

1 **Genomic adaptations to chemosymbiosis in the deep-sea seep-dwelling**
2 **tubeworm *Lamellibrachia luymesi* (Siboglinidae, Annelida)**

3 Yuanning Li^{1,2*}, Michael G. Tassia¹, Damien S. Waits¹, Viktoria E. Bogantes¹, Kyle T.
4 David¹, Kenneth M. Halanych^{1*}

5 ¹ Department of Biological Sciences & Molette Biology Laboratory for Environmental
6 and Climate Change Studies, Auburn University, Auburn, AL, 36849. USA

7 ² Department of Ecology and Evolutionary Biology, Yale University, 165 Prospect St,
8 New Haven, CT 06511. USA

9 * Corresponding author: yuanning.li@yale.edu; ken@auburn.edu

10

11 Classification: Biological Sciences: Evolution

12

13

14

15

16

17

18

19

20

21 **Abstract**

22 Genetic mechanisms allowing organisms to maintain host-symbiont associations
23 at the molecular level are still mostly unknown. In the case of bacterial-animal
24 associations, most genetic studies have focused on adaptations and mechanisms of the
25 bacterial partner. The gutless tubeworms (Siboglinidae, Annelida) are obligate hosts of
26 chemoautotrophic endosymbionts (except for *Osedax* which houses heterotrophic
27 *Oceanospirillales*). Whereas several siboglinid endosymbiont genomes have been
28 characterized, genomes of hosts remain unexplored. Here, we present and characterize
29 the genome of the cold-seep dwelling tubeworm *Lamellibrachia luymesi*, one of the
30 longest-lived invertebrates. With a haploid genome size of ~688 Mb and overall
31 completeness of ~95%, we discovered that *L. luymesi* lacks many genes essential in
32 amino acid biosynthesis obligating them to products provided by the symbionts. In
33 comparison, the host carries hydrogen sulfide to thiotrophic endosymbionts using
34 hemoglobin. Interestingly, we found a large expansion of hemoglobin B1 genes many of
35 which possess a free cysteine residue which is hypothesized to function in sulfide-
36 binding. Moreover, sulfide-binding mediated by zinc ions is not conserved across
37 tubeworms, suggesting the hemoglobin structure and the sulfide-binding mechanism is
38 potentially more complex than previously thought. Our comparative analyses also
39 suggest the Toll-like receptor pathway may be essential to host immunity and
40 tolerance/sensitivity to symbionts and pathogens. Last, we identified several genes
41 known to play an important role in longevity. These results help elucidate previously
42 unknown links and potential genetic mechanisms related to the evolution of holobionts,
43 adaptations to reducing environments, and likely extend to other chemosynthetic
44 symbiosis.

45 Keywords: chemosynthetic symbiosis, cold seep, comparative genomics, nutrition
46 mode, hemoglobins, Toll-like receptor, aging

47

48 **Significance**

49 Symbioses between bacteria and animals are ubiquitous and ecosystems (e.g.,
50 seeps, hydrothermal vents, and organic falls) driven by chemoautotrophy have received
51 considerable attention because of the non-photosynthetic energy source. However,
52 genomic machinery that led to evolutionary success of these chemosynthetic
53 environments is poorly understood, especially for hosts. By characterizing the genome
54 of the seep-dwelling tubeworm *Lamellibrachia luymesi*, we provide genetic evidence of
55 how animals adapted to extreme environments and maintain chemosynthetic symbiosis.
56 Host genome adaptations include loss of biosynthesis pathways, expansion of
57 hemoglobin gene families, innate immunity mechanisms related to host-symbiont
58 recognition, and genes related to longevity. Our findings can be extended to other taxa
59 and shed light on the mechanisms that establish and promote symbiosis, especially in
60 chemosynthetic systems.

61

62

63

64

65

66

67

68

69

70

71

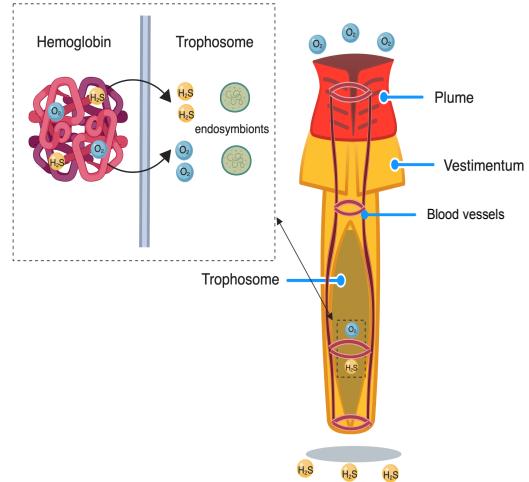
72 **Introduction**

73

A



B



74

75 **Fig. 1.** *Lamellibrachia luymesi*. (A). Seep habitat in the Gulf of Mexico. (B). Diagram of
76 adult *L. luymesi* worm to model O₂ and H₂S transport to symbionts in trophosome by
77 hemoglobin molecules. The hemoglobin model was created with the help of Biorender
78 (<https://biorender.com/>).

79 Recent advances in understanding the dominance of microbes on the planet has
80 placed new emphasis on elucidating mechanisms that promote microbe-animal
81 symbioses. Although considerable work has been undertaken on adaptations of
82 microbial genomes to facilitate animal symbiosis (such as corals, termites, humans),
83 examples of how animal host genomes have adapted to symbioses are limited (1).
84 Vestimentiferan tubeworms inhabit some of Earth's most extreme environments, such
85 as deep-sea hydrothermal vents and cold seeps, and are obligate dependents on
86 symbiosis for survival. These animals lack a digestive tract and rely on sulfide-oxidizing
87 bacteria symbionts for nutrition and growth. At some seeps, tubeworms, such as
88 *Lamellibrachia luymesi* in the Gulf of Mexico, are so abundant that they transform the
89 habitat (Fig. 1A) and thus facilitate biodiversity promoting adaptive radiations and
90 evolutionary novelties (2). Given the obligate nature of the symbiosis between
91 tubeworms and their gammaproteobacterial chemoautotrophic endosymbiont, one may
92 reasonably expect adaptations in several cellular mechanisms and pathways (e.g.

93 nutrition, gas exchange, self-defense/self-recognition) to promote efficacy in the
94 symbiotic relationship.

95 Siboglinid hosts acquire their symbionts from the surrounding environment and
96 store them in a specialized tissue called the trophosome (3). The chemosynthetic
97 symbionts are known to use a variety of molecules (e.g. H₂S, O₂, H₂) for final electron
98 receptors facilitating a variety of fixation pathways (4). Primarily, vestimentiferan
99 symbionts use both reverse TCA cycle (rTCA) and the Calvin cycle for carbon fixation
100 providing a nutrient source for the host (4, 5). To date, metabolic studies have primarily
101 focused on mechanisms and pathways found in symbionts and studies from the host's
102 perspective are limited.

