key genes/gene clusters coding for FREP2 and FREP3, *B. glabrata toll-like receptor (BgTLR)*,
174 *Phox*, Guadeloupe Resistance Complex 1 genes (*GRC*) – referred to from here as the polymorphic
175 transmembrane cluster 1 (*PTC1*), polymorphic transmembrane cluster 2 (*PTC2*), *RADres*, heat
176 shock protein 90 (*HSP90*), *Granulin (GRN)*, *BgTEP*, Catalase (*cat*), *Biomphalysin*, *Glabralysin*,
177 OPM-04 (Knight marker), superoxide dismutase 1 (*sod1*), *Peroxiredoxin (prx4)*, qRS-5.1 and
178 qRS-2.1 (Supplementary Table 2).

179 A total of 919 tRNA and 107 rRNA genes were predicted in the *B. sudanica* nuclear
180 genome (Supplementary Table 3, Supplementary Table 4 and Supplementary File 1). The number
181 of tRNA genes identified in the *B. sudanica* genome is comparably higher than that observed in
182 the genomes of *B. pfeifferi* (n=514 (8)) and *B. glabrata* (n=510 (6)). As is the case for *B. pfeifferi*,
183 one selenocysteinyl tRNA (tRNA-SeC) is present in the genome of *B. sudanica* (Supplementary
184 Table 3 and Supplementary File 1), meaning this species is capable of synthesizing selenocysteine
185 containing polypeptides, or selenoproteins (53,54). The tRNA-Sec gene has not been identified in
186 *B. glabrata* (8). Overall, fewer rRNA genes were predicted in this genome assembly of *B. sudanica*
187 compared to that of *B. pfeifferi* (107 and 757, respectively), which could be a result of some rRNA
188 genes being misassembled in *B. sudanica*.

189 As with other *Biomphalaria*, repetitive elements composed a large proportion of the *B.*
190 *sudanica* genome (40.3%) (Supplementary Figure 1). About 87% of protein coding genes overlap

191 with at least one annotated repeated element in their gene model. The overlap is primarily within
192 introns and untranslated regions (UTRs); however, 1576 genes have repeat elements within their
193 predicted coding sequence (CDS) (Supplementary Table 5). The repeat regions largely comprise
194 unknown repeat elements, in addition to an abundance of unclassified long interspersed nuclear
195 elements (LINE), LINE/retrotransposable element Bovine B (RTE-BovB) and unclassified DNA
196 transposons (Supplementary Figure 1, Supplementary Table 6).

197

198 *Mitochondrial genome annotation and trimming processes*

199

200 The mitochondrial genome comprises the same gene content (13 genes, 3 rRNA, 22 tRNA),
201 and synteny as its congeners (55) (Supplementary Table 7). The mitochondrial genomes of
202 gastropods are unique in the fact that they have acquired transcriptional processes during their
203 evolutionary history that are not often observed in vertebrates (56). Therefore, while an annotated
204 mitochondrial genome of *B. sudanica* has been previously published (55), our novel contribution
205 here is to validate gene boundaries and explore the transcription processes using long-read
206 transcriptomic data.

207 Raw q20 PacBio IsoSeq reads were mapped to the mitochondrial genome with minimap2
208 (46,47). To explore the trimming process of the primary mitochondrial transcript, the intermediary
209 pre-mRNA mitochondrial transcripts, i.e., that cover multiple features within the mitochondrial
210 genome, were recovered and counted (Supplementary Figure 2). In line with the tRNA punctuation
211 model (57), we established that pre-mRNAs of the *B. sudanica* mitochondrial genome are trimmed
212 at the tRNA genes (with the potential exception of *atp6/atp8*), with the minus strand being
213 processed 3'-to-5' while the plus strand shows an odd mixture of both directions.

214 In *B. sudanica*, there are three mitochondrial gene re-arrangements in comparison to a
215 typical animal mitochondrial genome that affect the transcription processes (see (58)). First, it is
216 typical for *atp6* and *atp8* to be adjacent and remain together in the mature mRNA; however, in
217 *Biomphalaria*, including *B. sudanica*, these genes are separated by a tRNA gene (*trnN*, see
218 Fragment 5; Supplementary Figure 2 and (55)). In the case of *atp6/atp8*, monocistronic
219 transcription of each of these genes following the cleavage of *trnN* is expected. Despite this, only
220 a few monocistronic transcripts of each *atp6* (n=16) and *atp8* (n=1) were identified from the raw
221 RNA reads, in contrast with the far more abundant untrimmed intermediaries (n=145;

222 Supplementary Figure 2). Considering that the ancestral condition is the translation of both
223 proteins from a bicistronic mRNA, it is a tempting hypothesis that this is still the case in
224 *Biomphalaria* despite the extra trimming points.

225 Second, it appears that, *nad4l* is either a non-functional pseudogene in the *B. sudanica*
226 mitochondrial genome or is only expressed in very low levels. The adjacency and bicistronic
227 transcription of *nad4/nad4l*, is well conserved in invertebrate and vertebrate mitochondrial
228 genomes, yet in many molluscan lineages, including *B. sudanica*, these genes are nonadjacent
229 (Supplementary Figure 2) (56). Additionally, the transcript data demonstrates that *nad4* terminates
230 on an abbreviated stop codon (T--) as was experimentally supported in its South American sister
231 species *B. glabrata* (59), and was abundantly represented in the transcriptome (n=83
232 monocistronic transcripts from Fragment 3; Supplementary Figure 2). On the other hand, *nad4l*
233 transcripts were rare, represented by only four intermediary reads that were attached to transcripts
234 for neighboring gene Cytochrome B (*cob*) (Fragment 1, Supplementary Figure 2).

235 Lastly, *nad6/nad5/nad1* are found adjacent to one another and seem to be translated as a
236 polycistronic mRNA (Fragment 1, Supplementary Figure 2), since their gene boundaries predicted
237 by MITOS2 (60) overlap; there are no clear cuts in the read coverage; and pure monocistronic
238 reads were only recovered in small numbers for *nad5* (n=16) and *nad1* (n=2). Given the lack of a
239 clear trimming point in this region, the observed monocistronic and bicistronic reads (e.g.
240 *nad6/nad5* or *nad5/nad1*) are likely the product of partial RNA degradation, implying that the
241 three genes are translated into proteins from the same mRNA molecule.

242

243 *Location signals and transmembrane domains*

244

245 Location signals for exportation, i.e. the signal peptides or mitochondrial targeting
246 peptides, were identified with SignalP v6.0 (61) and TargetP v2.0 (62) in 3,339 genes (5,016
247 isoforms) and 69 genes (111 isoforms), respectively (Supplementary Table 8). Transmembrane
248 domains were predicted in 4,922 genes (8,728 isoforms), and the location signals analysis suggests
249 that 835 of these were firmly anchored to the plasma membrane, organellar membranes or vesicles,
250 and seven to the mitochondria (Supplementary File 4).

251 An additional 82 proteins (146 isoforms) that showed no location signals in SignalP and
252 TargetP were identified by SecretomeP (63) as potentially secreted through an alternative pathway.

253 However, since 34 of these had at least one transmembrane domain, we suspect many of these are
254 false negatives for either the signal or the mitochondrial targeting peptide. Furthermore four genes
255 have signal peptides predicted in some but not all their isoforms (genes BSUD.7093, BSUD.10729,
256 BSUD.12693 and BSUD.24440). This might represent cases of functional isoforms generated
257 through alternative splicing that have different locations in the cell, as observed in other species
258 (64).

259 Two identical secreted proteins worth pointing out, BSUD.4529 (contig 217) and
260 BSUD.14556 (contig 559), were identified as orthologs to the precursor protein of peptide P12 in
261 *B. glabrata* (BGLB027975), which has been shown to trigger behavior modifications in *S. mansoni*
262 miracidia, and thus is potentially an attractant (65,66). Compared to the *B. glabrata* P12, the
263 ortholog in *B. sudanica* contains a non-synonymous change within the 13 aa region, changing the
264 5th amino acid from Glycine to Valine (DITSVLDPEVADD).

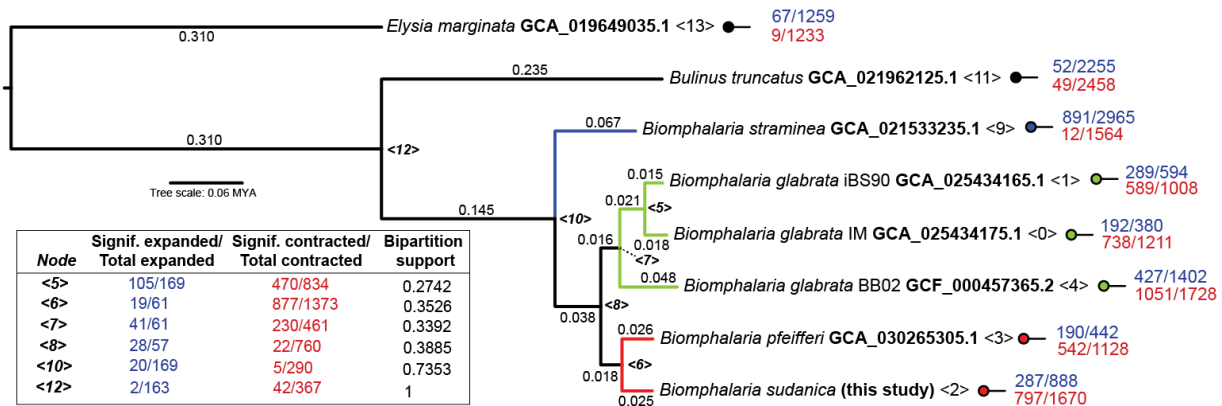
265

266 *Gene family evolution in Biomphalaria species*

267

268 The evolutionary dynamics of *Biomphalaria* protein families among *B. sudanica*, *B. glabrata*, *B.*
269 *pfeifferi*, *B. straminea*, *Bulinus truncatus* and *Elysia marginata* were estimated by first identifying
270 orthology in Phylogenetic Hierarchical Orthogroups (HOG) with Orthofinder v2.5.4 (67).
271 Orthofinder clustered ~87% of identified genes into 29,664 orthogroups with representatives in at
272 least 2 species. Orthogroups were distributed across 31,723 HOGs (Supplementary Table 9).
273 HOGs were taken as an approximate estimation of protein families found on these genomes, and
274 significant changes in their size estimated with CAFE 5 (68). The phylogeny utilized for the
275 estimation was the species tree generated by Orthofinder for the HOG definition, including 3541
276 single copy orthogroups with representatives in all the genomes, and the root calibrated to be 20
277 Million Years based on the appearance of *Bulinus* in the fossil record 19-20 MYA (69) (see Figure
278 1). Gene Ontology (GO) terms were assigned to each HOG with eggNOG-mapper v2.0 (70,71)
279 assuming that a GO assigned to one member applied to the whole HOG (Supplementary Table
280 10), significantly enriched GO terms among the expanded and contracted HOGs on each node
281 were identified with topGO (72) (Supplementary Table 11) and then summarized with REVIGO
282 (73) to allow further interpretation (Supplementary Table 12 Supplementary File 5).

283 Significant gene family expansions/contractions were identified throughout the analyzed
 284 phylogeny. Lineage-specific trends were identified within the *Biomphalaria* genus, in comparison
 285 to outgroups, between South American and African *Biomphalaria* species, within the African
 286 lineages (*B. sudanica* and *B. pfeifferi*), and between lines of *B. glabrata* (Figure 1, Supplementary
 287 Table 12 and Supplementary File 5). Expansions in immune related gene families were detected
 288 in the common ancestor of all *Biomphalaria* species relative to the previous node representing the
 289 split from *Bulinus truncatus*, including those involved in acute inflammatory response, glomerular
 290 filtration and regulation of cell adhesion (see node 10, Figure 1 and Supplementary File 5 and
 291 Supplementary Table 12).
 292



293
 294 Figure 1. Species tree generated in Orthofinder using the Species Tree of All Genes (STAG)
 295 algorithm (67,74). Root is time calibrated to be 20 Million Years Ago based on appearance of
 296 *Bulinus* in the fossil record (69). Node support values represent the bipartition proportions in each
 297 of the individual species tree estimates. Branch lengths represent the average number of
 298 substitutions per site across all the individual trees inferred from each gene family. The number of
 299 (significant/total) gene families expanded (blue) and contracted (red) in the ancestral populations
 300 of the *Biomphalaria* species, and outgroups *Elysia marginata* and *Bulinus truncatus* as determined
 301 in CAFE 5 (68) are shown for each internal and terminal node (<0> to <13>).
 302

303 In the branch leading to the common ancestor of the African species, *B. sudanica* and *B.*
 304 *pfeifferi*, we identified substantially more gene family contractions than expansions (see node 6,
 305 Figure 1). Within this lineage we found the expansion or contraction of gene families associated
 306 with the circulatory system (e.g. GO:0001525, GO:0001987, GO:0007512 and GO:0061337),

307 compound transport (e.g. GO:0098660, GO:0006811, GO:0035459 and GO:0006855), protein
308 maturation (e.g. GO:0006508, GO:0036211 and GO:0016579), metabolic pathways and regulation
309 (e.g. GO:1901293, GO:0006210, GO:0006212 and GO:0006164), development and growth (e.g.
310 GO:0050767, GO:0001558 and GO:0036342), and protection towards environmental stressors
311 including chemical stressors (e.g. GO:0009620, GO:0009636, GO:0010033 and GO:0010035)
312 (Supplementary Table 12 for full details). Several genes associated with chemotaxis (e.g.
313 GO:0006935, GO:0009410 and GO:0071310), which are often associated with immunity, were
314 contracted in the African *Biomphalaria* (Supplementary Table 12). Regarding the *B. sudanica*
315 lineage, gene families belonging to 287 GO terms were significantly expanded, including genes
316 associated with the regulation of hippo signaling (See Node 2, Group 4 in Supplementary File 5
317 and Supplementary Table 12) that account for multiple immune response processes. However,
318 many GO terms associated with the immune response were also identified in contracted gene
319 families in *B. sudanica*, including defense responses and regulation of immune system processes
320 suggesting that a complex evolutionary trend took place in this species.

321

322 *Identification and phylogeny of variable immunoglobulin and lectin domain-containing molecules*
323 *(VIgGs): FREPs and CREPs*

324

325 Genes putatively belonging to FREP, C-type lectin-related protein (CREP) or galectin-
326 related protein (GREP) families were predicted based on the presence of a secretion signal and in
327 conjunction with either fibrinogen (FBD), C-type lectin, or Galectin as predicted by InterProScan
328 v5.56-89.0 (51) (Supplementary File 3) and hmmsearch v3.3.2 (75) searches with custom IgSF
329 profiles (as described in (43)).

330 Following this selection pipeline, 246 genes distributed among 140 HOGs were determined
331 to have key domains that made them FREP, CREP or GREP candidates (Supplementary Table
332 13). Upon cross-checking our candidate protein domains, none of the seven initial GREP
333 candidates (i.e. Galectin domain-containing proteins) were identified to contain putative IgSF or
334 other immunoglobulin domains (Supplementary Table 13).

335 Unlike the FBD-containing proteins (i.e. candidate FREPs), proteins bearing C-type lectin
336 domains showed a considerable structural diversity (i.e. displaying several domains non-related
337 with the CREP family). Some of these genes may participate in the snail immune response, as they

338 often had signatures compatible with IgSF domains that were identified by InterPro (accessions:
339 IPR003598, IPR003599, IPR007110, IPR013098, IPR013783 and IPR036179), but were not
340 determined as members of the CREP family following previously described criteria (76).
341 Considering the hallmarks of a CREP or FREP gene (signal peptide present plus one or two IgSF
342 domains), a total of 10 potential CREP and 57 FREP (Supplementary Table 14) genes were
343 identified (Table 1 and Supplementary Table 13).

344

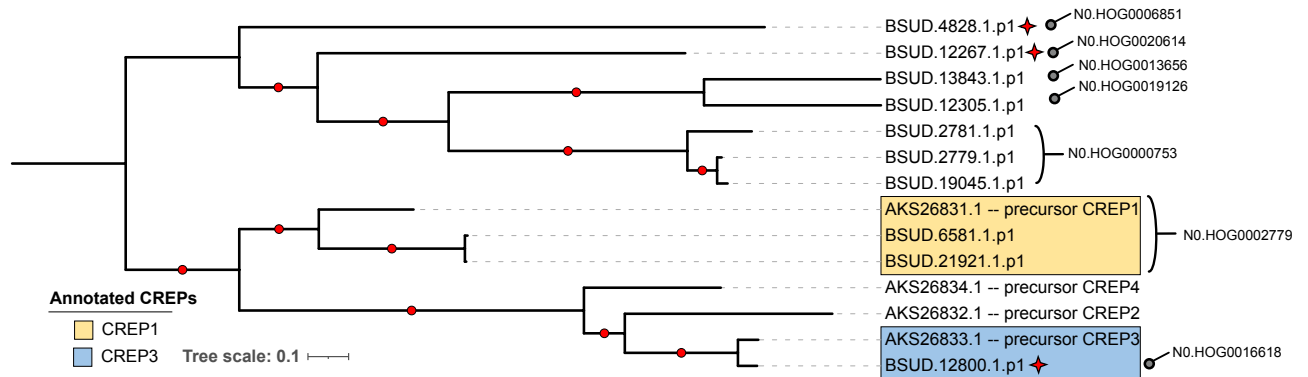
345 **Table 1.** Summary of variable immunoglobulin and lectin domain-containing molecules (VIgLs),
346 fibrinogen-related proteins (FREPs), C-type lectin-related protein (CREPs) and galectin-related
347 proteins (GREPs) identified in *Biomphalaria sudanica* (this study), *B. glabrata* and *B. pfeifferi* (8).
348 No GREPs were identified in either of the African *Biomphalaria* species *B. sudanica* or *B. pfeifferi*.
349 *One FREP in *B. sudanica* contained only a partial fibrinogen domain, and is considered truncated
350 (see BSUD.19120.1). Summaries of *B. sudanica* FREP and CREP gene compositions can be found
351 in Supplementary Table 14.

Complete VIgLs	<i>B. sudanica</i>	<i>B. glabrata</i>	<i>B. pfeifferi</i>
Fibrinogen-related proteins (FREPs)	57*	39	55
C-type lectin-related proteins (CREPs)	10	4	11
Galectin-related proteins (GREPs)	0	1	0

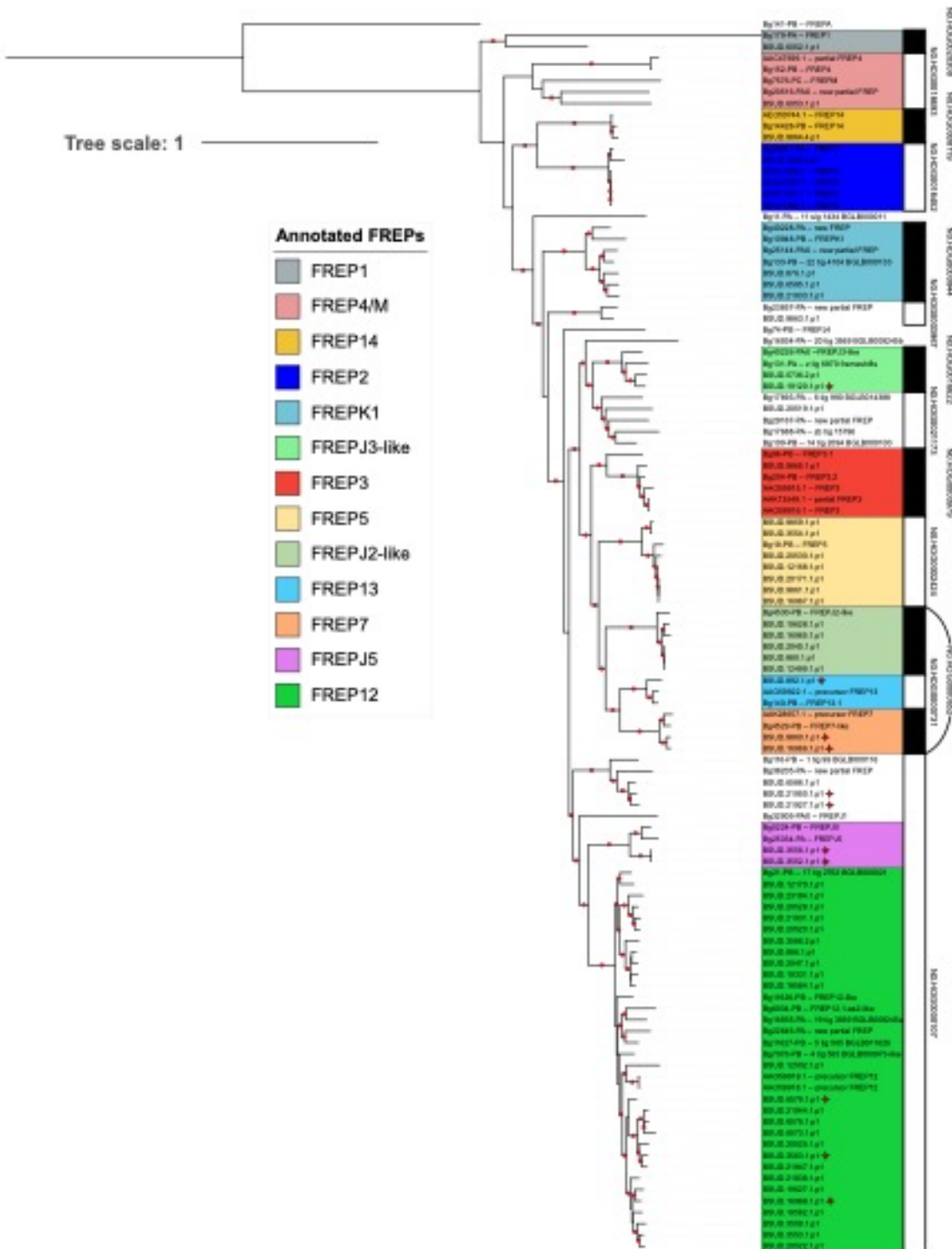
352

353 The selected full CREP and FREP genes were aligned with a set of reference sequences
354 from *B. glabrata* (Supplementary Table 15), their phylogeny relations estimated by maximum
355 likelihood and annotated based on their position in the phylogeny (Figure 2 and Figure 3). Thirty-
356 nine of these reference sequences were reported and curated previously (43), while 17 were
357 annotated as members of these gene families (FREP/CREP) and available on the National Center
358 for Biotechnology Information (NCBI) (Supplementary Table 15). The CREP family is divided
359 into two main monophyletic groups, one of which includes all the reference sequences (Figure 2).
360 In this group containing reference sequences, the HOGs associated with CREPs identified for *B.*
361 *sudanica* were N0.HOG0002779 and N0.HOG0016618, which are closely related with the
362 references CREP1 and CREP3, respectively. The other monophyletic group, which comprise five

363 HOGs, is formed exclusively by divergent *B. sudanica* CREPs, possibly representing new subtypes
364 within the CREP family.
365
366



367 **Figure 2.** Maximum likelihood tree of C-type lectin-related proteins (CREPs) identified from
368 *Biomphalaria sudanica* in the current study (see Supplementary Table 14) and four CREPs
369 identified previously from *B. glabrata* (see Supplementary Table 15) organized within hierarchical
370 orthogroups (HOGs) as determined by Orthofinder (67). Branch lengths represent the number of
371 substitutions per site. Nodes with bootstrap values >75 (estimated with 1000 replicates of non-
372 parametric bootstrap) are signified by a red dot on the branch before bipartition. Red stars indicate
373 three CREPs with unusual features, namely that they include weak hits for secondary
374 immunoglobulin domains, which may overlap with C-lectin domains as well as containing a large
375 interdomain region (see Supplementary Table 14).



376
 377 **Figure 3.** Maximum likelihood tree of the 57 fibrinogen-related proteins (FREPs) identified from
 378 *Biomphalaria sudanica* in the current study (see Supplementary Table 14) amongst reference
 379 sequences for FREPs identified previously from *B. glabrata* (see Supplementary Table 15)
 380 organized within hierarchical orthogroups (HOGs) as determined by Orthofinder (67). Branch

381 lengths represent the number of substitutions per site. Nodes with bootstrap values >75 (estimated
382 with 1000 replicates of non-parametric bootstrap) are signified by a red dot on the branch before
383 bipartition. Red stars indicate FREPs with unusual features according to our annotation
384 summarized in Supplementary Table 14, including those containing weak hits for additional
385 immunoglobulin (IgSF) domains (e.g. BSUD.16968), IgSF rearrangements (e.g. BSUD.21927) or
386 containing partial fibrinogen domains (e.g. BSUD.19120).

387

388 FREP genes were grouped in thirteen HOGs, eleven of which are closely related to
389 classified reference sequences (Figure 3). Thus, most FREPs in *B. sudanica* could be classified as
390 a known subtype within the family, except for FREP genes belonging to N0.HOG0000107,
391 N0.HOG0000652, N0.HOG0021173 and N0.HOG0003967. Genes from the HOG
392 N0.HOG0000107 is by far the largest FREP subgroup in *B. sudanica* identified in this work, and
393 it is subdivided in three monophyletic lineages, the largest lineage comprising 25 genes is
394 associated with several reference sequences identified as FREP12. The sister clade to the FREP12
395 group is the smaller monophyletic FREPJ5 group and basal to these are a group of three genes
396 from *B. sudanica* and two unclassified *B. glabrata* reference sequences, representing an
397 undescribed FREP lineage. Another HOG containing FREPJ2 and FREP7 (N0.HOG0000652) is
398 remarkable since it is paraphyletic (Figure 3). HOGs N0.HOG0021173 and N0.HOG0003967,
399 each comprise a single FREP gene of *B. sudanica*, and were not associated with FREP genes
400 previously assigned to a subgroup within the family. Lastly, HOG N0.HOG0018693 possesses
401 reference genes annotated for FREP4 and FREPM, yet unlike N0.HOG0000107 or
402 N0.HOG0000652, they do not form clear monophyletic groups with strong bootstrap support.

403 When comparing the numbers of CREP and FREP genes assigned to HOGs across the eight
404 analyzed genomes (*Biomphalaria* spp., *Bulinus truncatus*, *Elysia marginata*), it is evident that
405 FREPs are not only more numerous, but also their HOGs, at least for *B. sudanica*, contain more
406 genes that did not meet our selection criteria (i.e. containing well supported functional domains,
407 see Methods) to be considered for phylogenetic comparison (Supplementary Table 16). For
408 example, 12 of the 42 candidate genes initially identified in N0.HOG0000107 (FREPJ5/FREP12)
409 were rejected because of not containing complete FREP characters, suggesting these are not typical
410 FREP genes. Secondly, the number of FREP genes assigned to each HOG varies considerably
411 across *B. glabrata* and the African *Biomphalaria* lineages. For instance, N0.HOG0000107 is

412 relatively abundant in the African species *B. sudanica* (n=30) and *B. pfeifferi* (n=30), but this
 413 number varies considerably between *B. glabrata* lines iM (n=6), BB02 (n=0), and iBs90 (n=19)
 414 (Supplementary Table 16). In addition, only a single copy of a FREP13 classified gene
 415 (N0.HOG0000731) is present in all *Biomphalaria* genomes except for in *B. glabrata* iBS90 where
 416 it is absent, and in *B. glabrata* BB02 where it comprises 27 different genes. A full study of both
 417 CREP and FREP diversity and evolution is beyond the scope of this study, but these differences
 418 suggest that the FREP genes are under a complex dynamic of duplication and replacement, while
 419 the CREPs are more stable, at least among the laboratory strains for which there are genomes
 420 available.

421

422 *Intraspecific diversity of the B. sudanica genome: comparing four inbred lines*

423

424 In addition to the Bs111 reference genome, we conducted Illumina whole genome
 425 sequencing of four other inbred lines of *B. sudanica* and aligned them to the reference (Table 2).
 426 The line Bs163 was the most divergent, with an average genomic divergence value of 0.44% from
 427 the four other inbred lines (Bs111, Bs110, BsKEMRI and Bs5-2). The Bs163 line is also the most
 428 resistant to *S. mansoni* infection in the laboratory setting (77). Among the remaining four, pairwise
 429 divergences were similar and ranged between 0.28% and 0.41% with one exception: Bs110 and
 430 Bs111 were only 0.16% divergent showing high similarity. Among all five lines, median
 431 nucleotide diversity was 0.32% (95% CI = 0.06-0.83%). The heterozygosity per inbred line ranges
 432 from 0.11% to 0.21% (Table 2).

433

434 **Table 2.** Genome statistics from Illumina sequencing (paired end 150 bp) for four *Biomphalaria*
 435 *sudanica* inbred lines (Bs110, Bs163, Bs5-2 and BsKEMRI) aligned to the ~944.2 Mb
 436 *Biomphalaria sudanica* 111 reference genome.

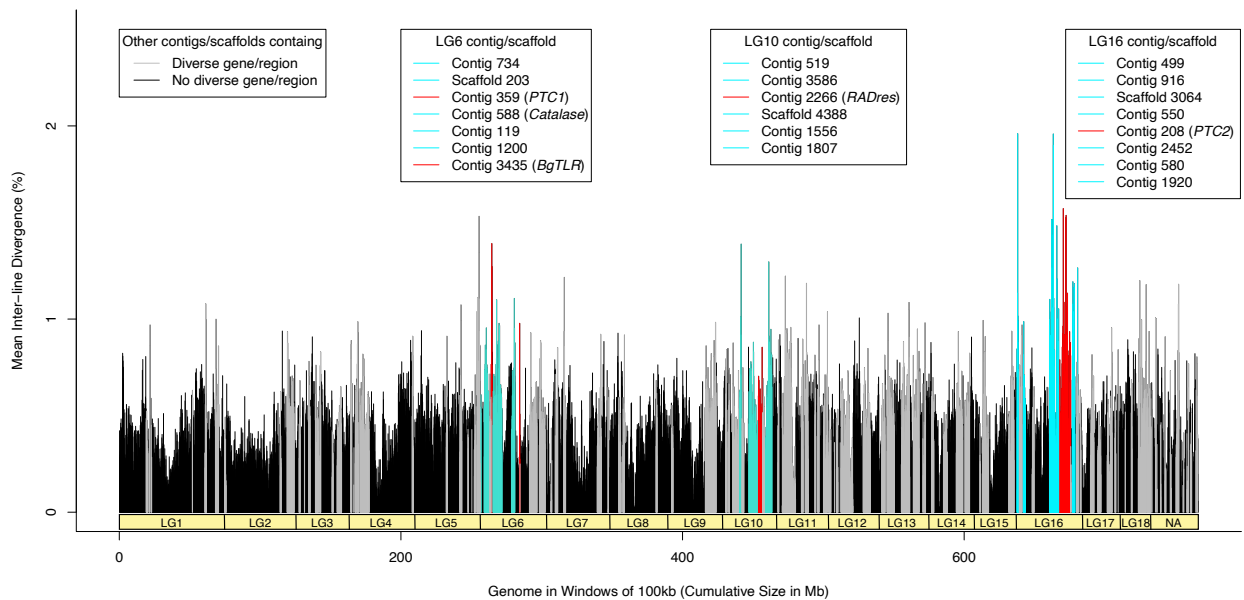
<i>Biomphalaria sudanica</i> inbred line	NCBI Accessions	Raw total sequences	Paired reads mapped	Bases mapped CIGAR	Coverage	Heterozygous Sites (percentage of genome)	Number of missing sites (percentage of genome)
Bs110	In prep.	174,760,210	86,796,593	23,247,544,501	24.6	0.13%	4.75%
Bs163	In prep.	207,521,490	102,626,199	27,823,100,537	29.4	0.11%	6.16%
Bs5-2	In prep.	548,701,568	272,013,259	73,112,262,010	77.2	0.17%	3.83%
BsKEMRI	In prep.	169,400,086	83,924,684	22,815,557,212	24.1	0.21%	5.31%

437 *Characterization and categorization of highly diverse genes and genomic regions in Biomphalaria*
438 *sudanica* genome

439

440 We followed a novel bioinformatic pipeline (see Methods section: *Assessment of highly*
441 *diverse genes and genome regions for novel pathogen recognition receptors*) to identify putative
442 immune related genes (not just those involved in schistosome immunity) in the *B. sudanica*
443 genome that may be under balancing selection, further narrowing these down to genes coding for
444 membrane-associated proteins that could represent PRRs. We first identified genes in genomic
445 regions that showed high divergence between the inbred lines of *B. sudanica*, hyperdiverse gene
446 selection being based on the highest nucleotide diversity across both the entire coding and/or
447 noncoding gene regions, and genes that were contained within the top 0.1-1% of the most
448 nucleotide diverse sliding windows between 10-100 kb. We found multiple gene clusters that were
449 of notably high diversity (Figure 4 and Figure 5). Following the bioinformatic pipeline, 1047
450 hyperdiverse genes (4.4% of all genes), from 184 different contigs/scaffolds were selected for
451 further analysis (Supplementary Table 17).

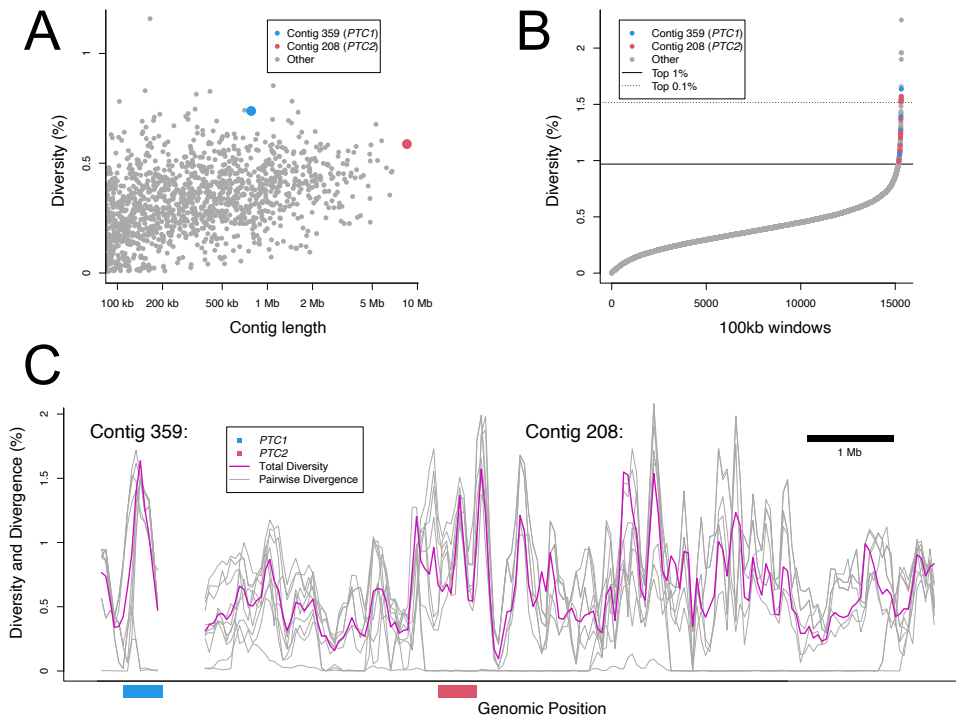
452



453

454 **Figure 4.** Mean nucleotide diversity of five *Biomphalaria sudanica* inbred line genomes in
455 windows of 100 kb (0 kb stagger). Linkage groups (LG1-LG18) for *B. sudanica* are inferred from
456 *B. glabrata* (6), showing the hypothesized chromosomal position of contigs in *B. sudanica*.
457 Notable clusters of highly diverse genomic regions and genes were seen in *B. sudanica* linkage

458 groups LG6, LG10 and LG16 (highlighted in turquoise and red peaks). Peaks in red represent
459 regions containing candidate immune loci that are orthologous to some of those previously
460 associated with *Schistosoma mansoni* resistance in *B. glabrata* (*PTC1*, *PTC2*, *Catalase*, *BgTLR*,
461 *RADres*, see Supplementary Table 2).
462
463



464
465 **Figure 5.** (A) Patterns of nucleotide diversity across genomes of five inbred lines of *Biomphalaria*
466 *sudanica*, highlighting two polymorphic transmembrane gene clusters that are orthologous to those
467 previously associated with *Schistosoma mansoni* resistance in *B. glabrata* (*PTC1* (40) and *PTC2*
468 (41)). (B) Genome-wide nucleotide diversity across overlapping 100-kb genomic windows
469 (starting at 0-kb and 50-kb intervals) with windows on contigs 359 (*PTC1*) and 208 (*PTC2*) that
470 occur in the top 1% of genome wide nucleotide diversity between inbred lines being colored blue
471 and red, respectively. (C) Genome-wide nucleotide diversity (purple line) and pairwise divergence
472 for each haplotype pair (grey lines) across contigs 359 (*PTC1*) and 208 (*PTC2*). *PTC1* and *PTC2*
473 regions are indicated by the blue and red bars, respectively. Similarly diverse regions span across
474 several megabases of contig 208 in *B. sudanica*. Even in diverse regions, pairwise divergence can
475 be near zero, indicating shared haplotypes.
476

477 Characterized genes in diverse regions were categorized based on their annotation
478 descriptions (BLAST results, InterPro, Pfam, DeepTMHMM) as to what predicted role in immune
479 function they may have. We determined 182 (17.4% of shortlist) proteins had no role in immune
480 function (group 1) and 245 (23.4%) proteins had a role (group 2), or 138 (13.2%) a potential role
481 (group 3), in immune function. Of the other shortlisted genes, 81 (7.7%) had an unknown function
482 but contained transmembrane domains (TMDs) (group 4), while the remaining 401 had an
483 unknown function (could not be interpreted due to an absence of information) without containing
484 a TMD (group 5). In performing comparative gene ontology analysis in topGO (72), several GO
485 terms including immune-relevant terms concerning molecular binding and receptor activity
486 classes, CD40 receptor complexes, and defense/inflammatory/immune responses as well as
487 regulation of innate immune responses, were significantly enriched ($p < 0.05$ Fisher's test) in this
488 subset of 1047 highly diverse genes compared to the full *B. sudanica* genome (see Supplementary
489 Table 18).

490 Although our novel pipeline was aimed at identifying immune genes under balancing
491 selection in *B. sudanica* that could be due to interactions with any snail pathogen, some of those
492 shortlisted were of relevance to candidate *S. mansoni* immune genes. Twelve of the 245 genes in
493 diverse regions that were suspected to have immune function (group 2) in *B. sudanica* were
494 orthologues to candidates for innate immune genes associated with *S. mansoni* resistance and
495 likely act as PRRs in *B. glabrata*, namely *PTC1* (n=5), *PTC2* (n=6), and *BgTLR* (n=1)
496 (Supplementary Table 17). In addition, four identified VIGLs were identified in highly diverse
497 gene regions (Supplementary Table 17). These were FREP12 (BSUD.12502), and two neighboring
498 FREPs on contig 777, BSUD.20519 (determined as an unknown FREP in N0.HOG0021173,
499 Figure 3) and BSUD.20520 (grouped within FREP12 N0.HOG0000107, Figure 3). In addition,
500 one full CREP (BSUD.13843, Figure 2 and Supplementary Table 14), was also identified within
501 the highly diverse genes/genomic regions analysis (Supplementary Table 17). Several other
502 immune suspected (group 2) proteins contained functional domains of interest, such as Fibrinogen-
503 related domains (FREDs), Fibronectin, C-type lectin (incomplete CREPs) and IgSF domains.

504

505 *Pathogen recognition receptor candidates from the highly diverse genes and genomic regions in*
506 *Biomphalaria sudanica genome*

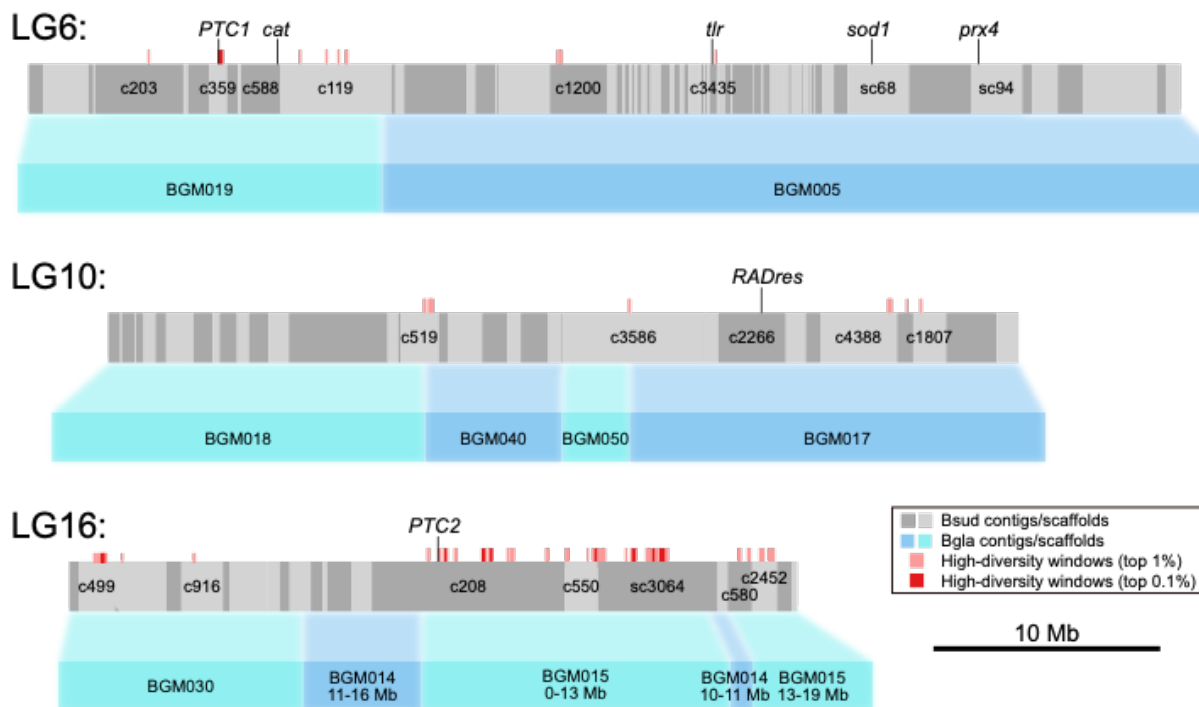
507

508 From the 1047 hyperdiverse genes identified (see above), a list of 20 proteins that contained
 509 at least one transmembrane domain with the highest coding nucleotide diversity between the five
 510 *B. sudanica* inbred line genomes were identified as potential PRR candidates (Table 3).

511
 512 **Table 3.** Top 20 diverse proteins (as determined by coding nucleotide diversity) of *Biomphalaria*
 513 *sudanica* that contain at least 1 transmembrane domain, and thus represent putative pathogen
 514 recognition receptors. * Non-truncated BSUD.4885 protein is 405 aa in length as determined by
 515 manual alignment.

Gene (contig/scaffold, Linkage Group)	Nucleotide Diversity % (coding)	Protein length (aa)	Inferred protein-receptor function / functional domains / protein family
BSUD.20937 (80, LG1)	5.7	453	Tumor Necrosis Factor (<i>TNF</i>)
BSUD.4884 (2266, LG10)	5.3	86	Unknown, linked (same contig) to <i>RADres</i> (78). Bacteriocin domain
BSUD.4885 (2266, LG10)	4.8	354*	Unknown, linked (same contig) to <i>RADres</i> (78). Domains: TMEM154 protein family, Rifin
BSUD.3983 (208, LG16)	4.0	817	Unknown, linked (same contig) to <i>PTC2</i> (41). Domains: SHP2-interacting transmembrane adaptor protein, SIT
BSUD.4096 (208, LG16)	3.7	270	Unknown, linked (same contig) to <i>PTC2</i> (41). Domains: Viral glycoprotein L
BSUD.20268 (7608, LG6)	3.3	75	Unknown. Domains of unknown function (DUF6768)
BSUD.9255 (379, LG11)	3.3	453	Unknown Domains of unknown function (DUF4781), Death domain
BSUD.15077 (580, LG16)	3.1	318	Cell adhesion molecule. Wnt and fibroblast growth factor (FGF) inhibitory regulator and Protocadherin domain. [SignalP = Membrane]
BSUD.4112 (208, LG16)	2.9	341	Cell adhesion molecule linked (same contig) to <i>PTC2</i> (41). Protocadherin domain. [SignalP = Membrane]
BSUD.18104 (69, LG15)	2.8	1154	G-protein coupled receptor. 7 transmembrane receptor (rhodopsin family), Leucine-rich repeat
BSUD.3984 (208, LG16)	2.7	324	Polymorphic transmembrane cluster 2 (<i>PTC2</i>) gene 9 (41). [SignalP = Membrane]
BSUD.14211 (540, LG5)	2.7	516	Cell adhesion molecule. F5/8 type C domain. [SignalP = Membrane, Galactose-binding-like domain]
BSUD.17807 (6862, LG13)	2.7	504	Cell adhesion molecule. Ephrin type-A receptor 2 transmembrane and immunoglobulin domain. [SignalP = Membrane]
BSUD.14415 (550, LG16)	2.6	483	Unknown. Wolframin EF-hand domain
BSUD.9838 (3979, LG4)	2.6	203	Transient receptor potential cation channel subfamily M member 2-like. Ion transport domain
BSUD.12903 (499, LG16)	2.6	647	Receptor with C-type lectin and immunoglobulin domain
BSUD.23928 (4388, LG10)	2.5	944	G-protein coupled receptor. 7 transmembrane receptor (rhodopsin family), Low-density lipoprotein receptor domain class A, Leucine-rich repeat
BSUD.8884 (359, LG6)	2.4	565	Guadeloupe Resistance Complex (<i>GRC/PTC1</i>) gene 2 (40,79). Fibronectin type III domain, TMEM154 protein family. [SignalP = Membrane]
BSUD.4003 (208, LG16)	2.4	473	Unknown, linked (same contig) to <i>PTC2</i> (41). Prodomain subtilisin 2, TMEM154 protein family. [SignalP = Membrane]
BSUD.4089 (208, LG16)	2.3	331	Receptor with C-type lectin domain. [SignalP = Membrane]

517 The majority (n=14) of the most diverse immune/PRR genes (Table 3) were clustered on
518 three linkage groups: 6, 10 or 16, which were also the regions of highest diversity overall (Figure
519 4 and Figure 6). These linkage groups also contain genes/clusters of genes that have previously
520 been inferred to have associations with schistosome resistance mechanisms in *B. glabrata* (*PTC1*,
521 Catalase (*cat*), *BgTLR*, *sod1*, *prx4*, in LG6, *RADres* in LG10 and *PTC2* in LG16, see Figure 6).
522 Direct orthologs to the candidate PRR genes *PTC2* gene 9 (BSUD.3984) and *PTC1* gene 2
523 (BSUD.8884) were identified in the topmost diverse PRR-like genes in the *B. sudanica* genome
524 (Table 3).
525



526
527
528 **Figure 6.** Orthology of *Biomphalaria sudanica* contigs/scaffolds (grey boxes, with different
529 shades of grey representing alternating contigs/scaffolds) to *B. glabrata* scaffolds (blue boxes, (6))
530 pertaining to linkage groups (LG) inferred from three *B. glabrata* linkage groups, highlighted here
531 because they are notably enriched for both diverse regions/genes and orthologous candidate
532 immune genes from *B. glabrata* (*PTC1*, Catalase (*cat*), *BgTLR*, *sod1*, *prx4*, *RADres* and *PTC2*),
533 positions of which are shown. Contigs/scaffolds with candidate genes and/or diverse regions (100
534 kb windows in top 1%, light red boxes, or top 0.1%, dark red boxes) are labeled. Synteny is

535 relatively high between species, except for a large rearrangement on LG16 (see Supplementary
536 Figure 3C).