103 Another key adaptation contributing to the ability of tubeworms to thrive in
104 chemosynthetic habitats involves hemoglobins (Hbs) that bind oxygen and sulfide
105 simultaneously and reversibly at two different sites (6) (Fig. 1B). To avoid the toxicity of
106 sulfide, siboglinids possess three different extracellular hemoglobins (Hbs): two
107 dissolved in the vascular blood, V1 and V2, and one in the coelomic fluid, C1 (7, 8).
108 Siboglinid Hbs consist of four heme-containing chains (A1, A2, B1, B2). Sulfur-binding
109 capabilities are hypothesized to be dependent on free cysteine residues at key positions
110 in Hbs, especially in the A2 and B2 chains (6). V1 Hb can form persulfide groups on its
111 four linker chains (L1-L4), a mechanism that can account for the higher sulfide-binding
112 potential of this Hb (6). However, a study suggested sulfide-binding affinity was
113 mediated by the zinc moieties bound to amino acid residues at the interface between
114 pairs of A2 chains in *Riftia* (9). Thus, it is still not clear which mechanism is primarily
115 responsible for sulfide-binding in siboglinids.

116 In contrast to rapidly growing vent-dwelling vestimentiferans (10), seep-dwelling
117 vestimentiferans have much slower growth rates, and are among the most long-lived
118 non-colonial marine invertebrates (up to 250 years) (11). Immunity has important
119 implications in aging (12), and is also a critical evolutionary driver of maintaining
120 symbiosis (13). However, little is known about genetic mechanisms relating immunity
121 and symbiosis. Because tubeworm endosymbionts are housed internally and their
122 establishment process resembles infection (3), tubeworm symbiosis provides a unique

123 opportunity to examine evolution of immunity functions associated with host-symbiont
124 relationships. However, Information on extremophile immunity and/or immune tolerance
125 is lacking.

126 Using comparative genomics, transcriptomic and proteomic analyses on the
127 tubeworm *Lamellibrachia luymesii*, we provide evidence for genetic pathways and novel
128 candidate genes which may underlie the extraordinary characteristics of tubeworm
129 symbioses. In particular, we focus on nutrition mode, hemoglobin evolution, immunity
130 function, and longevity.

131 **Results and Discussion**

132 **Genome features**

133 Using Illumina paired-end, mate-pair and 10X genomic sequencing (Table S1),
134 we the assembled genome of a single *Lamellibrachia luymesii* individual. The haploid
135 genome assembly size is ~688 Mb (Fig. S1) with ~500X coverage and N50 values of
136 373 Kb (scaffolds) and 24 Kb (contigs). Although N50 lengths and assembly quality of *L.*
137 *luymesii* are comparable to those of other annelids (e.g. *Capitella teleta*, *Helobdella*
138 *robusta*) (Tables S2, S3), the overall genome completeness measured by BUSCO (~
139 95%) is one of the highest among lophotrochozoans (Table S2). With the support of
140 RNA-seq data from three different tissues (Table S1), we estimated *L. luymesii* genome
141 contains 38,998 gene models. The genome also exhibits heterozygosity (0.6%) and
142 repetitive content (36.92%) similar to other lophotrochozoans (Fig. S2, Table S4). We
143 found 94 orthology groups (OGs) appear to have undergone a genomic expansion
144 compared to other lophotrochozoan genomes (Tables S5).

145 **Nutritional adaptations**

A

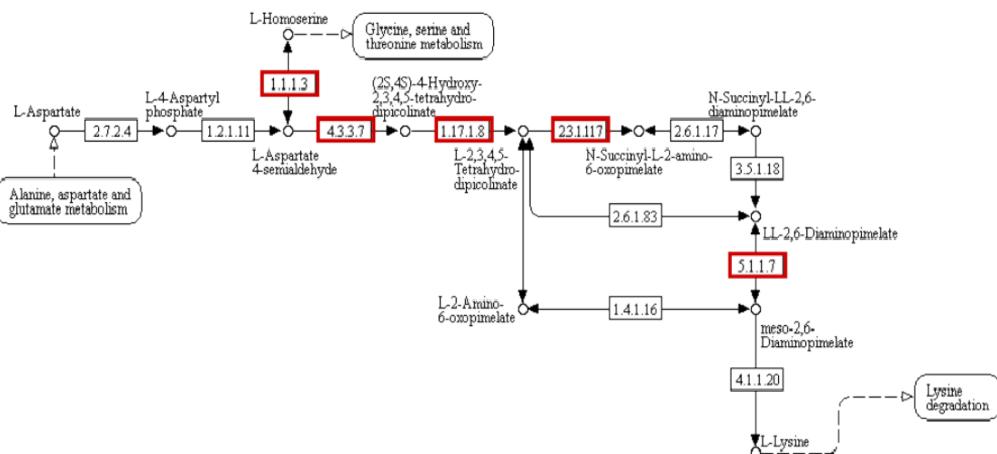
Lamellibrachia host	argA	argB	argC	argD	argE	argG	argH	argJ	ilvB	ilvG	ilvH	ilvA	ilvC	ilvD	leuA	leuB	leuCD	GlyA	Gly
Capitella teleta																			
Lamellibrachia symbiont																			
	Arg																		

Lamellibrachia host	aroA	aroB	aroC	aroE	aroQ	trpAB	trpC	trpD	trpE	trpF	metA	metE	metH	metK	thrB	thrC	dapA	dapB	dapD	dapF	lysA
Capitella teleta																					
Lamellibrachia symbiont																					
	Phe					Trp					Met				Thr						

Lamellibrachia host	hisA	hisB	hisD	hisG	hisH	tyrA2	proA	proB	proC	cysE	cysK	cysM	gltB	gltD	glnA	serA	serB	serC1	Dat	aspC	asnB
Capitella teleta																					
Lamellibrachia symbiont																					
	His					Tyr	Pro			Cys			Glu		Gln	Ser	Ser		Ala	Asp	Asn

■ Present □ Absent ■ Gene only absent in *Lamellibrachia* host but present in both *Capitella* and *Lamellibrachia* symbiont

B



146

147 **Fig. 2.** *Lamellibrachia luymesii* lacks amino acid biosynthesis genes. (A) Presence
148 (green) or absence (white boxes) of key genes associated with amino acid biosynthesis
149 in the genomes of *Capitella teleta*, *L. luymesii* and *L. luymesii* symbionts. Blue boxes
150 represents genes present in *C. teleta* and *L. luymesii* gammaproteobacterial symbionts
151 but absent in *L. luymesii*. (B) Example of Lysine biosynthesis pathway. Red boxes
152 indicate genes missing in *L. luymesii*. Figure was created with the help of KEGG
153 webserver.

154 Only 57 genes associated with amino acid biosynthesis were found in the *L.*
155 *luymesii* genome, of which eight were also identified in the proteomic analysis. In
156 contrast, the *Capitella teleta* (Capitellidae, Annelida) genome contains 90 such genes

157 (Fig. 2A; Supplementary Dataset 1), despite being a less complete and more
158 fragmented genome (Table S2). These gene were not clustered together in the
159 genomes suggesting that they were probably not missed due to random chance given
160 the completeness of sequencing. Interestingly, the *L. luymesii* symbiont genome
161 contains 110 genes, an essentially complete set for biosynthesis of all 20 proteinogenic
162 amino acids and of 11 vitamins/cofactors. Genes found in *C. telata*'s genome but
163 lacking in *L. luymesii* are involved in biosynthesis of 13 amino acids (e.g., five key
164 enzymes in the Lysine biosynthesis pathway Fig. 2B). As amino acids are essential for
165 protein biosynthesis in the host, the lack of many important amino acid synthesis-related
166 genes indicate that the host depends on symbionts for amino acids and cofactors.
167 Moreover, we found a large gene expansion of nutrient uptake ABC transport protein-
168 coding genes in *L. luymesii* compared with other lophotrochozoans (Table S5). These
169 findings are consistent with previous biochemical analyses which suggest that *Riftia* is
170 also dependent on its bacterial symbiont for the biosynthesis of polyamines that are
171 important for host metabolism and physiology (14).

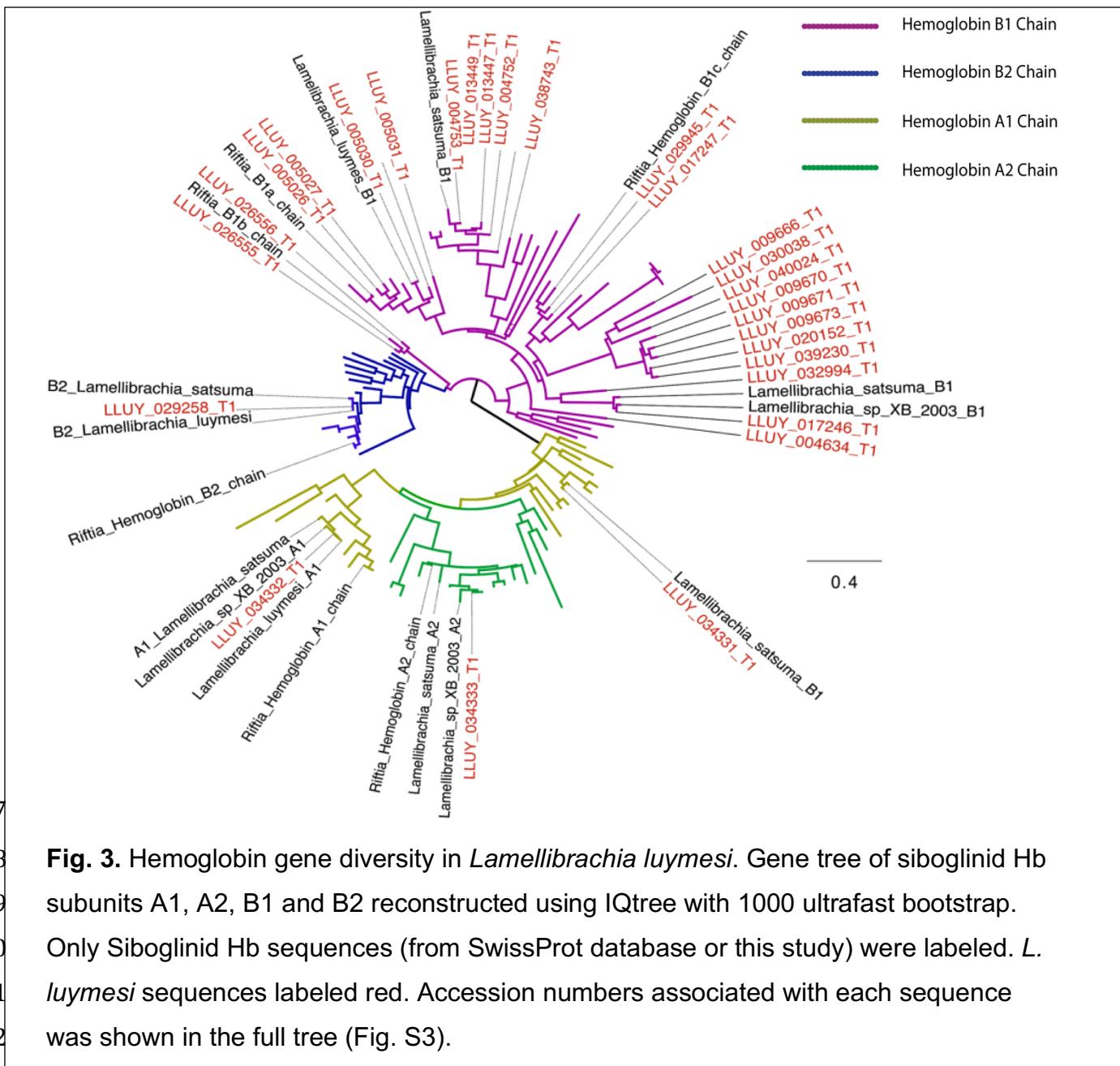
172 Obligate bacterial symbionts often lack genes that are commonly found in other
173 free-living bacteria, while retaining only those genes with functions essential to host
174 needs (e.g. in sponges, (15); in termites, (16)). However, there are known cases of loss
175 in essential gene functions in multicellular eukaryotes, but this phenomenon appears to
176 be more frequent in bacterial symbionts (1). Interestingly, thiotrophic symbionts of the
177 vesicomyid clam *Calyptogena magnifica* (17) and vent mussel *Bathymodiolus azoricus*
178 (18) have been suggested provide their host with products from amino acid
179 biosynthesis. Moreover, a recent study has suggested that the flatworm *Paracatenula*
180 itself does not store primary energy in host cells; rather, this function is performed by its
181 chemosynthetic symbionts (19). Although the tubeworms and bivalves under
182 examination in the aforementioned studies live in chemosynthetic environments, the
183 different hosts and bacteria represent disparate genomic backgrounds suggesting that
184 modification and loss of the amino acid biosynthesis pathways may be a convergent
185 adaptation in a variety of chemosynthetic symbioses between bacteria and animals.

186 In addition to the immediate release of fixed carbon and provision of amino acids
187 by symbionts, we have found proteomic evidence of a second possible nutritional mode
188 whereby the host directly digests symbionts, as shown by the detection of abundant
189 host-derived digestive enzymes in trophosome tissue (Table S6). Previous observations
190 indicated that symbionts could be digested by *Riftia* (20) but, direct evidence and
191 mechanisms related to symbiont digestion lacked characterization. We identified 15
192 host proteins related to lysosomal proteases that were both highly expressed and
193 detected as proteins in the trophosome tissue of host genome, such as Saposin and
194 multiple copies of Cathepsin (Table S6). Lysosomes, which contain an array of digestive
195 enzymes, are also thought to play an essential role in symbiont digestion with the
196 chemosynthetic mussel *Bathymodilus azoricus* (18). We additionally identified 19 major
197 proteasome components as proteins in the trophosome tissue, indicating a potential role
198 in protein degradation of symbiont digestion (Table S6). Host lysosomal proteases and
199 proteasome components likely facilitate degradation of symbionts and may play a role
200 maintaining appropriate population levels of symbionts within trophosome.

201 We also characterized ~ 200 bacterial proteins present in the same trophosome
202 tissue to further understand host-symbiont interactions. Key enzymatic genes,
203 RubisCO, and ATP citrate lyase (ACL) type II associated with carbon fixation cycles,
204 were identified in proteomic analysis from *L. luymesii* (Table S7). Our results corroborate
205 that both rTCA and Calvin cycle, pathways for carbon fixation might be common in all
206 vestimentiferan endosymbionts (4). Several key components related to sulfide and
207 nitrogen metabolic pathways were identified consistent with previous analyses (4, 5).