537

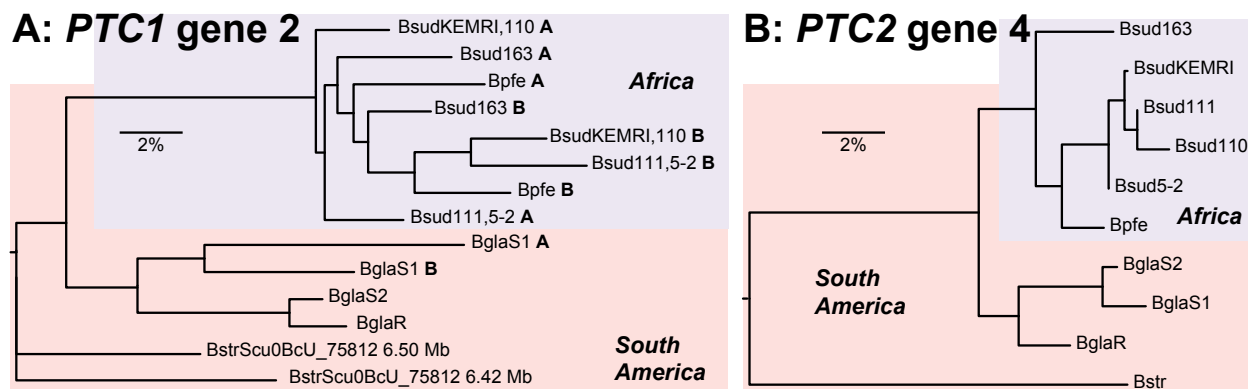
538 *Pathogen recognition receptor candidates within Linkage Group 6*

539

540 The order of genes in LG6 was highly conserved (Supplementary Figure 3A). Within LG6,
541 the hyperdiverse protein, BSUD.8884, was identified as a PRR candidate (Table 3), and is an
542 ortholog of the *GRC/PTC1* gene 2 that is tied to schistosome resistance in *B. glabrata* (40). This
543 gene shows a tandem duplication in African *Biomphalaria* species (tandem neighboring gene in
544 *B. sudanica* is BSUD.8885), and this duplication event has occurred independently of a similar
545 duplication seen in some *B. glabrata* haplotypes (Figure 7). A second shortlisted hyperdiverse
546 protein, BSUD.20268 in contig 7608 (Table 3), could not be well described upon comparisons to
547 other domains or proteins, perhaps due to its small (75 aa) size.

548 The ortholog to the *BgTLR* gene, BSUD.8256 (see *tlr* Figure 6), is also contained within
549 LG6 in a highly diverse region of the *B. sudanica* genome, although the gene itself is relatively
550 conserved within *B. sudanica* (0.1% coding nucleotide diversity). Additional immune-related
551 genes on this linkage group are listed in Supplementary Table 17.

552



553

554 **Figure 7.** Allelic phylogenies of *Biomphalaria sudanica*, *B. pfeifferi*, and *B. glabrata*, rooted with
555 *B. straminea*. (A) Polymorphic transmembrane cluster 1 (*PTC1*) gene 2 (i.e. *grctm2* (40)) shows a
556 tandem duplication (see A and B duplicates for relevant taxa indicated in bold) in the African
557 species (A=BSUD.8884 and B=BSUD.8885 for *B. sudanica*) that is clearly independent of similar
558 duplications seen in *B. glabrata* haplotype S1 and *B. straminea*. (B) *PTC2* gene 4 (BSUD.3979 in
559 *B. sudanica*) shows distinct haplotypes in each *Biomphalaria* species.

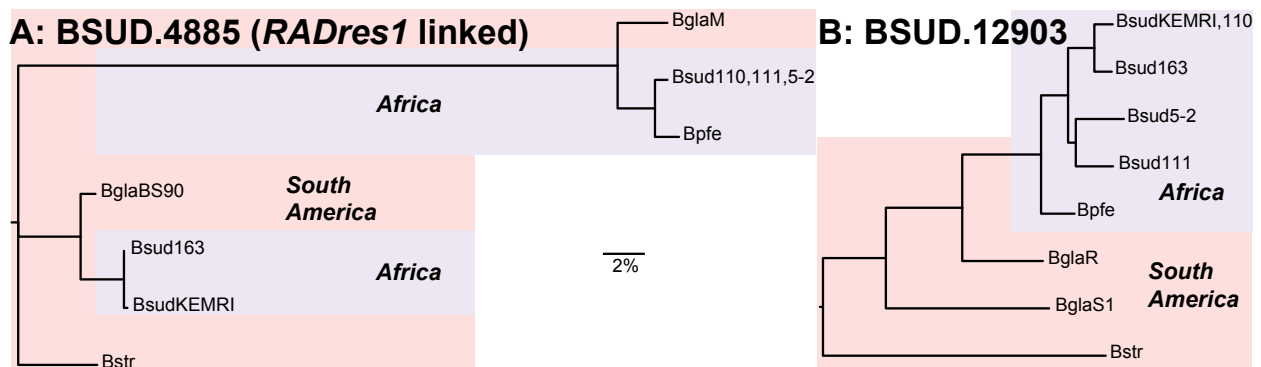
560

561 *Pathogen recognition receptor candidates within Linkage Group 10*

562

563 Synteny between *B. glabrata* and *B. sudanica* across LG10 was highly conserved
564 (Supplementary Figure 3B). Two proteins coded by neighboring genes located less than 1 Mb from
565 the orthologous site of schistosome resistance marker *RADres1* (78), BSUD.4884 and BSUD.4885
566 (contig 2266), showed extremely high (~5%) nucleotide diversity between the *B. sudanica* inbred
567 lines (Table 3). The predicted function of BSUD.4884, an 86 aa protein partially matching a
568 hypothetical protein recorded in *B. glabrata* (GenBank accession KAI8743980), could not be
569 confidently characterized due to an absence of significant matches to orthologous proteins or
570 domains (Supplementary Table 19). Comparisons with the annotated *B. glabrata* genome also
571 revealed that the neighboring protein, BSUD.4885, had been erroneously truncated on the
572 extracellular portion; the manually annotated gene codes for a 405 aa protein. The complete
573 BSUD.4885 protein was predicted to contain both TMEM154 and Rifin domains (see (80)),
574 features shared with some *PTC1* and/or *PTC2* genes (see BSUD.3980, BSUD.8873, BSUD.8874,
575 BSUD.8876 and BSUD.8884, Table 3 and Supplementary Table 17). Remarkably, haplotype
576 lineages of BSUD.4885 in *B. sudanica* also occur in *B. glabrata* (Figure 8), suggesting that either
577 this is a shared polymorphism within the genus that has survived the colonization of Africa, or that
578 there were originally two divergent tandem paralogs but one or the other is independently deleted
579 in every genome since no genome appears to have two copies of this gene.

580



581

582 **Figure 8.** Phylogenetic trees generated using RAxML (81) of exemplar genes showing unusually
583 high diversity in *Biomphalaria sudanica*, alongside orthologous alleles in *B. pfeifferi* and *B.*
584 *glabrata*, and rooted with *B. straminea*. (A) BSUD.4885 (contig 2266, linkage group (LG) 10) is

585 a gene with exceptionally high diversity, has a protein structure that shows similarities to genes in
586 the polymorphic transmembrane cluster 1 (*PTC1*) (40) and *PTC2* (41), and is closely linked to a
587 genomic marker, *RADres1*, that was previously shown to significantly influence schistosome
588 resistance in *B. glabrata* 13-16-R1 strain (78). Phylogeny shows an apparent trans-species
589 polymorphism of divergent haplotypes in African and South American snails. (B) BSUD.12903
590 (contig 499, LG16) is another pathogen recognition receptor candidate, in one of the most
591 polymorphic contigs in the *B. sudanica* genome, predicted to contain C-type lectin,
592 immunoglobulin, TMEM154 and alternatively expressed fibronectin III domains (see Figure 10).

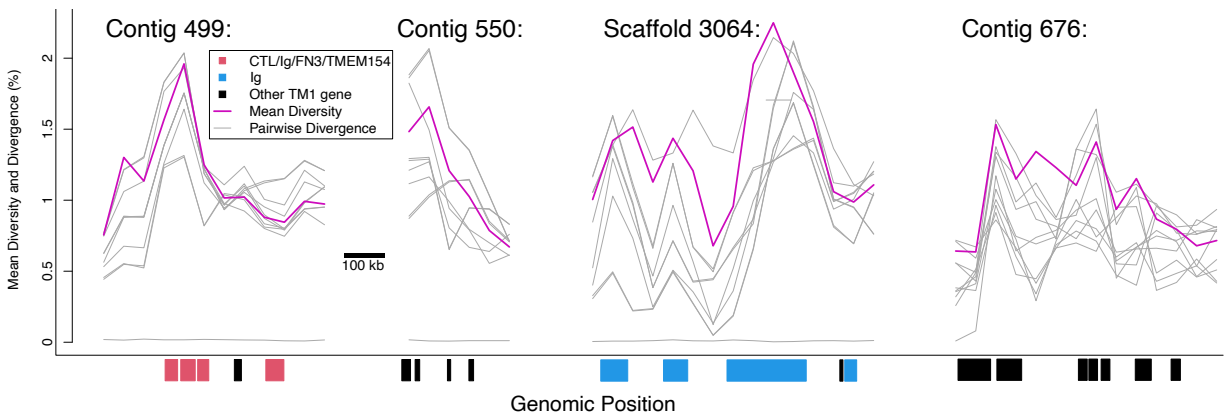
593

594 *Pathogen recognition receptor candidates within Linkage Group 16*

595

596 Linkage group 16 contains many of the most diverse genomic regions (Figure 9), such as
597 contig 208, contig 499, contig 550, contig 580, and scaffold 3064, which include highly diverse
598 individual genes (Supplementary Table 17). Six of the PRR candidates within LG16 originate from
599 contig 208 (BSUD.3983, BSUD.4096, BSUD.4112, BSUD.3984, BSUD.4003, BSUD.4089,
600 Table 3), the orthologous region to the *PTC2* region in *B. glabrata* (41). As the orthologous *PTC2*
601 genes are adjacent to these and many other diverse genes, several of which have single
602 transmembrane domains and/or homology to *PTC2* genes (e.g. BSUD.4003), the resistance-
603 associated *PTC2* region originally described in *B. glabrata* appears to be part of a much wider
604 region of contig 208 (Figure 5C; Figure 6) with distinctive evolutionary and structural features.
605 Furthermore, contigs neighboring contig 208 according to our inferred linkage map also contain
606 single transmembrane domain proteins with homology to *PTC2* genes (BSUD.14425,
607 BSUD.15084 and BSUD.23257 on contig 550, contig 580 and scaffold 3064, respectively).

608

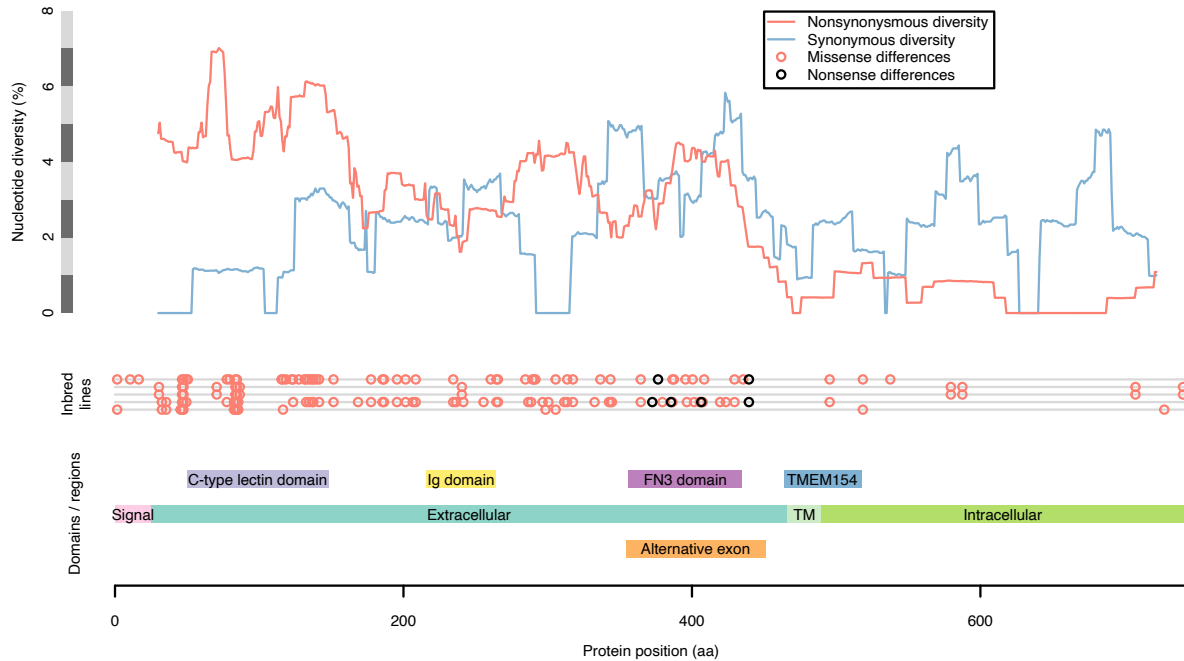


609
610 **Figure 9.** Mean nucleotide diversity (purple line) and pairwise divergence for each haplotype pair
611 (grey lines) across portions of four contigs in LG16 (contig 499 300-850kb, contig 550 0-250 kb,
612 scaffold 3064 2300-3000 kb) and LG5 (contig 676 150-850 kb) showing high diversity and clusters
613 of transmembrane genes. All plotted genes (shown in red, blue and black) encode single-pass
614 transmembrane proteins (TM1); other genes in these regions are not shown. Key functional
615 domains potentially involved in pathogen recognition include C-type lectin (CTL), fibronectin
616 type III (FN3) and immunoglobulin (Ig), and TMEM154, which is a membrane-spanning domain
617 also found in several polymorphic transmembrane cluster 1 (*PTCI*) genes (40). Genes shown in
618 red, including BSUD.12903 (Figure 10) all have at least three of these four functional domains,
619 while genes shown in blue have only Ig.

620
621 PRR candidates BSUD.4112 (contig 208) and BSUD.15077 (contig 580) were predicted
622 to be cell adhesion molecules since both contained a domain matching that of a protocadherin
623 (Table 3; Supplementary Table 19). Two other PRR candidates, BSUD.12903 (contig 499) and
624 BSUD.4089 (contig 208) (Table 3), both contain an extracellular C-type Lectin domain known to
625 function in immune responses to pathogens. For BSUD.12903 at least, nonsynonymous diversity
626 substantially exceeds synonymous diversity in the extracellular region of the protein where these
627 predicted functional domains (C-type lectin and IgSF) were present (Figure 10). BSUD.12903
628 occurs within a cluster of proteins representing one of the most diverse regions of the genome,
629 many of which contain combinations of C-type lectin, IgSF and FN3 domains (Supplementary
630 Table 17). BSUD.4089 encodes a longer isoform (331 aa) with a single transmembrane domain
631 and an extracellular C-type lectin domain, and a shorter isoform (67 aa) with just the secreted

632 peptide signal (confirmed through InterProScan, see Supplementary File 3) and a partial match to
633 the C-type lectin domain, which would result in a truncated domain.

634



635

636 **Figure 10.** Detailed view of protein BSUD.12903 containing the alternatively expressed exon of
637 741 aa, which demonstrates that in the extracellular region nonsynonymous diversity greatly
638 exceeds synonymous diversity (shown here in 50 aa sliding windows) in regions where functional
639 domains potentially involved in pathogen recognition are present (C-type lectin, immunoglobulin,
640 and fibronectin III (FN3) domains). Two inbred lines of *Biomphalaria sudanica* have multiple
641 nonsense variants (stop codon or frameshift) in a single exon containing a FN3 domain, which is
642 expressed in the *B. pfeifferi* ortholog (KAK0057508 (8)), suggesting that FN3 may not be
643 expressed in all *B. sudanica* lines. The FN3 exon also occurs in the *B. glabrata* ortholog
644 BGLB024560, where it is also variably either expressed (NCBI Accession XM_056010427) or
645 excluded (NCBI Accession SRX8534561). Both FN3 and/or TMEM154 domains are present in *B.*
646 *sudanica* genes BSUD.8884, BSUD.8874 and BSUD.8876 that are orthologous to *B. glabrata*
647 *PTC1* region genes associated with schistosome resistance: *grctm2*, *grctm3* and *grctm4* (40).

648

649 A large inversion (~15 Mb) between *B. sudanica* and *B. glabrata* iM line on LG16 appears
650 to reverse parts of scaffold 3064, contig 208 and contig 550 relative to *B. glabrata* scaffolds

651 BGM014 and BGM015 (Figure 6; Supplementary Figure 3C). BGM014 itself may be incorrectly
652 assembled since it is split across *B. glabrata* linkage groups LG3 and LG16.

653

654 **Discussion**

655

656 *A new genomic resource for vector biology*

657

658 Long-read PacBio HiFi DNA and RNA sequencing, supplemented by Illumina short-read
659 sequencing of RNA, was performed on an inbred line of *B. sudanica* sensu lato, originally collected
660 from shoreline habitats of the Kisumu region of Lake Victoria. This sequencing approach and
661 bioinformatic pipeline resulted in a well-supported genome annotation, which adds to the growing
662 repository of *Biomphalaria* species genomes. *Biomphalaria sudanica* remains a “neglected
663 vector” despite its established role in transmission of *S. mansoni* in lake and marsh habitats in the
664 African Rift Valley, where schistosomiasis transmission is among the highest in the world. Our
665 analyses focused on characterizing immune genes using a comparative framework with that of the
666 better studied South American vector of schistosomiasis, *B. glabrata*, and another African vector,
667 *B. pfeifferi*, for which genomic resources have recently been developed (8). We also developed a
668 novel bioinformatic pipeline to characterize regions of marked intraspecific diversity as a
669 mechanism to identify novel genes that may be involved in pathogen recognition and are under
670 balancing selection. These analyses and resources will facilitate future work to explore further and
671 uncover additional vector immune defense mechanisms against pathogens such as schistosomes,
672 with direct relevance to public health in Africa where the majority of *S. mansoni* infections occur.
673 The hope is that the scientific community will use these resources to support the development of
674 novel tools for schistosome control, potentially including gene-drive manipulation of the snail
675 vector (82), and development of surveillance tools, as is described as a priority in World Health
676 Organization guidelines (83).

677

678 *Comparative analysis of immune genes among species*

679

680 Based on the average number of substitutions per site, the genomes of *B. sudanica* and *B.*
681 *pfeifferi* are approximately 5.1% divergent from each other, which is comparable to that of the

682 isolates of *B. glabrata* compared here (between 3.3-8.7%) (see Figure 1). We found that *B.*
683 *sudanica* possesses orthologous genes to all those previously identified from *B. glabrata* as having
684 a potential function in immunity against schistosomes. However, the functions of these orthologs
685 have yet to be verified for *B. sudanica*.

686 The number of complete variable IgSF and lectin domain-containing molecules (FREPs
687 and CREPs) in *B. sudanica* was almost identical to that observed in *B. pfeifferi*, with both species
688 having a higher number than observed in *B. glabrata* (Table 1). The largest and most divergent
689 FREP family in *B. sudanica* is FREP12, two genes of which were not identified by the VIGL
690 annotation pipeline but instead were discovered with the pipeline to identify highly diverse
691 immune related proteins and PRRs. Within the African species, some FREP-like gene families
692 were enriched in *B. sudanica* and underrepresented in *B. pfeifferi*. This association is interesting
693 given the observation that *B. pfeifferi* is more susceptible to schistosome infections compared to
694 *B. sudanica* (84,85); however the bulk of enriched FREPs in *B. sudanica* (FREPK1, FREPJ3,
695 FREP5, FREPJ5, FREP12) are not known to play a role in schistosome resistance as based on
696 studies of *B. glabrata* with comparable studies for African species yet to be undertaken. It is also
697 apparent that FREP genes duplicate more readily in *B. sudanica*, resulting in more truncated
698 FREPs than CREPs.

699 It should be acknowledged that the annotation and nomenclature of FREPs is not trivial.
700 Diversification of FREPs by gene conversion, duplications, gene loss, exon shuffling, and other
701 structural rearrangements, rather than the slow accumulation of mutations, complicates
702 phylogenetic analysis and therefore classification and nomenclature of these genes (43). Due to
703 the complex molecular evolution of FREPs, the *B. sudanica* genome contains a diverse and
704 potentially variable suite of FREP sequences as seen in *B. glabrata* (86), independent of nucleotide
705 diversity at any particular gene. Further investigation into these hypervariable genes is necessary
706 to fully describe the molecular mechanisms that drive their diversity and interactions with
707 pathogens.

708 Elevated diversity at candidate innate immune gene regions, *PTC1* and *PTC2*, in both *B.*
709 *sudanica* and *B. glabrata* appears to be independently generated, since the genes remain
710 reciprocally monophyletic (Figure 7), which suggests ongoing selection favoring novel diversity
711 at these loci. Occasional transspecies polymorphisms, including one apparent at BSUD.4885

712 (Figure 8), are suggestive of a long-term balancing selection, strong enough to overcome what was
713 likely a narrow population bottleneck during the colonization of Africa.

714

715 *Intraspecific genomic diversity*

716

717 Hyperdiverse protein coding genes in *B. sudanica* that were shortlisted as potential PRRs
718 under balancing selection were enriched in molecules that have been associated with various
719 aspects of immune functions in other organisms. The hyperdiverse transmembrane proteins
720 included G-protein coupled receptors (GPCRs), TLRs, cluster of differentiation (CD) and cell
721 adhesion molecules as well as proteins containing functional domains such as C-type lectins, FN3,
722 IgSF, transient receptor potential (TRP) channels and many potentially immune related domains
723 that await further immunological characterization. These often show high extracellular nucleotide
724 non-synonymous diversity which overlaps with functional binding domains (see Figure 10). While
725 we have focused on candidate PPRs with high nucleotide diversity, we do not mean to discount
726 the importance of conserved genes, or the between-paralog diversity found in multi-gene families.
727 In particular, genes with well-studied immune roles in *B. glabrata*, such as FREPs (42) and
728 components of the oxidative burst pathway (87), undoubtedly also play essential roles in *B.*
729 *sudanica* immunity regardless of their allelic diversity.

730 The most diverse candidate protein (Table 3; BSUD.20937) has a single transmembrane
731 domain and a predicted tumor necrosis factor (TNF) domain. TNFs in invertebrates are far less
732 characterized than those in mammals; however, recent work in mollusks and crustaceans indicate
733 that just like in vertebrates, they have multiple roles in innate immunity and function in activation
734 of antimicrobial peptides, apoptosis and phagocytosis of hemocytes, and activation of immune
735 related enzymes such as lysozyme and phenoloxidase (88–90). We could not detect this gene in
736 two inbred lines (Bs163 and Bs5-2), presumably due to either deletion or extreme sequence
737 divergence, and thus our diversity measurement may even be an underestimate. The high diversity
738 of BSUD.20937 between *B. sudanica* inbred lines, and that of the other hyperdiverse
739 transmembrane genes selected in our bioinformatic pipeline, suggests that these may well be PRR
740 genes under balancing selection driven by variability in ligand binding sites that interact with
741 pathogens.

742 Regarding genome-wide patterns of diversity, the identified high diversity windows were
743 clustered on three linkage groups (6, 10, and 16), and near candidate immune genes. Linkage group
744 6 contains a cluster of highly diverse windows as well as several candidate immune genes
745 including *PTC1*, *cat*, *tlr*, *sod1*, and *prx4*. Linkage group 10 also harbors diverse windows as well
746 as resistance region *RADres*. The highest diversity genomic windows occur on linkage group 16,
747 site of the highly diverse gene cluster *PTC2* which we show is imbedded within a larger region of
748 diversity (Figure 4 and 6) that shows chromosomal rearrangement between *B. sudanica* and *B.*
749 *glabrata* (Supplementary Figure 3C). Such structural rearrangements may be enriched along with
750 single-nucleotide polymorphisms in hyperdiverse regions like *PTC2* and may even help to
751 maintain diversity if they are segregating within species; however we also note that the 8.5 Mb
752 contig 208 is the largest contig in our Bs111 assembly and thus it afforded the greatest power to
753 observe large rearrangements that may be undetectable elsewhere in the genome. Immune-relevant
754 genes may often be linked within chromosomal regions with distinct evolutionary dynamics,
755 including perhaps the maintenance of elevated adaptive diversity, as also noted for *B. glabrata*
756 (91,92).

757

758 *Molecular evolutionary legacy of colonization and adaptation in Africa*

759

760 *Biomphalaria* originated in South America, but at least one lineage crossed the Atlantic
761 based on previous estimates approximately 1.8 and 5 MYA (8–11) and diversified into the
762 contemporary African species. Such a transcontinental colonization, whenever it did take place, is
763 striking for a freshwater gastropod of limited vagility.

764 Since the transcontinental colonization of *Biomphalaria* is predicted as sufficiently
765 evolutionarily recent, its history can be explored more easily with the growing collection of
766 *Biomphalaria* genomes. Despite the divergence between *B. glabrata* and the African species being
767 between 10.5-10.8% (Figure 1), thousands of orthologs could be aligned easily across these species
768 and without saturation of neutral sites. These highly orthologous genomes allowed us to examine
769 the expansion and contraction of gene families, with particular interest being paid to the genomic
770 consequences of migration and adaptation of *Biomphalaria* to Africa. Originally, we hypothesized
771 that immune genes of this African *Biomphalaria* ancestor would be substantially expanded as it
772 encountered several new pathogens as it colonized Africa from South America; however, the

773 analysis revealed no significantly expanded immunity-related genes. In fact, the contraction of
774 genes outweighed expansions in the common ancestor of *B. sudanica* and *B. pfeifferi* compared to
775 the South American congeners, suggesting an overall simplification of the genome in *B. sudanica*
776 and *B. pfeifferi*. Perhaps this gene family contraction in African *Biomphalaria* was a maladaptive
777 result of a founder effect, given that the ancestor of African *Biomphalaria* must have descended
778 from a small number of colonizers, and/or perhaps many genes were freely lost without fitness
779 consequences in the absence of the neotropical pathogens to which they were adapted. Despite this
780 proteome simplification, the estimated genome size of *B. sudanica* was ~73 Mb larger than the
781 genome of outcrossing *B. glabrata*, and also ~173 Mb larger than that of *B. pfeifferi*, which is
782 hypothesized to have lost genes associated with mating given its reliance on selfing (8). This size
783 discrepancy could be a technical artifact from divergent allelic *B. sudanica* haplotypes getting
784 assembled as separate contigs, perhaps due to the higher heterozygosity in our Bs111 after only
785 three generations of selfing, compared to naturally inbred *B. pfeifferi*, which is a preferentially
786 selfing snail (93,94), and the long-term inbred lines of *B. glabrata*. Alternatively, if the increased
787 genome size of *B. sudanica* is biologically real, it could represent gene duplications and insertions
788 not picked up in our current analysis, which may have led to novel gene functions and expression
789 mechanisms favored in its African habitat. The mechanism could also be selectively neutral: the
790 high repeat content of all *Biomphalaria* genomes indicates that total genome size can vary
791 considerably with little direct effect on the content of coding genes.

792 As confirmed through our annotation of RNAs, a common feature in the two African
793 species is the presence of tRNA-SeC, demonstrating the ability of *B. sudanica*, like its sister
794 species *B. pfeifferi* (8), to synthesize polypeptides containing selenocysteine (53,54). Although
795 selenoproteins are present in many gastropod lineages (95), tRNA-SeC has not been identified in
796 *B. glabrata* (8), suggesting that the capacity to produce selenocysteinyl has been gained in the
797 African *Biomphalaria* species, or lost across South American lineages of *B. glabrata*. With the
798 growing repository of genome information for species basal to *Biomphalaria*, it will be possible
799 to investigate this further. Selenocysteine-containing proteins, i.e. selenoproteins, are proposed to
800 be involved in a wide range of bioactive processes in other invertebrates (96).

801 Interestingly, regarding rRNA genes, far fewer were predicted in *B. sudanica* compared to
802 *B. pfeifferi*, although this is likely due to misassembly of the highly repetitive tandem repeats that

803 rRNA genes are generally present in, causing near identical rRNA copies to be collapsed into
804 single copies, rather than biological differences.

805

806 *Transcriptomic analysis illuminates unusual mitochondrial transcription processes*

807

808 Molluscan mitochondrial genomes are unusual in terms of their genomic features and
809 transcriptomic processes (56,58), and the latter has not been explored for African *Biomphalaria*.
810 In *B. sudanica*, mitochondrial genes that were not separated by tRNA genes, namely *nad6*, *nad5*,
811 and *nad1*, were highly represented by polycistronic mRNA transcripts with no clear trimming
812 points. Arrangement and transcription processes of these genes therefore appear similar in other
813 invertebrate and molluscan species (58,97,98). In addition, *nad4l* lacked read coverage, with the
814 few transcripts that were observed being polycistronic with *cob*. Lack of expression in key genes
815 of the respiratory chain, such as the NADH dehydrogenase, is not unheard of within mollusks (99),
816 but regardless the sharp increase on read coverage at the start of *cob*, suggests this is also processed
817 by endonucleases.

818 Our transcript for *B. sudanica* across mitochondrial genes *atp6* and *atp8*, which in other
819 species form a single mRNA (97), are punctuated by *trnN* in *B. sudanica* as in *B. glabrata* (59).
820 However, the bulk of the recovered reads belong to the pre-mRNAs. Considering their
821 polycistronic status in other organisms, it is possible that translation could happen before the pre-
822 mRNA is fully resolved, which would add another layer of complexity to molluscan mitochondrial
823 genomes in that these genes are likely separated downstream into proteins through initiation on
824 the ribosome (as reviewed in (58)). Additionally, given that *atp8* has been found under relaxed
825 selective pressures and a high degree of variation in both *Biomphalaria* and *Bulinus* species
826 (55,100), it may play a less important role in the mitochondrial function within these planorbids.

827

828 **Conclusions**

829

830 The species *B. sudanica* is remarkable for two principal reasons: its dynamic evolutionary
831 history involving recent colonization and diversification in Africa, and its tragic impact on public
832 health as a schistosome vector. The genomic data presented here illuminate both aspects. Our
833 observations include expansions and contractions of gene families in African snails, as well as

834 numerous genomic regions of high diversity that contain genes that may play a role in host-
835 pathogen coevolution and their vectorial capacity for parasites such as schistosomes. In
836 combination with increasing genomic resources from other *Biomphalaria* isolates, this work will
837 facilitate an enhanced understanding of the biology of these snails and future mechanisms for
838 curbing transmission of schistosomiasis.

839

840 **Methods**

841

842 *Ethical approval and permitting information*

843

844 This project was undertaken following approval from the relevant bodies, including Kenya
845 Medical Research Institute (KEMRI) Scientific Review Unit (permit # KEMRI/RES/7/3/1),
846 Kenya's National Commission for Science, Technology, and Innovation (permit #
847 NACOSTI/P/22/148/39 and NACOSTI/P/15/9609/4270), Kenya Wildlife Service (permit #
848 WRTI-0136-02-22 and # 0004754), and National Environment, Management Authority (permit #
849 NEMA/AGR/159/2022 and # NEMA/AGR/46/2014) (Registration # 0178).

850

851 *PacBio whole genome sequencing and assembly of *Biomphalaria sudanica* line 111*

852

853 High molecular weight DNA was isolated from the headfoot tissue of a single adult (>8
854 mm in shell diameter) *Biomphalaria sudanica* snail from the inbred line "111" (Bs111) following
855 the Qiagen Blood & Tissue kit (Qiagen, MD, USA), the only modification to standard kit protocol
856 being the overnight lysis of tissue at 37°C. The inbred line, Bs111, was developed from snails
857 originally collected from the Kisumu region of Lake Victoria, Kenya in 2012 and bred via selfing
858 for 3 generations before the line was expanded through sibling mating and gDNA extracted for
859 sequencing. PacBio high-fidelity (HiFi) circular consensus long-read sequencing (101) was carried
860 out using two SMRT cells of a PacBio Sequel II at the University of Oregon. PacBio ccs reads
861 (raw reads available at NCBI Sequence Read Archive accession: In prep.) were assembled using
862 Flye (102), settings: flye --pacbio-hifi Bsud111_ccs.fasta -g 1g -i 0 -t 16 -o flye/Bsud111. Basic
863 statistics of the genome assembly were calculated in seqkits v0.16.1(103).

864

865 *Illumina whole genome sequencing, alignment and genotype calling of four genetic lines of B.*
866 *sudanica*

867

868 Short-read Illumina data for four additional *B. sudanica* genetic lines Bs5-2, Bs110, Bs163
869 and BsKEMRI, all originally collected between 2010-2012 from the Kisumu region of Lake
870 Victoria, Kenya, were also produced. Inbred lines were generated via selfing in isolation for 3
871 generations, after which, the lines were expanded through sibling mating. The BsKEMRI line snail
872 was not subject to intentional inbreeding but was maintained in the laboratory for ~10 years (2010-
873 2020) before gDNA was extracted for sequencing. High molecular weight DNA was extracted as
874 described above. Illumina pair-ended reads (150 bp) were generated using the HiSeq 3000 at the
875 Center for Quantitative Life Sciences at Oregon State University to obtain 15-20x coverage for
876 each sample. The FASTQ files had adapters removed with Cutadapt v1.15 (104) and then reads
877 were trimmed with Trimmomatic v0.30 (45) with options: LEADING:20 TRAILING:20
878 SLIDINGWINDOW:5:20 MINLEN:50, and aligned using BWA version 0.7.12 (105) using
879 command: `bwa mem -P -M -t 4`, with the PacBio assembly of *Bs111* (see above). Aligned genome
880 data and raw reads are available in archive NCBI Sequence Read Archive (Accession: In prep.).

881 Genotypes for each snail genome were called against the *Bs111* reference genome in
882 variant call format (vcf) produced using BCFtools v1.9 (106). Basic statistics of the alignments to
883 the genome were calculated using vcftools v0.1.17 (107).

884

885 *Nuclear genome annotation of B. sudanica '111' from long- and short-read RNA sequence data*

886

887 The bioinformatic pipeline for the annotation of *B. sudanica* was performed using a
888 custom script (see: github.com/J-Calvelo/Annotation-Biomphalaria-sudanica). To obtain the best
889 possible genome annotation for *B. sudanica*, combined long- and short-read RNA-seq data was
890 used. RNA for long- and short-read sequencing was obtained from multiple samples of pooled
891 *Bs111* snails during two developmental stages (juvenile and adult) or under different stressors
892 (heat stress or exposure to *S. mansoni* NMRI strain from the NIH-NIAID Schistosomiasis
893 Resource Centre (108)). Full details of samples used, and RNA extraction protocol and pooling
894 are contained in the Supplementary Material (Supplementary File 6). The equal mass pool of eight
895 samples was processed for sequencing using the standard IsoSeq protocol (PacBio protocol: PN

896 101-763-800 v2). The NEBNext Single Cell/Low Input cDNA Synthesis and Amplification
897 Module kit (New England BioLabs Inc. MA, USA) was used following manufacturers protocols
898 for cDNA synthesis and amplification from 300 ng of RNA, and the library was generated using
899 the SMRTbell Express Template Prep Kit 2.0 (Pacific Biosciences, CA, USA). A Pronex bead
900 purification (Promega, WI, USA) was performed on the amplified cDNA product following a size
901 selection of \approx 2kb transcripts. The final cDNA library was sequenced on one SMRT cell on the
902 PacBio Sequel II at GC3F, University of Oregon.

903 Short-read Illumina sequences of RNA were obtained from two samples, one of a single
904 unchallenged Bs111 and the second a pool of three Bs111 adult snails that had been challenged to
905 *S. mansoni* miracidia 24 hours prior (Supplementary File 6). Library preparation and directional
906 mRNA sequencing were performed at Novogene Corporation Inc. (CA, USA). Library preparation
907 was performed using the NEBNext Ultra II Directional RNA Library Prep Kit (New England
908 BioLabs Inc. MA, USA) for strand-specific Illumina libraries. Libraries were then sequenced on
909 the NovaSeq 6000 platform (Illumina, CA, USA) and 150 bp pair-ended reads were generated to
910 provide \sim 12 G of sequence data.

911

912 *Bioinformatic analysis of long- and short-read RNA-seq data*

913

914 Quality of the long-read data was first assessed with longQC software (44). Then, reads
915 were processed with Lima (github.com/PacificBiosciences/barcoding), with the options "--css and
916 --dump-clips", and isoseq3 refine tool (github.com/PacificBiosciences/IsoSeq), with the option "-
917 -require-polya", to recover all complete sequenced transcripts (both 5' and 3' adapters and poly-A
918 tail should be present to consider a transcript as complete). Complete transcripts were then mapped
919 to the assembled genome with minimap2 (46,47) with options "-ax splice:hq -uf". In parallel,
920 short-read quality was verified using FASTQC software (109) and Illumina's adapter sequences
921 were removed along with low-quality bases using Trimmomatic (45), with options
922 "ILLUMINACLIP:TruSeq3-PE-2.fa:2:30:10 SLIDINGWINDOW:5:20 MINLEN:25". The
923 cleaned reads were then mapped to the genome using STAR (48) with options --
924 genomeChrBinNbits 13 --genomeSAindexNbases 13, as suggested by the authors for the genome
925 assembly. Alignment quality was measured with Bamtools v2.5.2 (110). Transcripts were then
926 characterized using both types of reads with StringTie2 v2.2.1 in "--mix" mode (49) plus the "--

927 conservative” option. Transcriptome completeness was evaluated with BUSCO (52) using
928 Mollusca as the reference. Nuclear coding sequences longer than 150 bases (50 amino acids) were
929 predicted with the utility tool TransDecoder.Predict v5.5 of the Trinity platform (50), with the
930 option --single_best_only, making use of homology information from UniProt’s Swiss-Prot (111)
931 and Pfam v35 protein domains (112). Similarity searches were done with BLASTp, with the “-
932 evaluate 1e-5” parameter (113) and hmmscan program with default parameters (114), for both
933 databases, respectively.

934

935 *Functional annotation*

936

937 Predicted proteins determined by nuclear coding sequences of TransDecoder.Predict v5.5
938 (50) were functionally annotated using eggNOG-mapper v2.0 (database downloaded 7/22/2022
939 (70,71)) for Gene Ontology (GO Terms) assignment based on curated orthogroups, and
940 InterProScan v5.56-89.0 (51) for protein identification based on domain prediction.

941

942 *Repeat identification and masking of nuclear genome*

943

944 Repetitive elements were annotated and masked using Earl Grey v 1.2 (115) pipeline. In
945 short, repeat elements were identified with RepeatMasker v4.1.2 (116) with the “-norna, -nolow,
946 -s” option and using Dfam v3.6 (117) curated database for the group Eumetazoa as reference.
947 Then, RepeatModeler v2.0.3 (118), RECON v1.08 (119) and RepeatScout v1.0.5 (120) generated
948 a *de-novo* repeat library that was refined by a “BLAST, Extract and Extend” process to combine
949 fragmented detections into a single repeat candidate. These repeats were finally used by a second
950 run of RepeatMasker v4.1.2 (116) to identify novel repeat elements. Summarizing plots were
951 generated with ggplot2 (121).

952

953 *Non-coding RNA annotation*

954

955 Transfer RNAs (tRNAs) and ribosomal genes (rRNA) were predicted using tRNAscan-SE
956 v2.0.9 (122) and Barrnap v0.9 (github.com/tseemann/barrnap), respectively, options for both
957 configured for eukaryotes.

958

959 *Mitochondrial Genome Annotation*

960

961 The contig encoding the mitochondrial genome was identified by BLAST searches of the
962 genome assembly (-perc_identity 80) (113) against the known mitochondrial protein coding genes
963 from *B. glabrata* (NC_005439.1 (59)). One retrieved, the DNA sequence (contig 7934)
964 independently annotated the retrieved DNA sequence with MITOS2 (60) web server
965 (mitos2.bioinf.uni-leipzig.de/index.py: accessed July 2022). Genes not automatically annotated by
966 mitos2 (*nad4l* and *trnK*, see Results) were localized with BLASTn (113) searches with respective
967 *B. glabrata* mitochondrial genes (BLASTn with options -task blastn and -task blastn-short for
968 *nad4l* and *trnK*, respectively).

969 RNA PacBio reads were mapped to the mitochondrial genome following the same
970 procedure as for the nuclear genome annotation (see above). Transcript alignments to each
971 mitochondrial gene were inspected on Integrative Genomics Viewer (IGV) (123) to manually
972 improve gene annotation making the assumptions that a) translation started on the first viable start
973 codon and b) stop codons were complete unless read coverage suggested a premature RNA end
974 with a truncated stop codon, similar to the assumptions described previously (58). Reads belonging
975 to intermediary pre-mRNA were counted with htseq-count (124) using the options "--
976 nonunique=all --samout" and custom scripts (see: [github.com/J-Calvelo/Annotation-](https://github.com/J-Calvelo/Annotation-Biomphalaria-sudanica)
977 *Biomphalaria-sudanica*). Using the pre-mRNA read data, a schematic representation of the post-
978 transcriptional trimming/modification processes of the primary transcript was then manually
979 generated using Inkscape (inkscape.org).

980 In line with other already published mitochondrial genomes for *B. sudanica* (NCBI RefSeq:
981 NC_038060.1 (55)) and other *Biomphalaria* species (NCBI RefSeq: NC_038061.1 and
982 NC_038059.1, GenBank Accession: MG431965.1 (55)) the start coordinate for the mitochondrial
983 genome was set to the first codon of the ND5 gene (Supplementary Figure 4). However, for manual
984 inspection of read alignments on IGV, a custom origin set to position 7222 (between two tRNA
985 molecules, tRNA-s1 and tRNA-s2, lacking both gene predictions and read coverage from the
986 PacBio RNA sequence dataset) was used.

987

988 *Gene family analysis and evolutionary position of B. sudanica*

989

990 A general study of the evolutionary dynamics (copy number and selection pressure) of
991 *Biomphalaria* protein families was carried out based on the longest proteins predicted for each
992 gene, in combination with 3 *B. glabrata* strains: BB02 (NCBI RefSeq: GCF_000457365.2 (5))
993 iBS90 and iM (GenBank Accession: GCA_025434165.1 and GCA_025434175.1 (6)), and the
994 species *B. pfeifferi* (NCBI BioProject Accession: GCA_030265305.1 (8)), *B. straminea* (GenBank
995 Accession: GCA_021533235.1 (7)), with the planorbid snail *Bulinus truncatus* (GenBank
996 Accession: GCA_021962125.1 (125)), and the Plakobranchidae *Elysia marginata* (GenBank
997 Accession: GCA_019649035.1 (126)) as an outgroup. With the goal of homogenizing criteria and
998 avoid discrepancies between the reported mRNA sequence and their proteins, protein sequences
999 taken from other works were predicted from the reported cDNA using getorf from the EMBOSS
1000 v6.6.0.0 package (127). The longest ORF among all isoforms (min size 150 bases) of each gene
1001 were selected, ties were broken by picking the most upstream candidates for each gene. Orthology
1002 relationships were estimated in HOGs with Orthofinder v2.5.4 (67) and treated as putative protein
1003 families. A species tree was generated in Orthofinder using the Species Tree from All Genes
1004 (STAG) algorithm (67,74), to confirm that the topology of the tree matches those previously
1005 reported for these *Biomphalaria* species (9,128).

1006 Significant expansions/contractions above the background were identified with CAFE 5
1007 (68). The Enriched GO terms among significantly expanded/contracted HOGs from CAFE 5 were
1008 determined using the topGO R package v2.48.0 (72), with the Weight01 algorithm, node size=10
1009 and statistical test of Fisher (significant p-value ≤ 0.05). To this end, GO Terms were assigned to
1010 each species' genes as described above, and those annotations were assigned to each HOG.
1011 Summary plots of gene lists were produced in REVIGO to allow further interpretation (73).

1012

1013 *Cellular location of proteins: secreted, mitochondrial translocated and transmembrane proteins*

1014

1015 Location signals for exportation or mitochondrial translocation of proteins were predicted
1016 with SignalP v6.0 (61) and TargetP v2.0 (62), respectively. Isoforms with negative results for both
1017 tools were additionally analyzed with SecretomeP v1.0 (63) to identify additional secreted proteins
1018 based on their biochemical characteristics, that is, proteins exported through an alternative route
1019 or with location signals missed due to annotation errors. For all three analyses only hits with a

1020 probability/score above 0.95 were considered significant. Lastly, TMDs were predicted with
1021 DeepTMHMM v1.0.13 (129).

1022

1023 *VIgL (FREP/CREP/GREP) identification and analysis*

1024

1025 Based on the presence or absence of known domains and features of FREPs, CREPs and
1026 GREPs, as reviewed previously (43), these VIgLs were identified following these criteria: 1)
1027 Evidence of secretion (see section: *Cellular location of proteins*), 2) presence of IgSF domains,
1028 which were identified using hmmsearch v3.3.2 (75) using the domain profiles generated previously
1029 (43), and 3) evidence of their respective 3' terminal domain in the InterPro annotation: FBD
1030 (IPR002181, IPR014716, IPR020837, IPR036056), C-lectin (IDs: IPR001304, IPR016186
1031 IPR016187, IPR018378) or Galectin (IPR001079, IPR015533, IPR044156, IPR000922,
1032 IPR043159). For a summary of the expected hits for each family, see Dheilly et al. (76). Lastly, a
1033 BLAST search between *B. sudanica* transcriptome and the *B. glabrata* FREP, CREP, and GREP
1034 genes reported by Lu et al. (43) and an additional 17 reference sequences (13 FREP and 4 CREP
1035 from NCBI (Supplementary Table 15) were performed.

1036 Full FREP, CREP, and GREP candidates were defined as genes with a SP, at least one
1037 IGsF, and their respective 3' terminal signature domain (FBD, C-lectin, or Galectin domain
1038 respectively). Classification into subfamilies was carried out according to their phylogenetic
1039 relationship with the reference sequences. For each candidate gene the longest isoform reported
1040 for each gene and reported as a full candidate for each family were aligned in MAFFT v7.310
1041 (130) and positions with more than 20% gaps removed with trimal v1.4.rev22 (131) (options -gt
1042 0.8 -st 0). Then their phylogenetic relationships were estimated by maximum likelihood with IQ-
1043 TREE v.2.2.0.3 (132). The best substitution model was selected by ModelFinder (133) using the
1044 Bayesian Information Criterion, by setting “-m MFP” as part of IQ-TREE options. The DNA
1045 substitution models VT+I+G4 and JTT+F+R6 models were selected as the best fitting for the
1046 CREP and FREP sequence alignments, respectively. Node support was estimated with 1000
1047 replicates of non-parametric bootstrap using IQ-TREE options “-b 1000”. Trees were re-rooted at
1048 the midpoint and inspected in iTOL v6 (134).