208
209
210
211
212
213
214

215

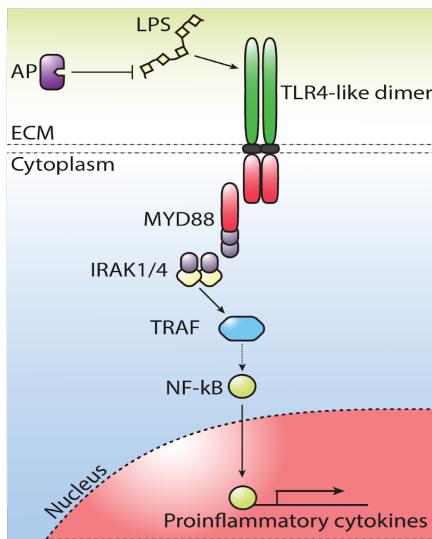
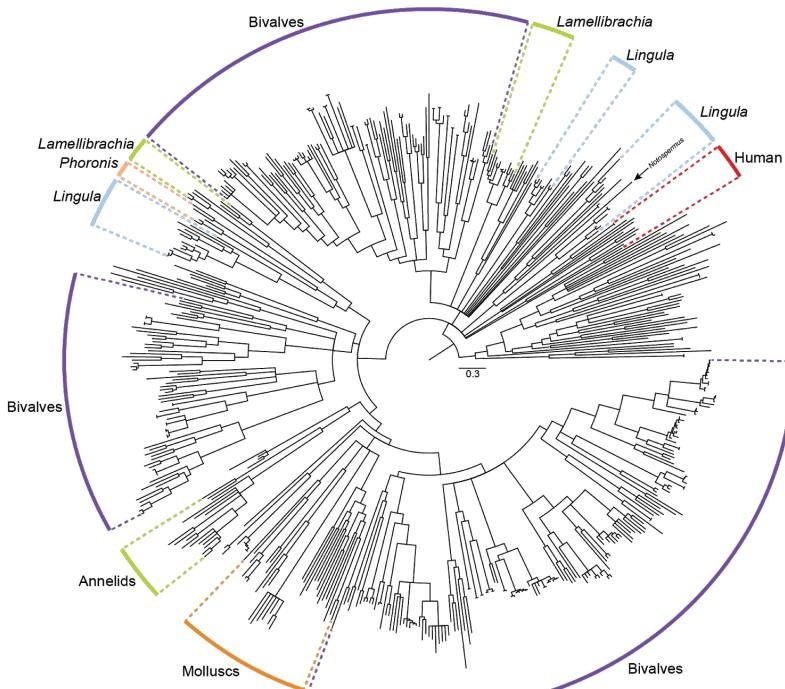
216 **Hemoglobin evolution**

223 Mechanisms of Hb sulfide-binding affinity in tubeworm siboglinids are still not
224 clear after 20 years of study. We collected all available Hb sequences from siboglinids
225 and their close relatives and processed them through a phylogenetic framework (Figs.
226 3, S3). Importantly, we are able to identify most Hbs and linkers from transcriptomic
227 and proteomic results (Table S7). Consistent with (6, 8, 9) a single copy of A2 and B2

228 Hb was identified in all siboglinids which possesses a conserved-free cysteine (i.e.,
229 cysteine residues not involved in disulfide bridges) at position 77 and 67, respectively.
230 With exception of A2 and B2 Hbs in the earthworm *Lumbricus terrestris*, homologous
231 cysteine residues were identified in 3 annelids (*Cirratulus spectabilis*, *Sabella pacifica*,
232 and *Sternapsis* sp.) from sulfide-free environments and *Arenicola marina* living in
233 sulfide-rich environments (Fig. S4). These results support the hypothesis that free
234 cysteine residues in A2 and B2 Hbs were present in all annelids and potentially involved
235 in H₂S detoxification process (21).

236 Surprisingly, we found a significant expansion of B1 Hbs, 25 copies, in *L. luymesii*
237 whereas most siboglinids and their close relatives only possess one copy indicated by
238 previous studies (Fig. 3B), except for *Riftia pachyptila* where three B1 Hbs were
239 identified (21). Noticeably, we found that 8 copies of *L. luymesii* B1 Hb sequences also
240 contains a free cysteine at position 77, the same position as free cysteine in A2 Hbs.
241 Five out the 8 copies were highly expressed in the trophosome, and one copy was
242 identified at the protein level (Table S8).

243 Instead of free cysteines mediating H₂S binding, another hypothesis suggested
244 that zinc moieties bound to amino acid residues at the interface between pairs of A2
245 chains influence H₂S binding (9). The Zn²⁺-binding site contained within A2 chain is
246 composed of three His residues (B12, B16, and G9) (9). However, none of these sites
247 are conserved across siboglinids, or even in vestimentiferans (Fig. S5) calling into
248 question the role of the zinc sulfide binding mechanism for H₂S transport.

A**B**

250

Fig. 4. Toll-Like Receptors (TLRs) in *Lamellibrachia luymesii*. (A) Putative TLR4-like pathway likely essential for immunity and response to symbionts and pathogens. AP: alkaline phosphatase; LPS: lipopolysaccharide. (B) Toll-Like Receptor gene tree from selected lophotrochozoan genomes and human reconstructed using IQtree with 1000 ultrafast bootstraps. All internal nodes possess $\geq 95\%$ bootstrap support.

256 Immune interactions between hosts and symbionts is a key evolutionary driver
257 that has potential implications in aging (12). The genetic machinery and functionality of
258 the immune system in chemosynthetic symbioses have not been extensively
259 characterized. Toll-Like Receptor (TLRs) provides a core cellular and molecular
260 interface between invading pathogens and recognition of host-microbial symbiosis (13)
261 (Fig. 4A). Consistent with previous analyses (Luo et al. 2018), we found that TLR gene
262 families experienced expansion within lophotrochozoan lineages (Fig. 4B; Table S9).
263 Within *L. luymesi*, 33 unique TLR proteins were identified compared to 5 in *Capitella*
264 *telata*, suggesting TLR genes have additional functions in tubeworms.

265 A substantial subset of TLR sequences recovered from *L. luymesi* best identify
266 as TLR4 by primary sequence identity and domain structures. In mammals, TLR4
267 recognizes and binds lipopolysaccharide (LPS; a major cell-membrane component of
268 Gram-negative bacteria which include tubeworm symbionts). LPS-bound TLR4 then
269 initiates a signal-transduction pathway that activates NF- κ B, a transcription factor that
270 promotes the expression of pro-inflammatory cytokines (22) (Fig. 4). *Lamellibrachia*
271 *luymesi* encodes seven TLR4-like proteins, which is in contrast to the one sequence
272 found in other annelid genomes suggesting potential for increased sensitivity to Gram-
273 negative bacteria in *L. luymesi*. Interestingly, we also found genomic expansions of
274 tumor necrosis factor receptors (TNFRs) and TNFR-associated factors (TRAFs) (Table
275 S5) which play vital roles in activation and the downstream responses of NF- κ B, further
276 supporting a specialized/expanded role for TLR4-like signaling. Whereas some other
277 components of the innate immunity (e.g. RIG-1-like receptor signaling pathway which
278 recognizes virus-derived nucleotide present in the cytoplasm) showed no indication of
279 gene expansion, the NLRP gene family (which plays a key role in an innate immunity
280 recognition of infectious pathogens and regulates inflammatory caspases) and Sushi
281 domain-containing genes (potential recognition and adhesion between hosts and
282 symbionts, (18) showed expansion relative to other lophotrochozoans. (Table S9).

283 The initial physical encounter between tubeworms and symbionts occurs in an
284 extracellular mucus secreted by pyriform glands of newly settled larvae (3). Within these
285 mucus matrices, symbionts can attach to the host using extracellular components

286 secreted from symbionts, such as LPS. The symbiont's colonization process induces
287 massive apoptosis of host skin tissue as symbionts travel from host epidermal cells into
288 trophosome (3). Recognition of lipopolysaccharide (LPS) by TLR4 can result in the
289 induction of signaling cascades that lead to activation of NF- κ B and the production of
290 proinflammatory cytokines (13). Although the mechanism by which host distinguishes
291 between symbionts and pathogens in most symbioses is still not clear, alkaline
292 phosphatase has been shown to be involved in the maintenance of homeostasis of
293 commensal bacteria in the squids, mouse, and zebrafish (23). The commensal
294 bacterially-derived LPS signaling via TLR4 yields an upregulation of intestinal alkaline
295 phosphatase and prevents inflammatory responses to resident microbiota. Importantly,
296 we also identified 8 copies of alkaline phosphatase, whereas only one copy was found
297 in each of the *Capitella teleta* and *Helobdella robusta* genomes, further supporting a
298 potential mechanism of tolerating Gram-negative bacteria and facilitating symbiotic
299 colonization. Thus, although further analysis is warranted, a TLR4-like signaling
300 pathway may be central for host immunity and in distinguishing between symbionts and
301 pathogens (Fig. 4A).

302 **Aging**

303 Deep-living vestimentiferans are long lived, and in addition to innate immunity,
304 our analyses of gene family expansion highlighted families that may play a direct role in
305 aging. We found expansion of interleukin 6 receptors (IL6R) which are the key
306 component of the main signaling pathway implicated in aging (24). Superoxide
307 dismutases (SODs) have important function role in cells to protect against oxidative
308 damage induced by metabolism and are implicated in aging and redox balancing. We
309 found genomic expansions of CuZn-superoxide dismutase (SOD1) genes and Mn-
310 superoxide dismutase (SOD2) in *L. luymesii*'s genome compared to other
311 lophotrochozoans (Fig. S6). Most lophotrochozoan genomes contain one or two copies
312 of SOD1 and SOD2, but *L. luymesii* has 5 copies of each gene (Fig. S6). Three of 5
313 SOD2 genes were recovered in transcriptomic and proteomic data (Table S6). Previous
314 studies suggested that overexpression of SOD1 or SOD2 could significantly extend
315 lifespan in mammals, fruit flies and *C. elegans* (25) and SOD gene product may help

316 symbionts overcome host cellular immune responses (26). Consistent with previous
317 studies, we also be able to identify symbionts' SOD gene as proteins in proteomic
318 analysis. Thus, SODs from both bacteria and tubeworms may play a central role for
319 overcoming oxidative damage and essential for extreme longevity for seep-living
320 vestimentiferans.

321 **Conclusion**

322 We characterize of the genome for the deep-sea seep-living tubeworm
323 *Lamellibrachia luymesi*. This report provides critical insight that hosts, like their bacterial
324 partners, may lose essential genomic components when their life-history strategy relies
325 on symbiotic interactions. Analyses show that *Lamellibrachia luymesi* has lost key
326 genes for amino acid biosynthesis making it necessarily dependent on endosymbionts.
327 Additionally, expansions have occurred in a number of gene families (e.g., TLRs, SODs,
328 Hemoglobins) that have been implicated in bacterial symbiosis. Evolutionarily,
329 increasing the number of paralogs provides opportunity for neofunctionalization or
330 subfunctionalization allowing more refined gene-gene interactions to promote symbiotic
331 efficacy. This balance of gene family expansion and gene loss may be a hallmark of
332 how genomic machinery adapts and develops interdependence across of variety of
333 bacterial-animal symbioses.

334 **Methods.**

335 **Organismal collection**

336 *Lamellibrachia luymesi* was collected from seeps in the Mississippi Canyon in the
337 Gulf of Mexico (N 28°11.58', W 89°47.94' 754m depth), using the *R/V Seward Johnson*
338 and *Johnson Sea Link* in October 2009. Samples were frozen at 80°C following
339 recovery.

340 **Genome sequencing and assembly**

341 Using vestimentum tissue of one individual, high molecular weight genomic DNA
342 was extracted using the DNeasy Blood & Tissue Kit (Qiagen). Four TruSeq paired-end
343 and two Nextra mate-pair genomic DNA libraries were generated and sequenced by

344 The Genomic Services Lab at the Hudson Alpha Institute for Biotechnology in
345 Huntsville, Alabama on an Illumina HiSeq platform (Table S1). Additionally, Hudson
346 Alpha constructed and sequenced a Chromium 10X sequencing library (10X genomics)
347 on an Illumina HiSeqX platform.

348 Our genome assembly workflow is shown in Fig. S7. Paired-end and 10X raw
349 reads were checked with FastQC v0.11.5 (27) and quality filtered (Q score >30) with
350 Trimmomatic v0.36 (28). Genome size, level of heterozygosity, and repeat content were
351 determined using kmer histograms generated from the paired-end libraries in Jellyfish
352 v2.2.3 (29) and GenomeScope (30) (Fig. S1). Mate-pair reads were trimmed and sorted
353 using NxTrim v0.3.1 (31), and only “mp” (true mate-pair reads) and “unknown” (mostly
354 large insert size reads) reads were used for downstream scaffolding analysis.

355 Given high heterozygosity in non-model species, all reads were assembled using
356 Platanus v1.2.4 (32) with a kmer size of 32. Scaffolding was conducted by mapping PE
357 and MP reads to Platanus contigs using SSPACE v3.0 (33). Gaps in scaffolds were
358 filled with GapCloser v1.12 (34) and redundant allele scaffolds were removed using
359 Redundans v0.13c (default settings; (35). Genome assembly quality was assessed with
360 QUAST v3.1 (36) and genome completeness with BUSCO v3 (37) using the
361 Metazoa_odb9 database (978 Busco genes). We also assemble the genome using 10X
362 data in Supernova 1.2.0 (Weisenfeld et al. 2017), but the genome quality and
363 completeness was inferior to the Platanus assembly (Fig. S7) and there for ignored.

364 **Transcriptome assembly and analysis**

365 Total RNA was extracted via Trizol (Thermo Fisher Scientific) from the plume,
366 vestimentum and trunk/trophosome tissue of the same *L. luymesi*. RNA-Seq was
367 carried out by Hudson Alpha on using Illumina HiSeq platform. After quality checking
368 and trimming of raw sequencing reads, transcripts were assembled in Trinity v2.4.0
369 (38). Transcript isoforms with high similarity ($\geq 95\%$) were removed with CD-HIT-EST
370 v4.7 (39). Transcripts were verified and abundance estimated by read mapping with
371 Bowtie v2.2.9 (40) and RSEM v1.2.26 (41).