1049

1050 *Identification of previously identified schistosome resistance genes*

1051

1052 Several candidate immune genes have been identified in *B. glabrata* that could be involved
1053 in the response of snails to *S. mansoni*, inferring resistance (Supplementary Table 2). To test
1054 whether polymorphisms were shared within *B. sudanica* and between species of *Biomphalaria* in
1055 some of these candidate PRRs and immune genes that show the strongest support (BSUD.12903,
1056 *PTC1*, *PTC2*, *RADres*), phylogenies were generated using RAxML with -m GTRCAT
1057 (Stamatakis, 2006).

1058

1059 *Assessment of highly diverse genes and genome regions for novel pathogen recognition receptors*

1060

1061 Intraspecific genomic diversity between the five *B. sudanica* inbred lines was determined
1062 by interrogating polymorphisms between the 5 inbred line vcf file (see above) using statistical
1063 measures in vcftools v0.1.17 (107). High-diversity genes were determined by calculating
1064 nucleotide diversity across both the entire gene (including noncoding regions such as introns and
1065 UTRs) and in coding regions only, determined by the Bs111 nuclear genome annotation, of which
1066 the top 1% were selected for further analysis (Supplementary Figure 5). In addition, we identified
1067 genes occurring in or near diverse windows, as follows. We calculated nucleotide diversity across
1068 the genome in sliding windows of 10 kb (starting every 2.5 kb), 30 kb, (starting every 7.5 kb) or
1069 100 kb (starting every 25 kb), and identified the top 0.1% of 10 kb windows, top 0.3% of 30 kb
1070 windows, and top 1% of 100 kb windows. We then identified all genes that occur within (or overlap
1071 partially with) 100 kb of the midpoint of any of these diverse windows.

1072 For all unique genes identified following these criteria, the largest protein-coding amino
1073 acid sequence for each gene was summarized from the annotation (Supplementary File 7). Where
1074 available (those with matches), transcripts of coding regions from each gene (transcript determined
1075 from Bs111 RNAseq data) were matched with their annotated UniProt description (uniprot.org/,
1076 database as of May 25th, 2022) and InterPro v5.56-89.0 description (51) (Supplementary Table 19
1077 and Supplementary Table 20). All amino acid sequences were further characterized by searching
1078 for orthologous proteins using BLASTp (113) on the NCBI protein database
1079 (ncbi.nlm.nih.gov/protein, database as of October 1st, 2022) (Supplementary Table 21). An e-value
1080 cut-off of 1e-50 was used for amino acid sequence matches. The proteins producing the most
1081 significant alignments (based on the lowest e-value and highest percentage identity from

1082 BLASTp), and those predicted using UniProt and Pfam to the query amino acid sequence were
1083 recorded and used to determine each peptide's function/biological process and key protein
1084 domains if this information was available and informative.

1085 Based on the annotated protein, functional domain presence and structure (i.e. presence of
1086 TMDs), each protein was then placed in immune and non-immune gene related categories, such
1087 as in similar analyses (135,136), and placed under the broader groups of 1. non-immune function
1088 suspected, 2. immune-related function, 3. potentially immune-related function, 4. unknown protein
1089 function but containing TMD(s), 5. unknown protein function because sufficient information could
1090 not be obtained and without TMD.

1091 To create a shortlist (n=20) of candidate PRRs under balancing selection, proteins with the
1092 highest nucleotide diversity, and categorized as either in group 2, 3 or 4 that also contained at least
1093 one transmembrane domain were shortlisted. The genes and in some cases the genomic regions
1094 surrounding these were assessed in more detail, including assessing positions of functional
1095 domains and intra/extracellular regions relative to nucleotide diversity (Supplementary Figure 5).

1096 To determine if any enriched GO terms, particularly those associated with immunity, were
1097 amongst the proteins identified as potential PRRs, topGO v2.48.0 (72) in R v4.3.1 (137) with the
1098 Weight01 algorithm, node size=10 and statistical test of Fisher (significant p-value ≤ 0.05) was
1099 used. The test was performed by comparing the *B. sudanica* whole genome gene family
1100 composition to three separate lists of the highly diverse genes identified: 1) all 1047 highly diverse
1101 genes identified in the *B. sudanica* genome (as listed Supplementary Table 17), 2) 245 genes
1102 classified as immune suspected genes from the highly diverse genes (group 2 in Supplementary
1103 Table 17), and 3) 242 genes shortlisted as the protein containing a transmembrane domain (TMD)
1104 and categorized in group 2, 3 (potential role in innate immunity) or 4 (an unknown function but
1105 contained TMD(s)) (Supplementary Table 17).

1106

1107 *Inferring linkage groups and synteny with Biomphalaria glabrata linkage groups*

1108

1109 We compared synteny and orthology with *B. glabrata* using the iM assembly (6). We
1110 defined orthologous genes as reciprocal best BLASTp hits of protein sequences. We defined
1111 orthologous contigs/scaffolds as those sharing the most orthologous genes. For *B.*
1112 *sudanica* contigs/scaffolds that were orthologous to scaffolds mapped in the 18 *B. glabrata* iM

1113 LGs, we assigned them to 18 LGs corresponding to the linkage groups and ordered them to mirror
1114 the orthologous order. A single iM scaffold (BGM014) is split between linkage groups, so we
1115 considered the distinctly-mapped portions of it separately. For linkage groups of particular interest,
1116 we generated dot plots of sequence similarity based on our previous approach (138). We used these
1117 to refine the order of contigs/scaffolds, to identify contigs/scaffolds with little sequence similarity
1118 despite possessing nominally orthologous genes, and to characterize inversions and other
1119 chromosomal rearrangements.

1120

1121 **Declarations**

1122 The authors have nothing to declare.

1123

1124 **Ethics approval and consent to participate**

1125 Not applicable

1126

1127 **Consent for publication**

1128 Not applicable

1129

1130 **Availability of data and materials**

1131 The datasets generated and/or analyzed during the current study are available in the NCBI
1132 BioProject repository, accession number In prep.

1133 The bioinformatic pipeline for the genome annotation is available at [github.com/J-](https://github.com/J-Calvelo/Annotation-Biomphalaria-sudanica)
1134 [Calvelo/Annotation-Biomphalaria-sudanica](https://github.com/J-Calvelo/Annotation-Biomphalaria-sudanica).

1135

1136 **Competing interests**

1137 Not applicable

1138

1139 **Funding**

1140 Funding for this project was provided by National Institutes of Health, National Institute of Allergy
1141 and Infectious Disease R01AI141862, which supported TP, JAT and MLS. ESL was funded by
1142 NIH grant R37AI101438.

1143

1144 **Authors' contributions**

1145 TP, JCal., JAT, SRB and MLS collected, analyzed and interpreted the majority of the data
1146 contributing to this manuscript. TP, JCal., JAT, AI and MLS were major contributors to writing
1147 the initial manuscript. JCal. developed the bioinformatic pipeline for the annotation of the genome.
1148 JAT developed and conceptualized the bioinformatic pipelines for the analysis and description of
1149 hyperdiverse genes under balancing selection. RB and JCay. collected and analyzed data regarding
1150 the hyperdiverse immune genes in the genome. SRB and MSB generated and assembled/aligned
1151 whole genome sequence data from the inbred lines and assisted in generating transcriptome
1152 sequence data. MLS, TP, JMS, assisted in the maintenance and development of inbred snail lines.
1153 FGH, GO, FR, KA, BM, MOdh. and MOdi. assisted in manuscript preparation and editing. MOdi
1154 assisted in project planning. ESL, MRL and LL assisted in the comparative work with
1155 African *Biomphalaria* species, and LL assisted with the FREPs identification pipeline and
1156 analysis. MLS and MOdi. were responsible for funding acquisition to complete this project. All
1157 authors reviewed and approved the final manuscript.

1158

1159 **Acknowledgements**

1160

1161 The NMRI schistosomes were provided by the NIAID Schistosomiasis Resource Center of
1162 the Biomedical Research Institute (Rockville, MD) through NIH-NIAID Contract
1163 HHSN272201700014I. We also acknowledge Hannah Tavalire for her work in developing the
1164 inbred lines of *B. sudanica*.

1165

1166 **References**

- 1167 1. Hotez PJ, Daar AS. The CNCDs and the NTDs: Blurring the Lines Dividing
1168 Noncommunicable and Communicable Chronic Diseases. *PLoS Negl Trop Dis*.
1169 2008;2(10):e312.
- 1170 2. WHO. Schistosomiasis: Key facts [Internet]. 2023 [cited 2023 Apr 14]. Available from:
1171 <https://www.who.int/news-room/fact-sheets/detail/schistosomiasis>
- 1172 3. Castillo MG, Humphries JE, Mourão MM, Marquez J, Gonzalez A, Montelongo CE.
1173 *Biomphalaria glabrata* immunity: Post-genome advances. *Dev Comp Immunol*.
1174 2020;104:103557.
- 1175 4. Mitta G, Gourbal B, Grunau C, Knight M, Bridger JM, Théron A. The Compatibility
1176 Between *Biomphalaria glabrata* Snails and *Schistosoma mansoni*: An Increasingly
1177 Complex Puzzle. *Adv Parasitol*. 2017;97:111–45.
- 1178 5. Adema CM, Hillier LW, Jones CS, Loker ES, Knight M, Minx P, et al. Whole genome
1179 analysis of a schistosomiasis-transmitting freshwater snail. *Nat Commun*. 2017;8:15451.
- 1180 6. Bu L, Zhong D, Lu L, Loker ES, Yan G, Zhang S-M. Compatibility between snails and
1181 schistosomes: insights from new genetic resources, comparative genomics, and genetic
1182 mapping. *Commun Biol*. 2022;5(1):940.
- 1183 7. Nong W, Yu Y, Aase-Remedios ME, Xie Y, So WL, Li Y, et al. Genome of the ramshorn
1184 snail *Biomphalaria straminea*—an obligate intermediate host of schistosomiasis.
1185 *Gigascience*. 2022;11:giac012.
- 1186 8. Bu L, Lu L, Laidemitt MR, Zhang S-M, Mutuku M, Mkoji G, et al. A genome sequence
1187 for *Biomphalaria pfeifferi*, the major vector snail for the human-infecting parasite
1188 *Schistosoma mansoni*. *PLoS Negl Trop Dis*. 2023 Mar 24;17(3):e0011208.
- 1189 9. DeJong RJ, Morgan JAT, Paraense WL, Pointier J-P, Amarista M, Ayeh-Kumi PFK, et al.
1190 Evolutionary Relationships and Biogeography of *Biomphalaria* (Gastropoda: Planorbidae)
1191 with Implications Regarding Its Role as Host of the Human Bloodfluke, *Schistosoma*
1192 *mansoni*. *Mol Biol Evol*. 2001 Dec 1;18(12):2225–39.
- 1193 10. Campbell G, Jones CS, Lockyer AE, Hughes S, Brown D, Noble LR, et al. Molecular
1194 evidence supports an African affinity of the Neotropical freshwater gastropod,
1195 *Biomphalaria glabrata*, Say 1818, an intermediate host for *Schistosoma mansoni*. *Proc R*
1196 *Soc London Ser B Biol Sci*. 2000;267(1460):2351–8.
- 1197 11. Woodruff DS, Mulvey M. Neotropical schistosomiasis: African affinities of the host snail
1198 *Biomphalaria glabrata* (Gastropoda: Planorbidae). *Biol J Linn Soc*. 1997;60(4):505–16.
- 1199 12. Martens E von. Conchylien aus dem obern Nilgebiet. *Malakozool Blätter*. 1870;17:32–6.
- 1200 13. Brown DS. *Freshwater Snails of Africa and their Medical Importance*. 2nd ed. London,
1201 UK: Taylor & Francis; 1994.
- 1202 14. Erko B, Balcha F, Kifle D. The ecology of *Biomphalaria sudanica* in Lake Ziway,
1203 Ethiopia. *Afr J Ecol*. 2006;44(3):347–52.
- 1204 15. Kazibwe F, Makanga B, Rubaire-Akiiki C, Ouma J, Kariuki C, Kabatereine NB, et al.
1205 Ecology of *Biomphalaria* (Gastropoda: Planorbidae) in Lake Albert, Western Uganda:
1206 snail distributions, infection with schistosomes and temporal associations with
1207 environmental dynamics. *Hydrobiologia*. 2006;568:433–44.
- 1208 16. Williams SN, Hunter PJ. The distribution of *Bulinus* and *Biomphalaria* in Khartoum and
1209 Blue Nile Provinces, Sudan. *Bull World Health Organ*. 1968;39(6):949.
- 1210 17. Loker ES, Moyo HG, Gardner SL. Trematode-gastropod associations in nine non-
1211 lacustrine habitats in the Mwanza region of Tanzania. *Parasitology*. 1981;83(2):381–99.

- 1212 18. Magendantz M. The biology of *Biomphalaria choanomphala* and *B. sudanica* in relation
1213 to their role in the transmission of *Schistosoma mansoni* in Lake Victoria at Mwanza,
1214 Tanzania. Bull World Health Organ. 1972;47(3):331.
- 1215 19. Platt RN, Le Clec'h W, Chevalier FD, McDew-White M, LoVerde PT, Ramiro de Assis
1216 R, et al. Genomic analysis of a parasite invasion: Colonization of the Americas by the
1217 blood fluke *Schistosoma mansoni*. Mol Ecol. 2022;31(8):2242–63.
- 1218 20. Morgan JAT, DeJong RJ, Adeoye GO, Ansa EDO, Barbosa CS, Brémond P, et al. Origin
1219 and diversification of the human parasite *Schistosoma mansoni*. Mol Ecol.
1220 2005;14(12):3889–902.
- 1221 21. King CH, Binder S, Shen Y, Whalen CC, Campbell CH, Wiegand RE, et al. SCORE
1222 Studies on the Impact of Drug Treatment on Morbidity due to *Schistosoma mansoni* and
1223 *Schistosoma haematobium* Infection. Am J Trop Med Hyg. 2020;103(1):30–5.
- 1224 22. Secor WE, Wiegand RE, Montgomery SP, Karanja DMS, Odiere MR. Comparison of
1225 school-based and community-wide mass drug administration for schistosomiasis control
1226 in an area of western Kenya with high initial *Schistosoma mansoni* infection prevalence: a
1227 cluster randomized trial. Am J Trop Med Hyg. 2020;102(2):318.
- 1228 23. Wiegand RE, Mwinzi PNM, Montgomery SP, Chan YL, Andiego K, Omedo M, et al. A
1229 persistent hotspot of *Schistosoma mansoni* infection in a five-year randomized trial of
1230 praziquantel preventative chemotherapy strategies. J Infect Dis. 2017;216(11):1425–33.
- 1231 24. Gouvras AN, Allan F, Kinung'Hi S, Rabone M, Emery A, Angelo T, et al. Longitudinal
1232 survey on the distribution of *Biomphalaria sudanica* and *B. choanomphala* in Mwanza
1233 region, on the shores of Lake Victoria, Tanzania: Implications for schistosomiasis
1234 transmission and control. Parasites and Vectors. 2017;10(1).
- 1235 25. Martens E von. Recente Conchylien aus dem Victoria Nianza (Ukerewe). Sitzungsberichte
1236 der Gesellschaft Naturforschender Freunde zu Berlin. 1879;103–5.
- 1237 26. Loker ES, Hofkin B V, Mkoji GM, Mungai B, Kihara JH, Koech DK. Distributions of
1238 freshwater snails in southern Kenya with implications for the biological control of
1239 schistosomiasis and other snail-mediated parasites. J Med Appl Malacol. 1993;5:1–20.
- 1240 27. Mutuku MW, Laidemitt MR, Beechler BR, Mwangi IN, Otiato FO, Agola EL, et al. A
1241 Search for Snail-Related Answers to Explain Differences in Response of *Schistosoma*
1242 *mansoni* to Praziquantel Treatment among Responding and Persistent Hotspot Villages
1243 along the Kenyan Shore of Lake Victoria. Am J Trop Med Hyg. 2019;101(1):65–77.
- 1244 28. Standley CJ, Wade CM, Stothard JR. A fresh insight into transmission of schistosomiasis:
1245 a misleading tale of *Biomphalaria* in Lake Victoria. PLoS One. 2011;6(10):e26563.
- 1246 29. Andrus PS, Stothard JR, Kabatereine NB, Wade CM. Comparing shell size and shape with
1247 canonical variate analysis of sympatric *Biomphalaria* species within Lake Albert and Lake
1248 Victoria, Uganda. Zool J Linn Soc. 2023 Jul 21;zlzd052.
- 1249 30. Standley CJ, Goodacre SL, Wade CM, Stothard JR. The population genetic structure of
1250 *Biomphalaria choanomphala* in Lake Victoria, East Africa: implications for
1251 schistosomiasis transmission. Parasit Vectors. 2014;7(1):524.
- 1252 31. Mutuku MW, Laidemitt MR, Spaan JM, Mwangi IN, Ochanda H, Steinauer ML, et al.
1253 Comparative vectorial competence of *Biomphalaria sudanica* and *Biomphalaria*
1254 *choanomphala*, snail hosts of *Schistosoma mansoni*, from transmission hotspots In Lake
1255 Victoria, Western Kenya. J Parasitol. 2021;107(2):349–57.
- 1256 32. Shultz AJ, Sackton TB. Immune genes are hotspots of shared positive selection across
1257 birds and mammals. Elife. 2019;8:e41815.

- 1258 33. Leffler EM, Gao Z, Pfeifer S, Ségurel L, Auton A, Venn O, et al. Multiple Instances of
1259 Ancient Balancing Selection Shared Between Humans and Chimpanzees. *Science*.
1260 2013;339(6127):1578–82.
- 1261 34. Takahata N, Satta Y, Klein J. Polymorphism and balancing selection at major
1262 histocompatibility complex loci. *Genetics*. 1992;130(4):925–38.
- 1263 35. Hughes AL, Yeager M. Natural selection at major histocompatibility complex loci of
1264 vertebrates. *Annu Rev Genet*. 1998;32(1):415–35.
- 1265 36. Ferrer-Admetlla A, Bosch E, Sikora M, Marquès-Bonet T, Ramírez-Soriano A, Muntasell
1266 A, et al. Balancing Selection Is the Main Force Shaping the Evolution of Innate Immunity
1267 Genes. *J Immunol*. 2008;181(2):1315–22.
- 1268 37. Chisholm ST, Coaker G, Day B, Staskawicz BJ. Host-Microbe Interactions: Shaping the
1269 Evolution of the Plant Immune Response. *Cell*. 2006;124(4):803–14.
- 1270 38. Van de Weyer A-L, Monteiro F, Furzer OJ, Nishimura MT, Cevik V, Witek K, et al. A
1271 Species-Wide Inventory of NLR Genes and Alleles in *Arabidopsis thaliana*. *Cell*.
1272 2019;178(5):1260-1272.e14.
- 1273 39. Lee D, Zdraljevic S, Stevens L, Wang Y, Tanny RE, Crombie TA, et al. Balancing
1274 selection maintains hyper-divergent haplotypes in *Caenorhabditis elegans*. *Nat Ecol Evol*.
1275 2021;5(6):794–807.
- 1276 40. Tennessen JA, Théron A, Marine M, Yeh J-Y, Rognon A, Blouin MS. Hyperdiverse Gene
1277 Cluster in Snail Host Conveys Resistance to Human Schistosome Parasites. *PLoS Genet*.
1278 2015;11(3):1–21.
- 1279 41. Tennessen JA, Bollmann SR, Peremyslova E, Kronmiller BA, Sergi C, Hamali B, et al.
1280 Clusters of polymorphic transmembrane genes control resistance to schistosomes in snail
1281 vectors. *Elife*. 2020;9:e59395.
- 1282 42. Adema CM, Hertel LA, Miller RD, Loker ES. A family of fibrinogen-related proteins that
1283 precipitates parasite-derived molecules is produced by an invertebrate after infection. *Proc*
1284 *Natl Acad Sci*. 1997;94(16):8691–6.
- 1285 43. Lu L, Loker ES, Adema CM, Zhang S-M, Bu L. Genomic and transcriptional analysis of
1286 genes containing fibrinogen and IgSF domains in the schistosome vector *Biomphalaria*
1287 *glabrata*, with emphasis on the differential responses of snails susceptible or resistant to
1288 *Schistosoma mansoni*. *PLoS Negl Trop Dis*. 2020;14(10):1–35.
- 1289 44. Fukasawa Y, Ermini L, Wang H, Carty K, Cheung M-S. LongQC: A Quality Control Tool
1290 for Third Generation Sequencing Long Read Data. *G3 Genes, Genomes, Genet*.
1291 2020;10(4):1193–6.
- 1292 45. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
1293 data. *Bioinformatics*. 2014;30(15):2114–20.
- 1294 46. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.
1295 2018;34(18):3094–100.
- 1296 47. Li H. New strategies to improve minimap2 alignment accuracy. *Bioinformatics*.
1297 2021;37(23):4572–4.
- 1298 48. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
1299 universal RNA-seq aligner. *Bioinformatics*. 2013;29(1):15–21.
- 1300 49. Shumate A, Wong B, Perteau G, Perteau M. Improved transcriptome assembly using a
1301 hybrid of long and short reads with StringTie. *PLOS Comput Biol*. 2022;18(6):1–18.
- 1302 50. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De novo
1303 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference

- 1304 generation and analysis. *Nat Protoc.* 2013;8(8):1494–512.
- 1305 51. Paysan-Lafosse T, Blum M, Chuguransky S, Grego T, Pinto BL, Salazar GA, et al.
- 1306 InterPro in 2022. *Nucleic Acids Res.* 2023;51(D1):D418–27.
- 1307 52. Manni M, Berkeley MR, Seppy M, Zdobnov EM. BUSCO: Assessing Genomic Data
- 1308 Quality and Beyond. *Curr Protoc.* 2021;1(12):e323.
- 1309 53. Serrão VHB, Silva IR, da Silva MTA, Scortecci JF, de Freitas Fernandes A, Thiemann
- 1310 OH. The unique tRNAs^{Sec} and its role in selenocysteine biosynthesis. *Amino Acids.*
- 1311 2018;50(9):1145–67.
- 1312 54. Commans S, Böck A. Selenocysteine inserting tRNAs: an overview. *FEMS Microbiol*
- 1313 *Rev.* 1999;23(3):335–51.
- 1314 55. Zhang S-M, Bu L, Laidemitt MR, Lu L, Mutuku MW, Mkoji GM, et al. Complete
- 1315 mitochondrial and rDNA complex sequences of important vector species of *Biomphalaria*,
- 1316 obligatory hosts of the human-infecting blood fluke, *Schistosoma mansoni*. *Sci Rep.*
- 1317 2018;8(1):1–10.
- 1318 56. Grande C, Templado J, Zardoya R. Evolution of gastropod mitochondrial genome
- 1319 arrangements. *BMC Evol Biol.* 2008;8(1):61.
- 1320 57. Ojala D, Montoya J, Attardi G. tRNA punctuation model of RNA processing in human
- 1321 mitochondria. *Nature.* 1981;290(5806):470–4.
- 1322 58. Ghiselli F, Gomes-dos-Santos A, Adema CM, Lopes-Lima M, Sharbrough J, Boore JL.
- 1323 Molluscan mitochondrial genomes break the rules. *Philos Trans R Soc B Biol Sci.*
- 1324 2021;376(1825):20200159.
- 1325 59. DeJong RJ, Emery AM, Adema CM. The mitochondrial genome of *Biomphalaria*
- 1326 *glabrata* (Gastropoda: Basommatophora), intermediate host of *Schistosoma mansoni*. *J*
- 1327 *Parasitol.* 2004;90(5):991–8.
- 1328 60. Bernt M, Donath A, Jühling F, Externbrink F, Florentz C, Fritsch G, et al. MITOS:
- 1329 Improved *de novo* metazoan mitochondrial genome annotation. *Mol Phylogenet Evol.*
- 1330 2013;69(2):313–9.
- 1331 61. Teufel F, Almagro Armenteros JJ, Johansen AR, Gíslason MH, Pihl SI, Tsirigos KD, et al.
- 1332 SignalP 6.0 predicts all five types of signal peptides using protein language models. *Nat*
- 1333 *Biotechnol.* 2022;40(7):1023–5.
- 1334 62. Almagro Armenteros JJ, Salvatore M, Emanuelsson O, Winther O, von Heijne G,
- 1335 Elofsson A, et al. Detecting sequence signals in targeting peptides using deep learning.
- 1336 *Life Sci Alliance.* 2019;2(5):e201900429.
- 1337 63. Bendtsen JD, Jensen LJ, Blom N, von Heijne G, Brunak S. Feature-based prediction of
- 1338 non-classical and leaderless protein secretion. *Protein Eng Des Sel.* 2004;17(4):349–56.
- 1339 64. Tress ML, Martelli PL, Frankish A, Reeves GA, Wesselink JJ, Yeats C, et al. The
- 1340 implications of alternative splicing in the ENCODE protein complement. *Proc Natl Acad*
- 1341 *Sci.* 2007;104(13):5495–500.
- 1342 65. Wang T, Wyeth RC, Liang D, Bose U, Ni G, McManus DP, et al. A *Biomphalaria*
- 1343 *glabrata* peptide that stimulates significant behaviour modifications in aquatic free-living
- 1344 *Schistosoma mansoni* miracidia. *PLoS Negl Trop Dis.* 2019;13(1):e0006948.
- 1345 66. Fogarty CE, Phan P, Duke MG, McManus DP, Wyeth RC, Cummins SF, et al.
- 1346 Identification of *Schistosoma mansoni* miracidia attractant candidates in infected
- 1347 *Biomphalaria glabrata* using behaviour-guided comparative proteomics. Vol. 13,
- 1348 *Frontiers in Immunology.* 2022.
- 1349 67. Emms DM, Kelly S. OrthoFinder: phylogenetic orthology inference for comparative

- 1350 genomics. *Genome Biol.* 2019;20(1):238.
- 1351 68. Mendes FK, Vanderpool D, Fulton B, Hahn MW. CAFE 5 models variation in
1352 evolutionary rates among gene families. *Bioinformatics.* 2020;36(22–23):5516–8.
- 1353 69. Pickford M. Freshwater and terrestrial Mollusca from the Early Miocene deposits of the
1354 northern Sperrgebiet, Namibia. *Mem Geol Surv Namibia.* 2008;20:65–74.
- 1355 70. Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J. eggNOG-
1356 mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the
1357 Metagenomic Scale. *Mol Biol Evol.* 2021;38(12):5825–9.
- 1358 71. Huerta-Cepas J, Szklarczyk D, Heller D, Hernández-Plaza A, Forslund SK, Cook H, et al.
1359 eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology
1360 resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.*
1361 2019;47(D1):D309–14.
- 1362 72. Alexa A, Rahnenfuhrer J. topGO: Enrichment Analysis for Gene Ontology. R package
1363 version 2.48.0. 2022.
- 1364 73. Supek F, Bošnjak M, Škunca N, Šmuc T. REVIGO Summarizes and Visualizes Long
1365 Lists of Gene Ontology Terms. *PLoS One.* 2011;6(7):e21800.
- 1366 74. Emms DM, Kelly S. STAG: Species Tree Inference from All Genes. *bioRxiv.*
1367 2018;267914.
- 1368 75. Eddy SR. Accelerated Profile HMM Searches. *PLOS Comput Biol.* 2011;7(10):e1002195.
- 1369 76. Dheilly NM, Duval D, Mouahid G, Emans R, Allienne J-F, Galinier R, et al. A family of
1370 variable immunoglobulin and lectin domain containing molecules in the snail
1371 *Biomphalaria glabrata*. *Dev Comp Immunol.* 2015;48(1):234–43.
- 1372 77. Spaan JM, Pennance T, Laidemitt MR, Sims N, Roth J, Lam Y, et al. Multi-strain
1373 compatibility polymorphism between a parasite and its snail host, a neglected vector of
1374 schistosomiasis in Africa. *Curr Res Parasitol Vector-Borne Dis.* 2023;3:100120.
- 1375 78. Tennessen JA, Bonner KM, Bollmann SR, Johnstun JA, Yeh J-Y, Marine M, et al.
1376 Genome-Wide Scan and Test of Candidate Genes in the Snail *Biomphalaria glabrata*
1377 Reveal New Locus Influencing Resistance to *Schistosoma mansoni*. *PLoS Negl Trop Dis.*
1378 2015;9(9):1–19.
- 1379 79. Allan ERO, Tennessen JA, Bollmann SR, Hanington PC, Bayne CJ, Blouin MS.
1380 Schistosome infectivity in the snail, *Biomphalaria glabrata* is partially dependent on the
1381 expression of Grctm6, a Guadeloupe Resistance Complex protein. *PLoS Negl Trop Dis.*
1382 2017;11(2):1–15.
- 1383 80. Goel S, Palmkvist M, Moll K, Joannin N, Lara P, R Akhouri R, et al. RIFINs are adhesins
1384 implicated in severe *Plasmodium falciparum* malaria. *Nat Med.* 2015;21(4):314–7.
- 1385 81. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
1386 large phylogenies. *Bioinformatics.* 2014 Jan 21;30(9):1312–3.
- 1387 82. Maier T, Wheeler NJ, Namigai EKO, Tycko J, Grewelle RE, Woldeamanuel Y, et al.
1388 Gene drives for schistosomiasis transmission control. *PLoS Negl Trop Dis.*
1389 2019;13(12):1–21.
- 1390 83. WHO. WHO guideline on control and elimination of human schistosomiasis. 2022.
- 1391 84. Mutuku MW, Lu L, Otiato FO, Mwangi IN, Kinuthia JM, Maina GM, et al. A
1392 Comparison of Kenyan *Biomphalaria pfeifferi* and *B. sudanica* as Vectors for *Schistosoma*
1393 *mansoni*, Including a Discussion of the Need to Better Understand the Effects of Snail
1394 Breeding Systems on Transmission. *J Parasitol.* 2017;103(6):669–76.
- 1395 85. Lu L, Zhang S-M, Mutuku MW, Mkoji GM, Loker ES. Relative compatibility of

- 1396 *Schistosoma mansoni* with *Biomphalaria sudanica* and *B. pfeifferi* from Kenya as assessed
1397 by PCR amplification of the *S. mansoni* ND5 gene in conjunction with traditional
1398 methods. Parasit Vectors. 2016;9(1):166.
- 1399 86. Moné Y, Gourbal B, Duval D, Du Pasquier L, Kieffer-Jaquinod S, Mitta G. A Large
1400 Repertoire of Parasite Epitopes Matched by a Large Repertoire of Host Immune Receptors
1401 in an Invertebrate Host/Parasite Model. PLoS Negl Trop Dis. 2010;4(9):e813.
- 1402 87. Larson MK, Bender RC, Bayne CJ. Resistance of *Biomphalaria glabrata* 13-16-R1 snails
1403 to *Schistosoma mansoni* PR1 is a function of haemocyte abundance and constitutive levels
1404 of specific transcripts in haemocytes. Int J Parasitol. 2014;44(6):343–53.
- 1405 88. Sun Y, Zhou Z, Wang L, Yang C, Jianga S, Song L. The immunomodulation of a novel
1406 tumor necrosis factor (CgTNF-1) in oyster *Crassostrea gigas*. Dev Comp Immunol.
1407 2014;45(2):291–9.
- 1408 89. Zheng Y, Liu Z, Wang L, Li M, Zhang Y, Zong Y, et al. A novel tumor necrosis factor in
1409 the Pacific oyster *Crassostrea gigas* mediates the antibacterial response by triggering the
1410 synthesis of lysozyme and nitric oxide. Fish Shellfish Immunol. 2020;98:334–41.
- 1411 90. Huang Y, Si Q, Du S, Du J, Ren Q. Molecular identification and functional analysis of a
1412 tumor necrosis factor superfamily gene from Chinese mitten crab (*Eriocheir sinensis*).
1413 Dev Comp Immunol. 2022;134:104456.
- 1414 91. Blouin MS, Bonner KM, Cooper B, Amarasinghe V, O'Donnell RP, Bayne CJ. Three
1415 genes involved in the oxidative burst are closely linked in the genome of the snail,
1416 *Biomphalaria glabrata*. Int J Parasitol. 2013;43(1):51–5.
- 1417 92. Tennessen JA, Bollmann SR, Blouin MS. A Targeted Capture Linkage Map Anchors the
1418 Genome of the Schistosomiasis Vector Snail, *Biomphalaria glabrata*. G3 Genes,
1419 Genomes, Genet. 2017;7(7):2353–61.
- 1420 93. Mimpfound R, Greer GJ. Allozyme variation among populations of *Biomphalaria pfeifferi*
1421 (Krauss, 1848) (Gastropoda: Planorbidae) in Cameroon. J Molluscan Stud.
1422 1990;56(4):461–7.
- 1423 94. Tian-Bi Y-N, Jarne P, Konan J-N, Utzinger J, N'Goran EK. Contrasting the distribution of
1424 phenotypic and molecular variation in the freshwater snail *Biomphalaria pfeifferi*, the
1425 intermediate host of *Schistosoma mansoni*. Heredity (Edinb). 2013;110(5):466–74.
- 1426 95. Baclaocos J, Santesmasses D, Mariotti M, Bierła K, Vetick MB, Lynch S, et al. Processive
1427 Recoding and Metazoan Evolution of Selenoprotein P: Up to 132 UGAs in Molluscs. J
1428 Mol Biol. 2019;431(22):4381–407.
- 1429 96. Budachetri K, Karim S. An insight into the functional role of thioredoxin reductase, a
1430 selenoprotein, in maintaining normal native microbiota in the Gulf Coast tick
1431 (*Amblyomma maculatum*). Insect Mol Biol. 2015;24(5):570–81.
- 1432 97. Gao S, Ren Y, Sun Y, Wu Z, Ruan J, He B, et al. PacBio full-length transcriptome
1433 profiling of insect mitochondrial gene expression. RNA Biol. 2016;13(9):820–5.
- 1434 98. Fourdrilis S, de Frias Martins AM, Backeljau T. Relation between mitochondrial DNA
1435 hyperdiversity, mutation rate and mitochondrial genome evolution in *Melarhaphe*
1436 *neritoides* (Gastropoda: Littorinidae) and other Caenogastropoda. Sci Rep. 2018 Dec
1437 19;8(1):17964.
- 1438 99. Yévenes M, Núñez-Acuña G, Gallardo-Escárate C, Gajardo G. Adaptive mitochondrial
1439 genome functioning in ecologically different farm-impacted natural seedbeds of the
1440 endemic blue mussel *Mytilus chilensis*. Comp Biochem Physiol Part D Genomics
1441 Proteomics. 2022;42:100955.

- 1442 100. Zhang S-M, Bu L, Lu L, Babbitt C, Adema CM, Loker ES. Comparative mitogenomics of
1443 freshwater snails of the genus *Bulinus* obligatory vectors of *Schistosoma haematobium*,
1444 causative agent of human urogenital schistosomiasis. *Sci Rep.* 2022;12(1):5357.
- 1445 101. Wenger AM, Peluso P, Rowell WJ, Chang P-C, Hall RJ, Concepcion GT, et al. Accurate
1446 circular consensus long-read sequencing improves variant detection and assembly of a
1447 human genome. *Nat Biotechnol.* 2019;37(10):1155–62.
- 1448 102. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using
1449 repeat graphs. *Nat Biotechnol.* 2019;37(5):540–6.
- 1450 103. Shen W, Le S, Li Y, Hu F. SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q
1451 File Manipulation. *PLoS One.* 2016;11(10):e0163962.
- 1452 104. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
1453 *EMBnet J.* 2011;17(1):10–2.
- 1454 105. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.
1455 *Bioinformatics.* 2009;25(14):1754–60.
- 1456 106. Li H. A statistical framework for SNP calling, mutation discovery, association mapping
1457 and population genetical parameter estimation from sequencing data. *Bioinformatics.*
1458 2011;27(21):2987–93.
- 1459 107. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant
1460 call format and VCFtools. *Bioinformatics.* 2011;27(15):2156–8.
- 1461 108. Lewis FA, Liang Y-S, Raghavan N, Knight M. The NIH-NIAID schistosomiasis resource
1462 center. *PLoS Negl Trop Dis.* 2008;2(7):e267–e267.
- 1463 109. Andrews S. FastQC: A Quality Control Tool for High Throughput Sequence Data
1464 [Internet]. 2010. Available from:
1465 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- 1466 110. Barnett DW, Garrison EK, Quinlan AR, Strömberg MP, Marth GT. BamTools: a C++ API
1467 and toolkit for analyzing and managing BAM files. *Bioinformatics.* 2011;27(12):1691–2.
- 1468 111. The UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. *Nucleic
1469 Acids Res.* 2020;49(D1):D480–9.
- 1470 112. Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, et al.
1471 Pfam: The protein families database in 2021. *Nucleic Acids Res.* 2020;49(D1):D412–9.
- 1472 113. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search
1473 tool. *J Mol Biol.* 1990;215(3):403–10.
- 1474 114. Wheeler TJ, Eddy SR. nhmmer: DNA homology search with profile HMMs.
1475 *Bioinformatics.* 2013;29(19):2487–9.
- 1476 115. Baril T, Imrie RM, Hayward A. Earl Grey: a fully automated user-friendly transposable
1477 element annotation and analysis pipeline. *bioRxiv.* 2022;2022.06.30.498289.
- 1478 116. Smit A, Hubley R, Green P. RepeatMasker Open-4.0. 2013.
- 1479 117. Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. The Dfam community resource of
1480 transposable element families, sequence models, and genome annotations. *Mob DNA.*
1481 2021;12(1):2.
- 1482 118. Smit A, Hubley R. RepeatModeler Open-1.0. 2008.
- 1483 119. Bao Z, Eddy SR. Automated de novo identification of repeat sequence families in
1484 sequenced genomes. *Genome Res.* 2002;12(8):1269–76.
- 1485 120. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large
1486 genomes. *Bioinformatics.* 2005;21(suppl_1):i351–8.
- 1487 121. Wickham H. ggplot2. *Wiley Interdiscip Rev Comput Stat.* 2011;3(2):180–5.

- 1488 122. Chan PP, Lin BY, Mak AJ, Lowe TM. tRNAscan-SE 2.0: improved detection and
1489 functional classification of transfer RNA genes. *Nucleic Acids Res.* 2021;49(16):9077–96.
- 1490 123. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al.
1491 Integrative genomics viewer. *Nat Biotechnol.* 2011;29(1):24–6.
- 1492 124. Putri GH, Anders S, Pyl PT, Pimanda JE, Zanini F. Analysing high-throughput
1493 sequencing data in Python with HTSeq 2.0. *Bioinformatics.* 2022;38(10):2943–5.
- 1494 125. Young ND, Stroehlein AJ, Wang T, Korhonen PK, Mentink-Kane M, Stothard JR, et al.
1495 Nuclear genome of *Bulinus truncatus*, an intermediate host of the carcinogenic human
1496 blood fluke *Schistosoma haematobium*. *Nat Commun.* 2022;13(1):977.
- 1497 126. Maeda T, Takahashi S, Yoshida T, Shimamura S, Takaki Y, Nagai Y, et al. Chloroplast
1498 acquisition without the gene transfer in kleptoplastic sea slugs, *Plakobranthus ocellatus*.
1499 *Elife.* 2021;10:e60176.
- 1500 127. Rice P, Longden I, Bleasby A. EMBOSS: The European Molecular Biology Open
1501 Software Suite. *Trends Genet.* 2000;16(6):276–7.
- 1502 128. Laidemitt MR, Buddenborg SK, Lewis LL, Michael LE, Sanchez MJ, Hewitt R, et al.
1503 *Schistosoma mansoni* Vector Snails in Antigua and Montserrat, with Snail-Related
1504 Considerations Pertinent to a Declaration of Elimination of Human Schistosomiasis. *Am J*
1505 *Trop Med Hyg.* 2020;103(6):2268–77.
- 1506 129. Hallgren J, Tsigos KD, Pedersen MD, Almagro Armenteros JJ, Marcatili P, Nielsen H,
1507 et al. DeepTMHMM predicts alpha and beta transmembrane proteins using deep neural
1508 networks. *bioRxiv.* 2022;2022.04.08.487609.
- 1509 130. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
1510 improvements in performance and usability. *Mol Biol Evol.* 2013;30(4):772–80.
- 1511 131. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated
1512 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics.*
1513 2009;25(15):1972–3.
- 1514 132. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: A Fast and Effective
1515 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol.*
1516 2015;32(1):268–74.
- 1517 133. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermini LS. ModelFinder:
1518 fast model selection for accurate phylogenetic estimates. *Nat Methods.* 2017;14(6):587–9.
- 1519 134. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new
1520 developments. *Nucleic Acids Res.* 2019;47(W1):W256–9.
- 1521 135. Dinguirard N, Cavalcanti MGS, Wu X-J, Bickham-Wright U, Sabat G, Yoshino TP.
1522 Proteomic Analysis of *Biomphalaria glabrata* Hemocytes During in vitro Encapsulation
1523 of *Schistosoma mansoni* Sporocysts. *Front Immunol.* 2018;9:2773.
- 1524 136. Wu X-J, Dinguirard N, Sabat G, Lui H, Gonzalez L, Gehring M, et al. Proteomic analysis
1525 of *Biomphalaria glabrata* plasma proteins with binding affinity to those expressed by
1526 early developing larval *Schistosoma mansoni*. *PLOS Pathog.* 2017;13(5):1–30.
- 1527 137. R Core Team. R: A Language and Environment for Statistical Computing. Vienna,
1528 Austria; 2018.
- 1529 138. Blouin MS, Bollmann SR, Tennessen JA. *PTC2* region genotypes counteract
1530 *Biomphalaria glabrata* population differences between M-line and BS90 in resistance to
1531 infection by *Schistosoma mansoni*. *PeerJ.* 2022;10:e13971.
- 1532

1533 **Supplementary Material**

1534

1535 **Supplementary Figures**

1536 Supplementary Figure 1. Repetitive elements in the *Biomphalaria sudanica* genome. ^ “Other”
1537 indicates a Simple Repeat, Microsatellite, RNA (.jpg file).

1538

1539 Supplementary Figure 2. Schematic presentation of the mitochondrial gene trimming process for
1540 the transcripts, showing regions of primary transcription (Fragment 1-6) on the plus and minus
1541 strand, and the trimming processes of these primary transcripts into pre-mRNA. Numbers above
1542 transcripts represent the RNA sequence depth from aligned PacBio IsoSeq ccs data (.pdf file).

1543

1544 Supplementary Figure 3. Dot plots of *Biomphalaria sudanica* linkage groups 6 (A), 10 (B) and 16
1545 (C), composed of multiple scaffolds determined using the *B. glabrata* iM line linkage map (Bu et
1546 al., 2022). Dots represent 600bp segments; dark blue is $\geq 97.5\%$ sequence similarity, light blue is
1547 $\geq 90\%$ sequence similarity (.pdf file).

1548

1549 Supplementary Figure 4. Mitochondrial genome of *Biomphalaria sudanica* with point of origin set
1550 to the start of the *nad5* gene (.pdf file).

1551

1552 Supplementary Figure 5. Schematic overview of the methods employed to delimit pathogen
1553 recognition receptors (PRR) and other immune genes of *Biomphalaria sudanica* under balancing
1554 selection, determined through the analysis of high intraspecific genetic diversity regions,
1555 potentially relevant to the resistance and susceptibility of this species to *Schistosoma mansoni*.

1556

1557 **Supplementary Tables**

1558 Supplementary Table 1. Summary metrics and NCBI identifiers for the generated genome (PacBio
1559 long-read) and transcriptome (PacBio long-read and Illumina short-read) sequences, and the
1560 complete annotation statistics (.xlsx file).

1561

1562 Supplementary Table 2. *Biomphalaria sudanica* orthologues to genes / genomic markers of 18
1563 candidate immune loci that have been identified for their involvement in *Schistosoma mansoni*

1564 immunity in *B. glabrata*. *The Knight marker, qRS-5.1, and qRS-2.1, truly represent a marker and
1565 not a coding gene, hence putative function and mechanism are unknown (.xlsx file).

1566

1567 Supplementary Table 3. Summary table of the 919 tRNA genes that were predicted by tRNAscan-
1568 SE 2.0 v2.0.9 (122) representing 24 tRNA types annotated in the *Biomphalaria sudanica* genome
1569 (.xlsx file).

1570

1571 Supplementary Table 4. Complete list showing the genome positions and prediction scores of 919
1572 tRNA genes predicted in the *Biomphalaria sudanica* genome by tRNAscan-SE 2.0 v2.0.9 (122)
1573 (.xlsx file).

1574

1575 Supplementary Table 5. Summary table of repeated elements composition of the *Biomphalaria*
1576 *sudanica* genome (.xlsx file).

1577

1578 Supplementary Table 6. Predicted coding genes that overlap with repeated elements. For each one
1579 it is specified the contig/scaffold, gene and coordinates, and each overlapping repeated element
1580 and the features they compromise: Intron region, 5'UTR, 3'UTR and/or CDS (.xlsx file).

1581

1582 Supplementary Table 7. Mitochondrial genome annotation of *Biomphalaria sudanica* and RNAseq
1583 read counts from the Illumina short-read datasets. For each feature, detail is given on the
1584 coordinates, encoding strand, source of the annotation and if their boundaries were manually
1585 adjusted based on the read mappings. The starting coordinate was selected to be the start of *nad5*
1586 in order to be consistent with the already published mitochondrial genomes of *Biomphalaria* (.xlsx
1587 file).

1588

1589 Supplementary Table 8. Summary table of the predicted location signals found in the
1590 *Biomphalaria sudanica* protein coding genes using the programs SignalP v6.0 (61) for the
1591 identification of Signal Peptides (SP), TargetP v2.0 (62) for the identification of mitochondrial
1592 transit peptide (mTP) and SecretomeP v1.0 (63) for proteins putatively secreted through other
1593 pathways. In addition to the total number of trans-membrane domains predicted by DeepTMHMM
1594 v1.0.13 (129) and the identified InterPro v5.56-89.0 (51) signals. Each isoform was putatively

1595 classified into 4 types: Secreted (SP identified or positive result with SecretomeP and no trans-
1596 membrane domains), Mitochondrial Cytoplasmic (mTP identified and no trans-membrane
1597 domains), Membrane (SP identified or positive result with SecretomeP and at least one trans-
1598 membrane domain), Mitochondrial Membrane (mTP identified and at least one trans-membrane
1599 domain). When either an SP or a mTP was identified, its location and probability is included in
1600 the column Location Signal Details (.xlsx file).

1601
1602 Supplementary Table 9. Phylogenetic Hierarchical Orthogroups (HOG) estimated by Orthofinder
1603 (67). For each HOG, detail is provided on their ID, the original Orthogroup ID and the node in the
1604 gene tree from which the HOG was determined together with the protein ID for the eight genomes
1605 analyzed: *B. glabrata* strain BB02 (GCF_000457365.2 (5)), *B. glabrata* strain iBS90
1606 (GCA_025434165.1 (6)), *B. glabrata* strain iM (GCA_025434175.1 (6)), and the species *B.*
1607 *pfeifferi* (GCA_030265305.1 (8)), *B. straminea* (GCA_021533235.1 (7)), with the planorbid snail
1608 *Bulinus truncatus* (GCA_021962125.1 (125)) and the Plakobranichidae *Elysia marginata*
1609 (GCA_019649035.1 (126)) (.xlsx file).