372 **Genome annotation**

373 Gene models were constructed following the Funannotate pipeline 1.3.0
374 (<https://github.com/nextgenusfs/funannotate>; Fig. S8) using information from the
375 genome assembly, transcriptome assembly, and SwissProt/Uniprot. For genome data,
376 repetitive regions were identified using RepeatModeler v1.0.8 (43) and soft-masked
377 using RepeatMasker v4.0.6 (44). For each transposable element (TE) superfamily,
378 relative ages of different copies were estimated by calculating Kimura distances
379 assuming that most of the mutations are neutral using repeatLandscape.R
380 (https://github.com/dunnlab/genome_annotation/blob/master/repeatLandscape.R).
381 RNA-Seq data combined into a single *de novo* assembly with Trinity and a spliced
382 alignment indexed against the genome assembly with HISAT 2.1.0 (45). The PASA
383 pipeline v2.3.3 (46) was used to identify high-quality gene models that were used to
384 train the *ab initio* gene predictor in AUGUSTUS v3.3 (47) and GenMark. Additionally,
385 SwissProt protein data was aligned to the genome assembly using Exonerate (48) and
386 *L. luymesi* transcripts aligned using Minimap2 v2.1 (49). tRNA genes were identified
387 with tRNAscan-SE v1.3.1 (50). Finally, EvidenceModeler 1.1.0 (51) was used to
388 combine all evidence of gene prediction from protein alignments, transcript alignments,
389 and *ab initio* predictions to construct high-quality consensus gene models. Functional
390 annotations of predicted gene models were performed using curated databases: KEGG
391 orthology was assigned using the KEGG Automatic Annotation server (52), domain
392 structure by InterProScan (53), and protein identity with the SwissProt database.
393 Secreted proteins were predicted using SignalP (54) and Phobius (55) in InterProScan.

394 **Proteomics characterization**

395 Proteomic analysis was performed by Proteomics & Metabolomics Facility at
396 Colorado State University. Briefly, trunk/trophosome tissue was cleaned and
397 homogenized. Protein in resulting supernatant was quantified by the Pierce BCA
398 Protein Assay Kit (ThermoFisher-Pierce). Absorbance was measured at 550nm and
399 using a Bovine Serum Albumin (BSA) control. 50 µg total protein was processed for in-
400 solution trypsin digestion (56). Tandem mass spectra were extracted, charge state
401 deconvoluted and deisotoped by ProteoWizard MsConvert (version 3.0). Spectra
402 searched against gene models of *L. luymesi* host (herein) and symbiont genomes ((4))
403 using Mascot (Matrix Science, London, UK; version 2.6.0) with a fragment ion mass

404 tolerance of 0.80 Da and a parent ion tolerance of 20 PPM. Search results assessed
405 with probabilistic protein identification algorithms (57) implemented in the Scaffold
406 software v. 4.8.4, (Proteome Software Inc., Portland, OR; (58). Protein identifications
407 required >99.0% probability (with Protein Prophet algorithm; (59) and presence of ≥ 1
408 identified peptide. Proteins that contained similar peptides and could not be
409 differentiated based on MS/MS analysis alone were grouped (SI methods).

410 **Gene family analysis**

411 Following all-to-all Diamond v1.0 (60) BLASTP searches against 22 selected
412 lophotrochozoan proteomes (Table S3), orthology groups (OGs) were identified using
413 Orthofinder with a default inflation parameter ($I=1.5$). Gene ontology annotation used
414 PANTHER v13.1 (61) with the PANTHER HMM scoring tool (pantherScore2.pl). Gene
415 family expansion and contraction was estimated using CAFÉ v2.1 (62). For each gene
416 family, CAFÉ generated a family-wide P value, with a significant P value indicating a
417 possible gene-family expansion or contraction event. Significantly expanded gene
418 families ($p < 0.05$) were then identified by InterProscan.

419 **Manual annotation of gene families**

420 In addition to the annotation pipeline mentioned above, we manually annotated
421 genes of interest herein: hemoglobin gene families, genes related to amino acid
422 synthesize, immunity, and longevity. These gene families were specifically selected *a*
423 *priori* based on our experience and review of available publications in the field. See *SI*/
424 *methods* for detailed procedure.

425

426 **Data Availability**

427 Raw reads, assembled genome sequences and annotation are accessible from
428 NCBI under BioProject numbers PRJNA516467, Sequence Read Archive accession
429 numbers SRR851910-SPR851919 and Whole Genome Shotgun project numbers
430 SDWI00000000. The genome annotations, proteomic results, scripts and data for the

431 analyses are available from the Github Repository at
432 <https://github.com/yzl0084/Lamellibrachia-genome>.

433 **Acknowledgments**

434 This study was supported by awards from the U.S. National Science Foundation
435 (NSF) (DEB-1036537 and IOS-0843473 to KMH, Scott Santos and DanThornhill). YL
436 was supported by the China Scholarship Council (CSC). We thank Chris Little, Maggie
437 Georgieva, Luke Parry, and Jason Flores for the helpful discussions. We thank Zack
438 and Ian Gilman for help with revising the manuscript. We thank Jon Palmer helped
439 troubleshoot the Funannoate pipeline. We thank Kitty Brown for help with proteomic
440 data interpretation. Bioinformatic analyses were conducted on the Auburn University
441 Molette Laboratory SkyNet server, Auburn University Hopper HPC system, and the
442 Alabama Supercomputer Authority. This is Molette Biology Laboratory contribution ###
443 and Auburn University Marine Biology Program contribution ###.

444 **Author Contributions**

445 YL and KMH designed research; YL, MGT, DSW, VEB, KTD and KMH
446 performed research and data analysis; YL, MGT and KMH wrote the paper. All authors
447 contributed to revise the paper.

448

449 **References**

- 450 1. Moran NA (2007) Symbiosis as an adaptive process and source of phenotypic
451 complexity. *Proc Natl Acad Sci U S A* 104 Suppl 1:8627–8633.
- 452 2. Boetius A (2005) Microfauna-macrofauna interaction in the seafloor: lessons from
453 the tubeworm. *PLoS Biol* 3(3):e102.
- 454 3. Nussbaumer AD, Fisher CR, Bright M (2006) Horizontal endosymbiont transmission
455 in hydrothermal vent tubeworms. *Nature* 441(7091):345.
- 456 4. Li Y, Liles MR, Halanych KM (2018) Endosymbiont genomes yield clues of
457 tubeworm success. *ISME J* 12(11):2785.