1610
1611 Supplementary Table 10. Gene Ontology (GO) terms assigned to each Phylogenetic Hierarchical
1612 Orthogroups (HOG) by eggNOG-mapper (70,71) (.xlsx file).

1613
1614 Supplementary Table 11. Enriched gene ontology (GO) terms found in the protein families that
1615 significantly expanded or contracted in the CAFE 5 (68) analysis. At each node it is specified the
1616 IDs and total count of the Phylogenetic Hierarchical Orthogroups (HOG) identified to be either
1617 expanded or contracted, followed by the enrichment results produced by v2.48.0 (72): the GO
1618 terms enriched, the total number of annotated HOGs with each GO terms, how many GO terms
1619 were observed, and the total number of GO terms expected given the list size and the p-value (.xlsx
1620 file).

1621
1622 Supplementary Table 12. Summary of the enriched gene ontology (GO) terms found in the protein
1623 families that significantly expanded or contracted in the CAFE 5 (68) analysis. For each group in
1624 the REVIGO (73) TreeMap summary (see Supplementary File 5) it is specified their member GO
1625 terms ID description, and the semantic similarity measurements calculated by REVIGO:

1626 Frequency of the term in the UniProt database (uniport.org/, database as of May 25th 2022), it's
1627 semantical Uniqueness in the whole list and their Dispensability if the term was nested with
1628 another GO term. In the latter case (dispensable terms), the GO term number it is nested with is
1629 reported, or displayed as “-” otherwise (.xlsx file).

1630

1631 Supplementary Table 13. *Biomphalaria sudanica* protein domain location and Signal Peptides
1632 (SP) found among the initial candidates identified for the protein families C-type lectin-related
1633 protein (CREP), Fibrinogen-related protein (FREP) or galectin-related protein (GREP),
1634 determined by the presence of C-type lectin, fibrinogen (FBD) and Galectin like domains,
1635 respectively. For each examined protein it is detailed to which of the three families it is considered
1636 a candidate for, and the details of each domain and SP, the analysis that identified it, the signature
1637 accession, score (e-value or probability depending on the analysis, domain coordinates, InterPro
1638 v5.56-89.0 (51) annotation and description (if applicable), and what type of signature it was
1639 classified as: C-type lectin like domain, FBD like domain, Galectin like domain, IgSF like domain,
1640 other immunoglobulin like domain or a secreted protein. If the protein was utilized in the
1641 phylogenetic analysis it is indicated with an “X” in the Final Selection column, and “.” otherwise
1642 (.xlsx file).

1643

1644 Supplementary Table 14. Summary table of *Biomphalaria sudanica* C-type lectin-related proteins
1645 (CREP) and Fibrinogen-related proteins (FREP) annotated, providing protein ID, CREP/FREP
1646 subfamily, protein length, signal peptide, immunoglobulin domain (IgSF) and the c-lectin domain
1647 (CREP) or fibrinogen (FBD) domain (FREP) positions, notes on each genes features and protein
1648 sequence. Coordinates given in the notes column specify the overlap of any remarkable features
1649 with domains C-lectin, FBD or IgSF.

1650

1651 Supplementary Table 15. Reference sequences of Fibrinogen-related protein (FREP) and C-type
1652 lectin-related protein (CREP) included in the phylogenetic analysis for both these protein families.
1653 Table details the sequence ID, source (reference), identifier provided within the family, length in
1654 amino acids and sequence. Sequences from Lu et al. (2020) marked with an X at the end of their
1655 Sequence ID were used with the modifications detailed in Lu et al. (2020), rather than the original
1656 genome (.xlsx file).

1657

1658 Supplementary Table 16. Total number of genes included on each of the Fibrinogen-related protein
1659 (FREP) and C-type lectin-related protein (CREP) Phylogenetically Hierarchical Orthogroups
1660 (HOGs) across the analyzed *Biomphalaria* species genomes and outgroups (*B. glabrata* strain
1661 BB02 (GCF_000457365.2 (5)), *B. glabrata* strain iBS90 (GCA_025434165.1 (6)), *B. glabrata*
1662 strain iM (GCA_025434175.1 (6)), and the species *B. pfeifferi* (GCA_030265305.1 (8)), *B.*
1663 *straminea* (GCA_021533235.1 (7)), *Bulinus truncatus* (GCA_021962125.1 (125)) and *Elysia*
1664 *marginata* (GCA_019649035.1 (126))). For *B. sudanica*, we specify both the raw counts (i.e.
1665 ‘Raw’, all the members of the HOG) and how many of these were selected for the phylogenetic
1666 analysis (i.e. ‘Selected’) (.xlsx file).

1667

1668 Supplementary Table 17. List of 1047 genes that were present in the most diverse regions of the
1669 *Biomphalaria sudanica* genome, providing information on genomic position, BLASTp hits,
1670 InterPro v5.56-89.0 (51) and Pfam v35 (112) domains and inferred protein function categorized
1671 into immune groups (as detailed in Methods) (.xlsx file).

1672

1673 Supplementary Table 18. Results of *Biomphalaria sudanica* gene family expansion and
1674 contraction analysis performed using topGO v2.48.0 (72) in R v4.3.1 (137) with the Weight01
1675 algorithm, node size=10 and statistical test of Fisher (significant p-value ≤ 0.05). The test was
1676 performed by comparing the *B. sudanica* whole genome gene family composition to three separate
1677 lists (column=List) of the highly diverse genes identified; 1) ‘All’ 1047 highly diverse genes
1678 identified in the *B. sudanica* genome (as listed Supplementary Table 17), 2) ‘Group-2’ 245 genes
1679 classified as immune suspected genes from the highly diverse genes (group 2, Supplementary
1680 Table 17), and 3) ‘Group-2-3-4’ 247 genes shortlisted as the protein containing a transmembrane
1681 domain (TMD) and categorized in Immune Group 2, 3 (potential role in innate immunity) or 4 (an
1682 unknown function but contained TMD(s)) (Supplementary Table 17). Each group contains three
1683 tabs, showing gene family expansions and contractions in each set of genes in respects to their
1684 Molecular Function (MF), Cellular Component (CC) and Biological Process (BP) (.xlsx file).

1685

1686 Supplementary Table 19. Pfam v35 (112) matches for function domains the 1047 most diverse
1687 genes identified in the *Biomphalaria sudanica* genome (.xlsx file).

1688

1689 Supplementary Table 20. UniProt (uniport.org/ database as of May 25th 2022) matches for the
1690 1047 most diverse genes identified in the *Biomphalaria sudanica* genome (.xlsx file).

1691

1692 Supplementary Table 21. BLASTp matches for the 1047 most diverse genes identified in the
1693 *Biomphalaria sudanica* genome (.xlsx file).

1694

1695 **Supplementary Files**

1696 Supplementary File 1. Complete functional annotation of *Biomphalaria sudanica*, including
1697 23,598 genes assigned to have an open reading frame (ORF) and therefore predicted protein by
1698 TransDecoder.Predict v5.5.0 using the Trinity platform (50), the 2,248 raw RNA genes without
1699 predicted proteins predicted by StringTie2 v2.2.1 (49), 919 tRNA's predicted in tRNAscan-SE 2.0
1700 v2.0.9 (122), and 107 rRNA's predicted in Barrnap v0.9 (github.com/tseemann/barrnap) (.gff file).

1701

1702 Supplementary File 2. Summary table of the 23,598 genes assigned to have an open reading frame
1703 (ORF) giving ORF type, length (amino acid) and prediction score from TransDecoder.Predict
1704 v5.5.0 using the Trinity platform (50) (.tab file).

1705

1706 Supplementary File 3. Summary InterProScan v5.56-89.0 (51) results for the predicted proteins. It
1707 includes the specific analysis, the signature accession and description identified, the start and stop
1708 location of the signature, its score, and associated InterPro Accession information (.xlsx file).

1709

1710 Supplementary File 4. Location of transmembrane domains identified on 4,922 genes (8,728
1711 isoforms) determined through DeepTMHMM v1.0.13 (129) (.gff3 file).

1712

1713 Supplementary File 5. File providing the REVIGO (73) TreeMap summaries of the GO terms
1714 enriched among the protein families, showing significant expansions and contractions for each
1715 node (excluding outgroups) of the species tree generated in Orthofinder using the Species Tree of
1716 All Genes (STAG) algorithm (67,74) (.pdf file).

1717

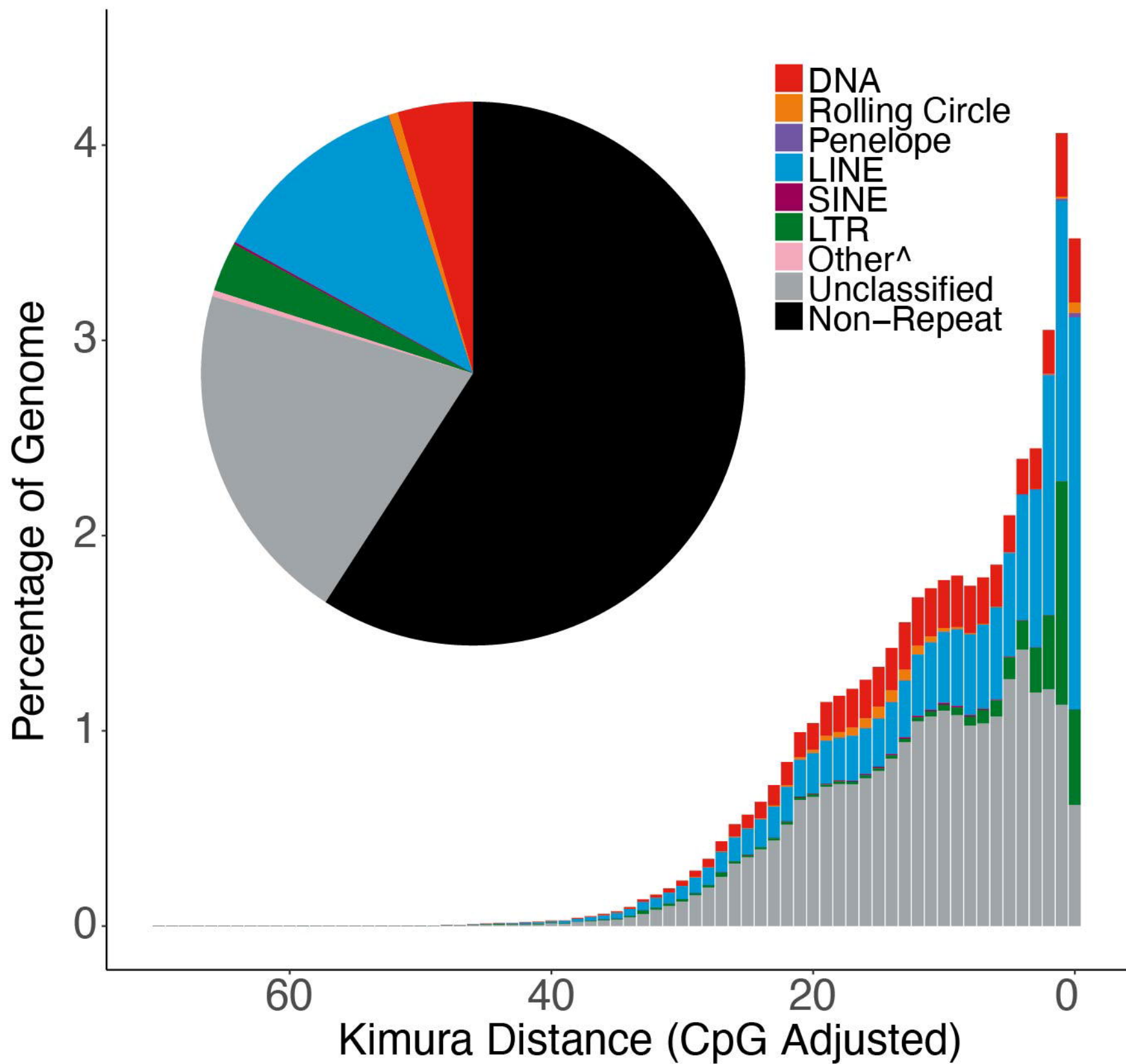
1718 Supplementary File 6. Description of RNA extraction methods, quantification, and pooling for
1719 sequencing (.docx file).

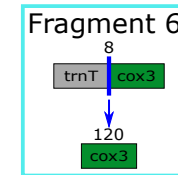
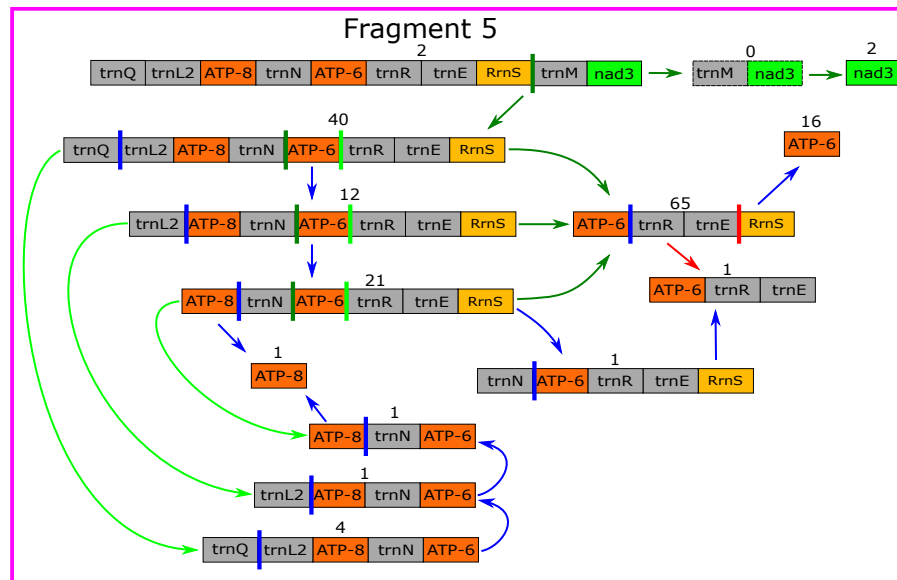
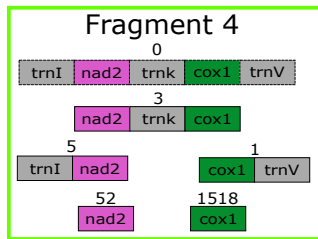
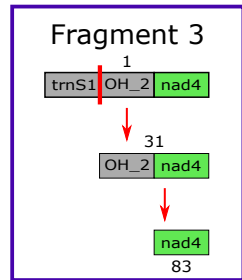
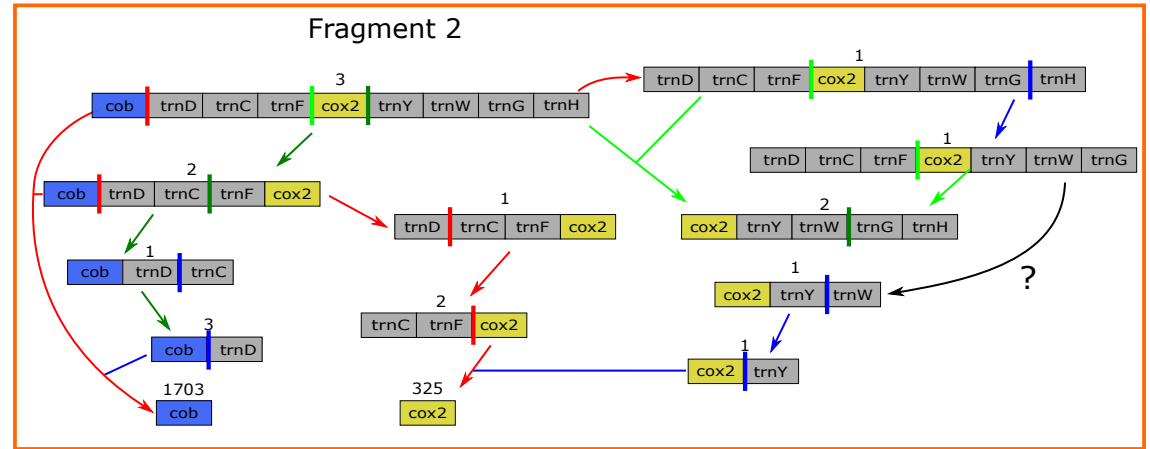
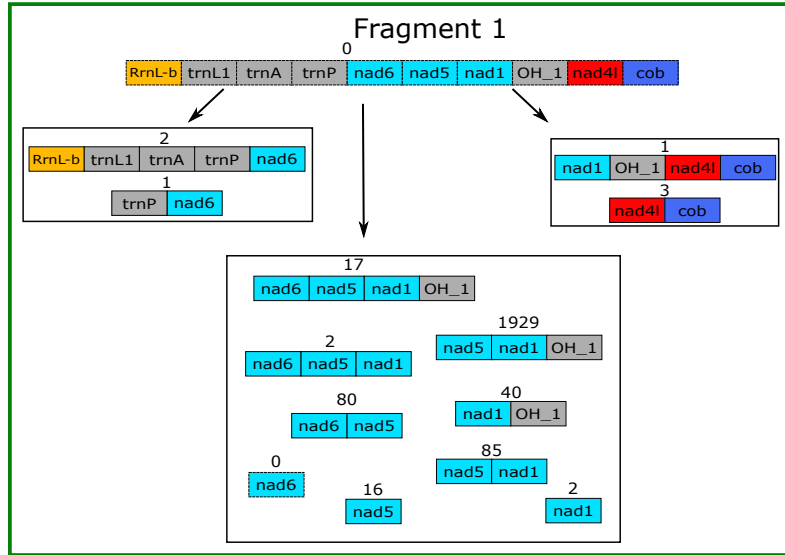
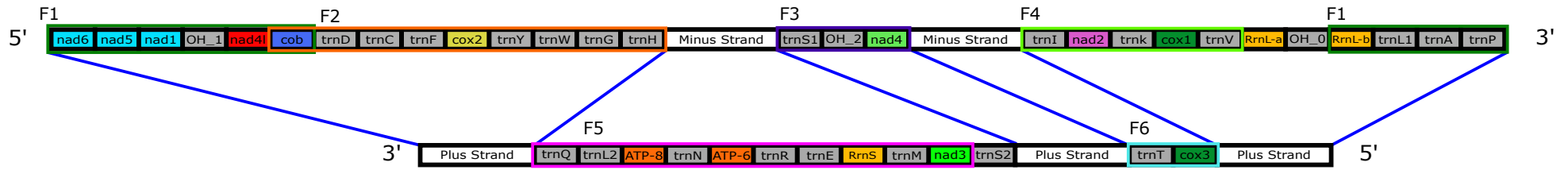
1720

1721 Supplementary File 7. Amino acid sequences of the 1047 shortlisted most diverse genes in the
1722 *Biomphalaria sudanica* genome, used for searches for novel pathogen recognition receptors (.fas
1723 file).