- 458 5. Markert S, et al. (2007) Physiological proteomics of the uncultured endosymbiont of
459 *Riftia pachyptila*. *Science* 315(5809):247–250.
- 460 6. Zal F, et al. (1997) Primary structure of the common polypeptide chain b from the
461 multi-hemoglobin system of the hydrothermal vent tube worm *Riftia pachyptila*: An
462 insight on the sulfide binding-site. *Proteins: Struct Funct Bioinf* 29(4):562–574.
- 463 7. Arp AJ, Childress JJ (1981) Blood function in the hydrothermal vent vestimentiferan
464 tube worm. *Science* 213(4505):342–344.
- 465 8. Zal F, Lallier FH, Green BN, Vinogradov SN, Toulmond A (1996) The multi-
466 hemoglobin system of the hydrothermal vent tube worm *Riftia pachyptila* II.
467 Complete polypeptide chain composition investigated by maximum entropy analysis
468 of mass spectra. *J Biol Chem* 271(15):8875–8881.
- 469 9. Flores JF, et al. (2005) Sulfide binding is mediated by zinc ions discovered in the
470 crystal structure of a hydrothermal vent tubeworm hemoglobin. *Proceedings of the*
471 *National Academy of Sciences* 102(8):2713–2718.
- 472 10. Lutz RA, et al. (1994) Rapid growth at deep-sea vents. *Nature* 371(6499):663.
- 473 11. Bergquist DC, Williams FM, Fisher CR (2000) Longevity record for deep-sea
474 invertebrate. *Nature* 403(6769):499.
- 475 12. Quesada V, et al. (2018) Giant tortoise genomes provide insights into longevity and
476 age-related disease. *Nature ecology & evolution*:1.
- 477 13. Chu H, Mazmanian SK (2013) Innate immune recognition of the microbiota
478 promotes host-microbial symbiosis. *Nat Immunol* 14(7):668.
- 479 14. Minic Z, Hervé G (2003) Arginine metabolism in the deep sea tube worm *Riftia*
480 *pachyptila* and its bacterial endosymbiont. *J Biol Chem*.
- 481 15. Tian R-M, et al. (2017) Genome Reduction and Microbe-Host Interactions Drive
482 Adaptation of a Sulfur-Oxidizing Bacterium Associated with a Cold Deep Sponge.
483 *mSystems* 2(2). doi:10.1128/mSystems.00184-16.
- 484 16. Tokuda G, et al. (2013) Maintenance of essential amino acid synthesis pathways in
485 the Blattabacterium cuenoti symbiont of a wood-feeding cockroach. *Biol Lett*
486 9(3):20121153.
- 487 17. Newton ILG, Girguis PR, Cavanaugh CM (2008) Comparative genomics of
488 vesicomyid clam (Bivalvia: Mollusca) chemosynthetic symbionts. *BMC Genomics*
489 9(1):585.
- 490 18. Ponnudurai R, et al. (2017) Metabolic and physiological interdependencies in the
491 *Bathymodiolus azoricus* symbiosis. *ISME J* 11(2):463.
- 492 19. Jäckle O, et al. (2019) Chemosynthetic symbiont with a drastically reduced genome

- 493 serves as primary energy storage in the marine flatworm *Paracatenula*. *Proc Natl
494 Acad Sci U S A.* doi:10.1073/pnas.1818995116.
- 495 20. Bright M, Keckeis H, Fisher CR (2000) An autoradiographic examination of carbon
496 fixation, transfer and utilization in the *Riftia pachyptila* symbiosis. *Mar Biol*
497 136(4):621–632.
- 498 21. Bailly X, et al. (2002) Evolution of the sulfide-binding function within the globin
499 multigenic family of the deep-sea hydrothermal vent tubeworm *Riftia pachyptila*. *Mol
500 Biol Evol* 19(9):1421–1433.
- 501 22. Park BS, Lee J-O (2013) Recognition of lipopolysaccharide pattern by TLR4
502 complexes. *Exp Mol Med* 45(12):e66.
- 503 23. Bates JM, Akerlund J, Mittge E, Guillemin K (2007) Intestinal alkaline phosphatase
504 detoxifies lipopolysaccharide and prevents inflammation in zebrafish in response to
505 the gut microbiota. *Cell Host Microbe* 2(6):371–382.
- 506 24. Maggio M, Guralnik JM, Longo DL, Ferrucci L (2006) Interleukin-6 in aging and
507 chronic disease: a magnificent pathway. *J Gerontol A Biol Sci Med Sci* 61(6):575–
508 584.
- 509 25. Melov S, et al. (2000) Extension of life-span with superoxide dismutase/catalase
510 mimetics. *Science* 289(5484):1567–1569.
- 511 26. Bright M, Bulgheresi S (2010) A complex journey: transmission of microbial
512 symbionts. *Nat Rev Microbiol* 8(3):218–230.
- 513 27. Andrews S, Others (2010) FastQC: a quality control tool for high throughput
514 sequence data.
- 515 28. Bolger AM, Lohse M, Usadel B (2014) Trimmomatic: a flexible trimmer for Illumina
516 sequence data. *Bioinformatics* 30(15):2114–2120.
- 517 29. Marçais G, Kingsford C (2011) A fast, lock-free approach for efficient parallel
518 counting of occurrences of k-mers. *Bioinformatics* 27(6):764–770.
- 519 30. Vulture GW, et al. (2017) GenomeScope: fast reference-free genome profiling from
520 short reads. *Bioinformatics* 33(14):2202–2204.
- 521 31. O'Connell J, et al. (2015) NxTrim: optimized trimming of Illumina mate pair reads.
522 *Bioinformatics* 31(12):2035–2037.
- 523 32. Kajitani R, et al. (2014) Efficient de novo assembly of highly heterozygous genomes
524 from whole-genome shotgun short reads. *Genome Res*:gr–170720.
- 525 33. Boetzer M, Pirovano W (2014) SSPACE-LongRead: scaffolding bacterial draft
526 genomes using long read sequence information. *BMC Bioinformatics* 15(1):211.
- 527 34. Luo R, et al. (2012) SOAPdenovo2: an empirically improved memory-efficient short-

- 528 read de novo assembler. *Gigascience* 1(1):18.
- 529 35. Pryszz LP, Gabaldón T (2016) Redundans: an assembly pipeline for highly
530 heterozygous genomes. *Nucleic Acids Res* 44(12):e113–e113.
- 531 36. Gurevich A, Saveliev V, Vyahhi N, Tesler G (2013) QUAST: quality assessment tool
532 for genome assemblies. *Bioinformatics* 29(8):1072–1075.
- 533 37. Waterhouse RM, et al. (2017) BUSCO applications from quality assessments to
534 gene prediction and phylogenomics. *Mol Biol Evol* 35(3):543–548.
- 535 38. Haas BJ, et al. (2013) De novo transcript sequence reconstruction from RNA-seq
536 using the Trinity platform for reference generation and analysis. *Nat Protoc*
537 8(8):1494.
- 538 39. Li W, Godzik A (2006) Cd-hit: a fast program for clustering and comparing large
539 sets of protein or nucleotide sequences. *Bioinformatics* 22(13):1658–1659.
- 540 40. Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nat
541 Methods* 9(4):357.
- 542 41. Li B, Dewey CN (2011) RSEM: accurate transcript quantification from RNA-Seq
543 data with or without a reference genome. *BMC Bioinformatics* 12(1):323.
- 544 42. Albertin CB, et al. (2015) The octopus genome and the evolution of cephalopod
545 neural and morphological novelties. *Nature* 524(7564):220.
- 546 43. Smit AFA, Hubley R (2008) RepeatModeler Open-1.0. Available fom [http://www
547 repeatmasker.org](http://www.repeatmasker.org).
- 548 44. Chen N (2004) Using RepeatMasker to identify repetitive elements in genomic
549 sequences. *Curr Protoc Bioinformatics* 5(1):4–10.
- 550 45. Kim D, Langmead B, Salzberg SL (2015) HISAT: a fast spliced aligner with low
551 memory requirements. *Nat Methods* 12(4):357.
- 552 46. Haas BJ, et al. (2003) Improving the *Arabidopsis* genome annotation using maximal
553 transcript alignment assemblies. *Nucleic Acids Res* 31(19):5654–5666.
- 554 47. Stanke M, et al. (2006) AUGUSTUS: ab initio prediction of alternative transcripts.
555 *Nucleic Acids Res* 34(suppl_2):W435–W439.
- 556 48. Slater GSC, Birney E (2005) Automated generation of heuristics for biological
557 sequence comparison. *BMC Bioinformatics* 6(1):31.
- 558 49. Li H (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
559 1:7.
- 560 50. Lowe TM, Eddy SR (1997) tRNAscan-SE: a program for improved detection of
561 transfer RNA genes in genomic sequence. *Nucleic Acids Res* 25(5):955.