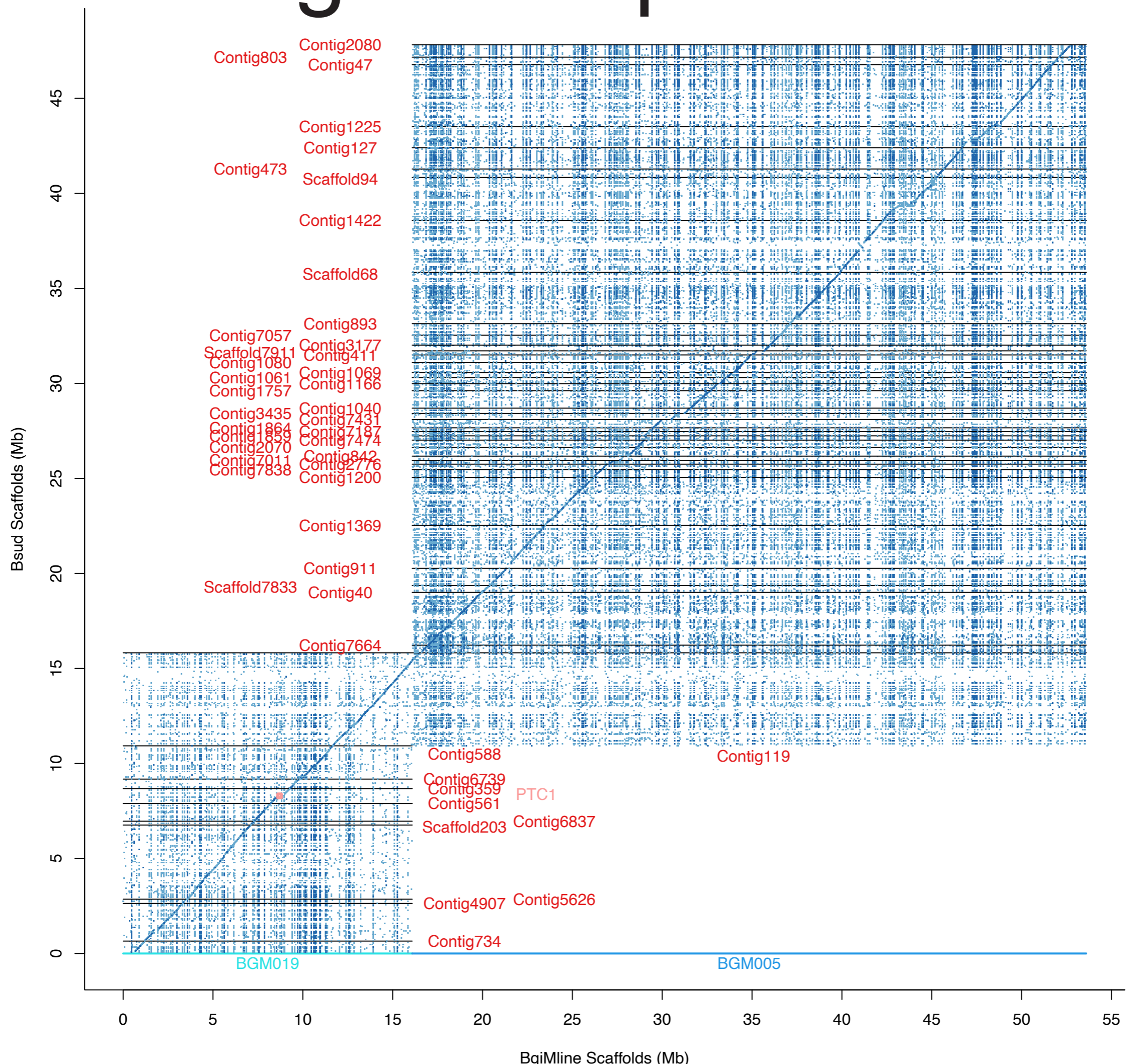
1724

1725

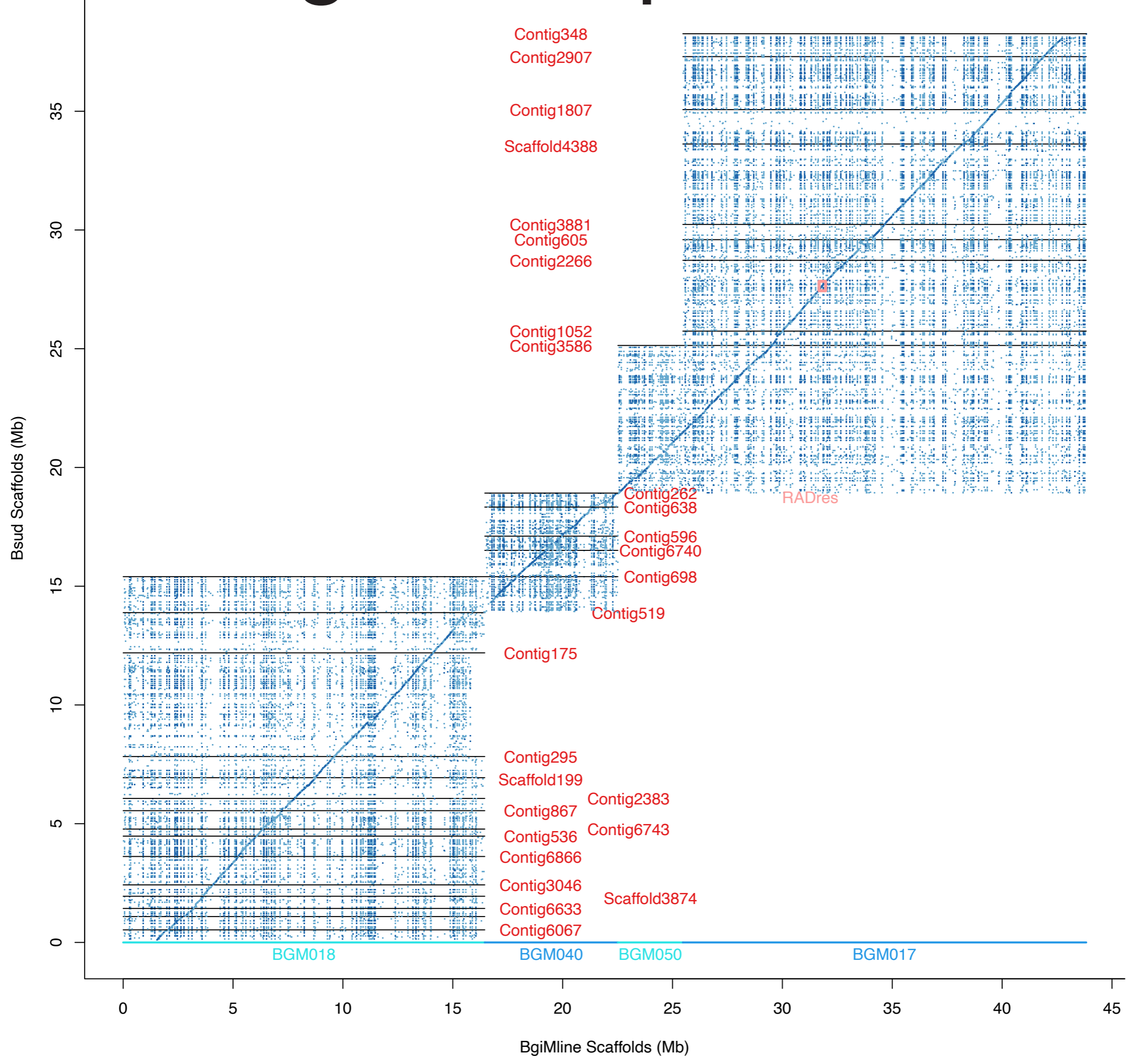




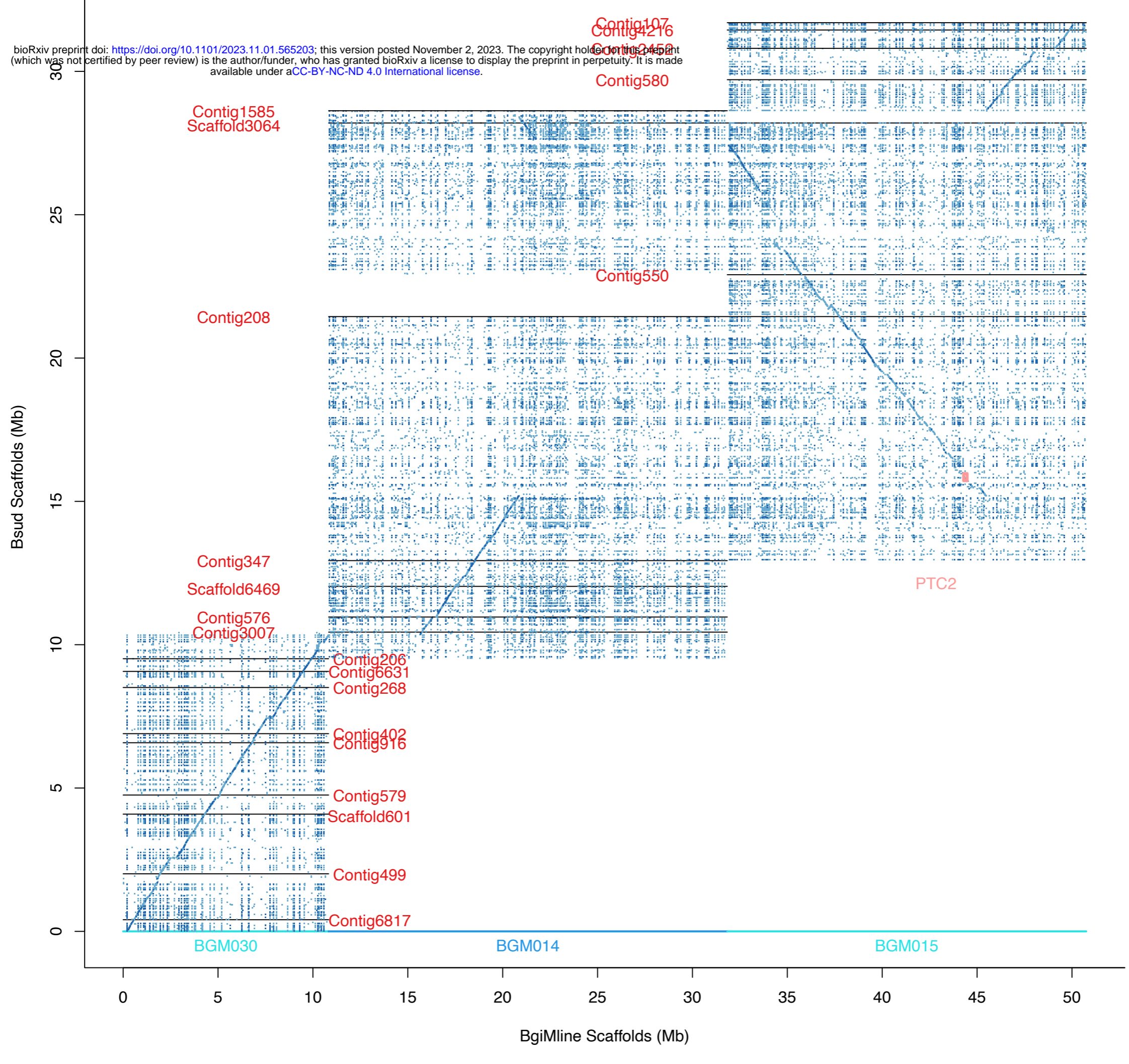
A. Linkage Group 6



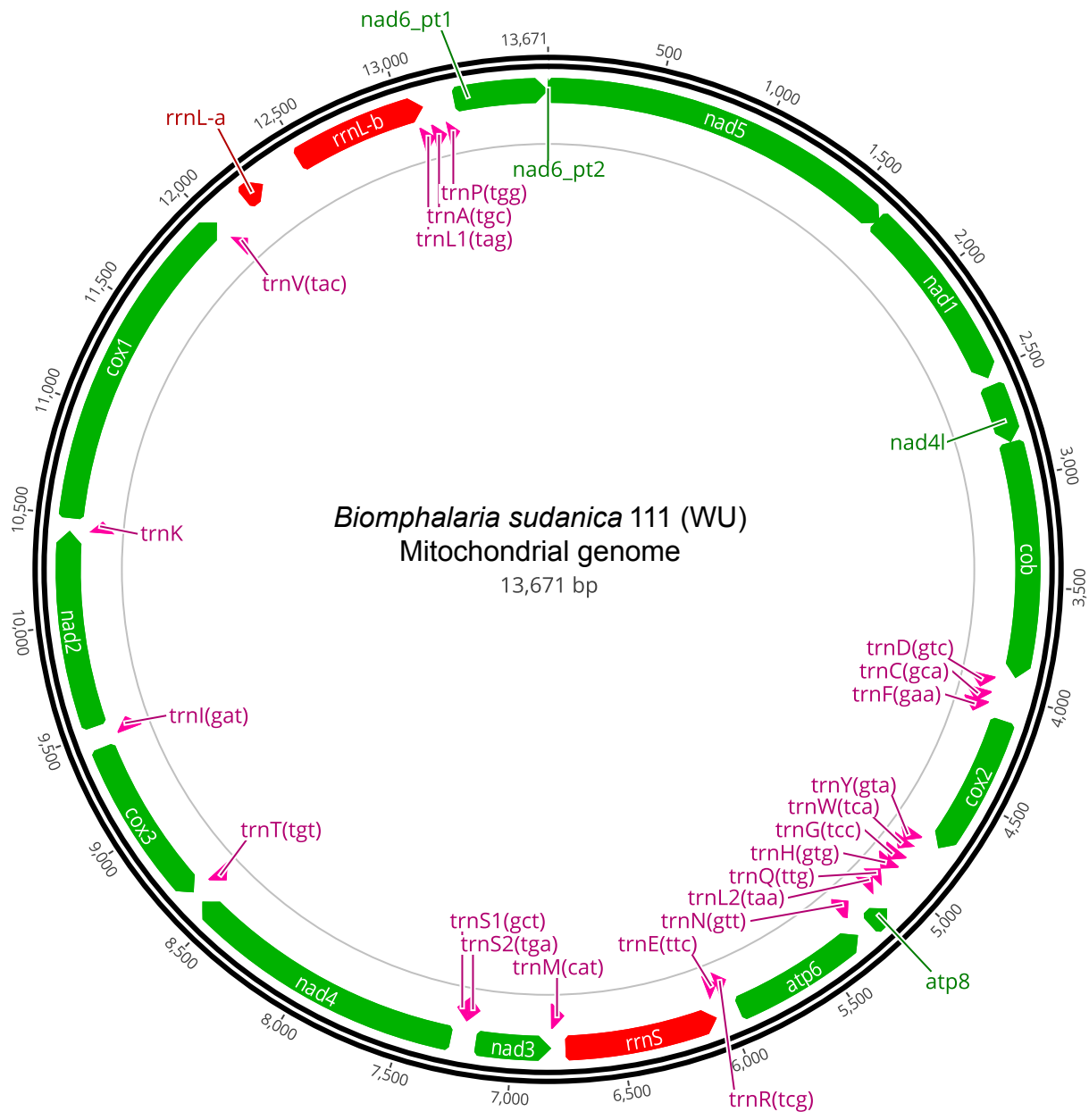
B. Linkage Group 10

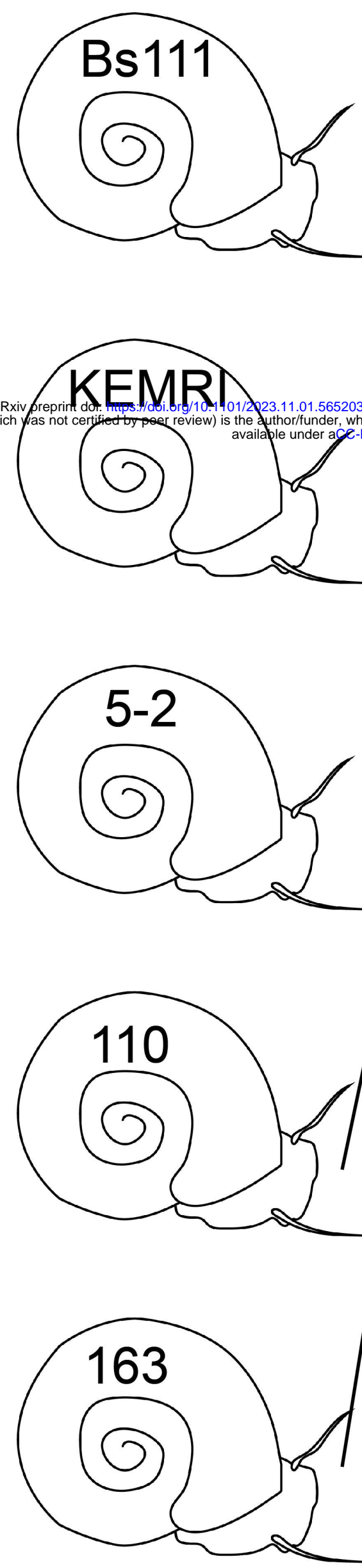


C. Linkage Group 16



bioRxiv preprint doi: <https://doi.org/10.1101/2023.11.01.565203>; this version posted November 2, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

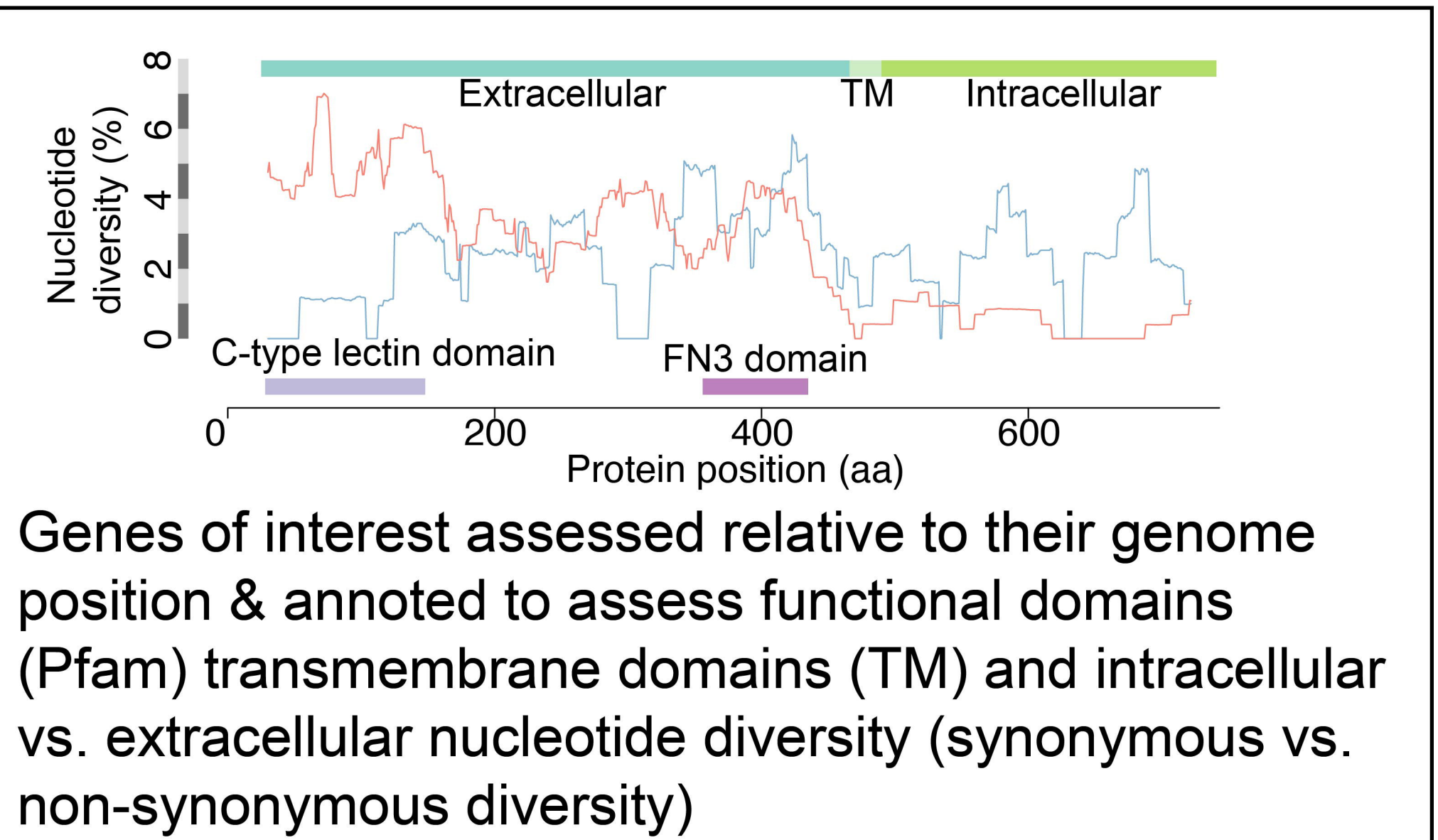
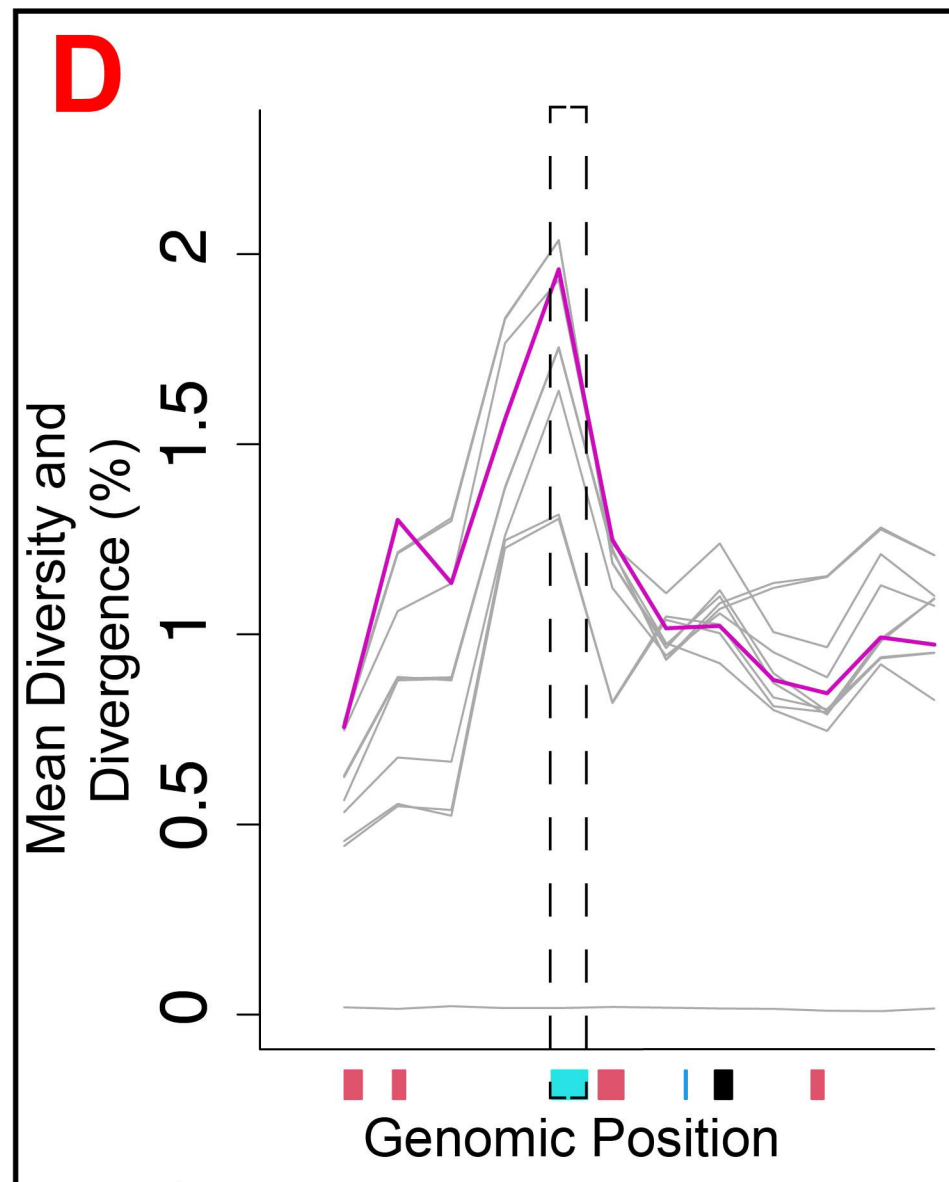
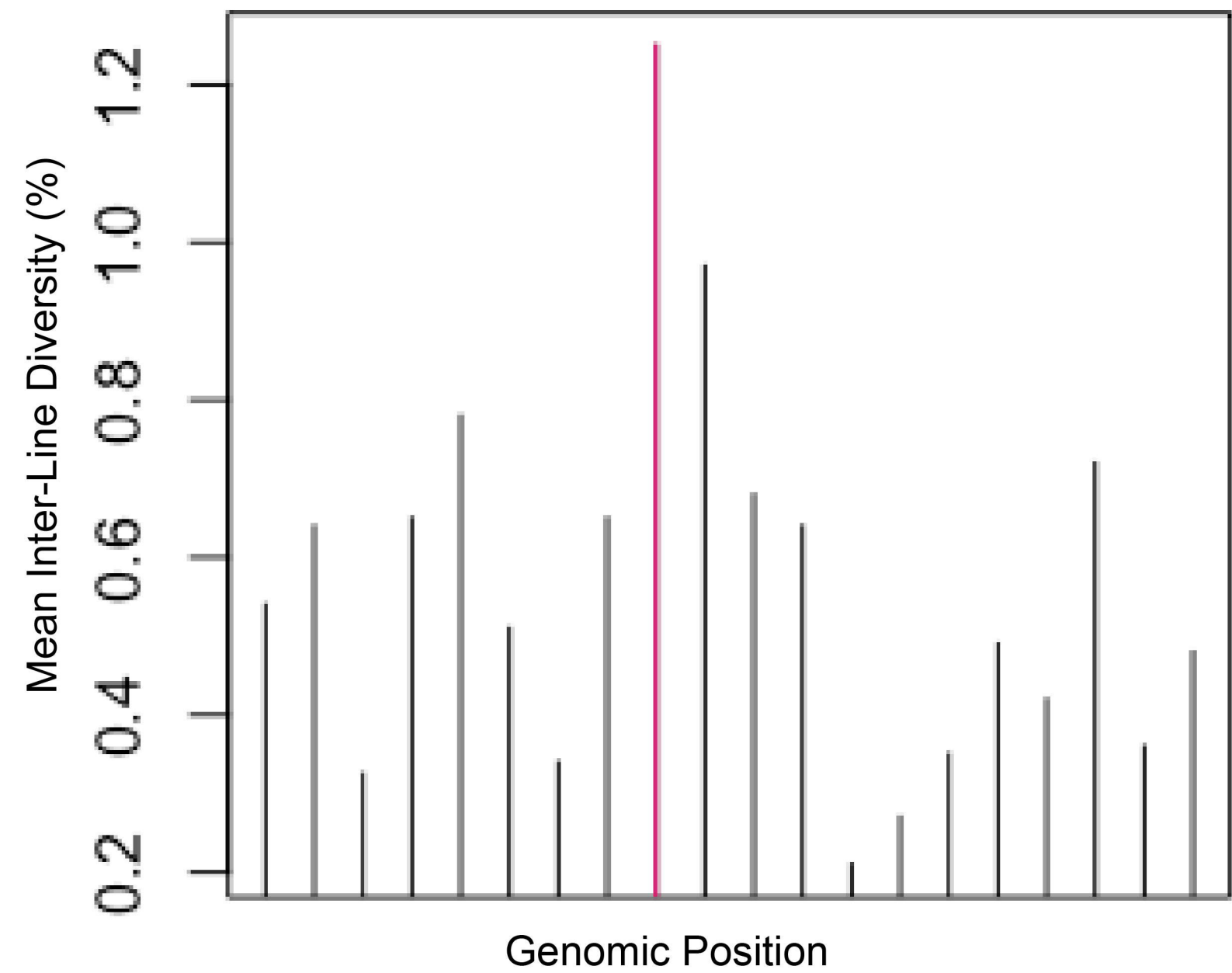




Genome sequencing (PacBio) and assembly for Bs111 reference genome. **A**

Genome sequencing (Illumina) of 4 other genetic lines and alignment to Bs111 reference genome. **B**

C Nucleotide hyper-diversity identified in annotated genes and in 10kb, 30kb or 100kb sliding windows.



E Discovery of candidate genes for pathogen recognition receptors (PRR) and other immune related genes that recognize pathogen-associated molecular patterns (PAMPs) and trigger innate immune response.