- 562 51. Haas BJ, et al. (2008) Automated eukaryotic gene structure annotation using
563 EVidenceModeler and the Program to Assemble Spliced Alignments. *Genome Biol*
564 9(1):1.
- 565 52. Moriya Y, Itoh M, Okuda S, Yoshizawa AC, Kanehisa M (2007) KAAS: an automatic
566 genome annotation and pathway reconstruction server. *Nucleic Acids Res*
567 35(suppl_2):W182–W185.
- 568 53. Zdobnov EM, Apweiler R (2001) InterProScan--an integration platform for the
569 signature-recognition methods in InterPro. *Bioinformatics* 17(9):847–848.
- 570 54. Petersen TN, Brunak S, von Heijne G, Nielsen H (2011) SignalP 4.0: discriminating
571 signal peptides from transmembrane regions. *Nat Methods* 8(10):785.
- 572 55. Käll L, Krogh A, Sonnhammer ELL (2007) Advantages of combined transmembrane
573 topology and signal peptide prediction—the Phobius web server. *Nucleic Acids Res*
574 35(suppl_2):W429–W432.
- 575 56. Schauer KL, Freund DM, Prenni JE, Curthoys NP (2013) Proteomic profiling and
576 pathway analysis of the response of rat renal proximal convoluted tubules to
577 metabolic acidosis. *American Journal of Physiology-Renal Physiology* 305(5):F628–
578 F640.
- 579 57. Keller A, Nesvizhskii AI, Kolker E, Aebersold R (2002) Empirical statistical model to
580 estimate the accuracy of peptide identifications made by MS/MS and database
581 search. *Anal Chem* 74(20):5383–5392.
- 582 58. Searle BC, Turner M, Nesvizhskii AI (2008) Improving sensitivity by probabilistically
583 combining results from multiple MS/MS search methodologies. *J Proteome Res*
584 7(1):245–253.
- 585 59. Nesvizhskii AI, Keller A, Kolker E, Aebersold R (2003) A statistical model for
586 identifying proteins by tandem mass spectrometry. *Anal Chem* 75(17):4646–4658.
- 587 60. Buchfink B, Xie C, Huson DH (2014) Fast and sensitive protein alignment using
588 DIAMOND. *Nat Methods* 12(1):59.
- 589 61. Mi H, et al. (2016) PANTHER version 11: expanded annotation data from Gene
590 Ontology and Reactome pathways, and data analysis tool enhancements. *Nucleic
591 Acids Res* 45(D1):D183–D189.
- 592 62. De Bie T, Cristianini N, Demuth JP, Hahn MW (2006) CAFE: a computational tool
593 for the study of gene family evolution. *Bioinformatics* 22(10):1269–1271.
- 594
- 595
- 596
- 597

PNAS

www.pnas.org

1
2
3 Supplementary Information for

4 **Genomic adaptations to chemosymbiosis in the deep-sea seep-dwelling**
5 **tubeworm *Lamellibrachia luymesi* (Siboglinidae, Annelida)**

6
7 Yuanning Li,*¹, Michael G. Tassia¹, Damien S. Waits¹, Viktoria E. Bogantes¹, Kyle T. David¹,
8 Kenneth M. Halanych¹

9
10 Corresponding authors: Yuanning Li, Kenneth M. Halanych
11 Email: yuanning.li@yale.edu; ken@auburn.edu

12
13 **This PDF file includes:**

14 Supplementary text
15 Figs. S1 to S8
16 Tables S1 to S10

17 **Other supplementary materials for this manuscript include the following:**

18 Dataset S1

19
20
21
22
23
24
25
26
27
28
29

30

31 **Supplementary Information Text**

32 **SI Methods**

33 **Proteomics characterization.**

34 Proteomic analysis of *Lamellibrachia luymesii* trunk/trophosome tissue was performed by
35 Proteomics & Metabolomics Facility at Colorado State University. Here we restate the protocol
36 provided by Colorado State University. 50 µg total protein was aliquoted from each sample and
37 processed for in-solution trypsin digestion as previously described (1). A total of 0.5µg of
38 peptides were then purified and concentrated using an on-line enrichment column (Waters
39 Symmetry Trap C18 100Å, 5µm, 180 µm ID x 20mm column). Subsequent chromatographic
40 separation was performed on a reverse phase nanospray column (Waters, Peptide BEH C18;
41 1.7µm, 75 µm ID x 150mm column, 45°C) using a 90 minute gradient: 5%-30% buffer B over 85
42 minutes followed by 30%-45%B over 5 minutes (0.1% formic acid in ACN) at a flow rate of 350
43 nanoliters/min. Peptides were eluted directly into the mass spectrometer (Orbitrap Velos Pro,
44 Thermo Scientific) equipped with a Nanospray Flex ion source (Thermo Scientific) and spectra
45 were collected over a m/z range of 400–2000 under positive mode ionization. Ions with charge
46 state +2 or +3 were accepted for MS/MS using a dynamic exclusion limit of 2 MS/MS spectra of
47 a given m/z value for 30 s (exclusion duration of 90 s). The instrument was operated in FT mode
48 for MS detection (resolution of 60,000) and ion trap mode for MS/MS detection with a
49 normalized collision energy set to 35%. Compound lists of the resulting spectra were generated
50 using Xcalibur 3.0 software (Thermo Scientific) with a S/N threshold of 1.5 and 1 scan/group.

51 Tandem mass spectra were extracted, charge state deconvoluted and deisotoped by
52 ProteoWizard MsConvert v3.0. Spectra from all samples were searched using Mascot (Matrix
53 Science, London, UK; version 2.6.0) against gene models of *Lamellibrachia* host and symbiont
54 genomes (derived from (2)) assuming the digestion enzyme trypsin. Mascot was searched with
55 a fragment ion mass tolerance of 0.80 Da and a parent ion tolerance of 20 PPM. Oxidation of
56 methionine and carbamidomethyl of cysteine were specified in Mascot as variable modifications.
57 Search results from all samples were imported and combined using the probabilistic protein
58 identification algorithms (3) implemented in the Scaffold software (version Scaffold_4.8.4,
59 Proteome Software Inc., Portland, OR) (4). Protein identifications were accepted if they could be
60 established at greater than 99.0% probability and contained at least 1 identified peptide. Protein
61 probabilities were assigned by the Protein Prophet algorithm (5). Proteins that contained similar
62 peptides and could not be differentiated based on MS/MS analysis alone were grouped to
63 satisfy the principles of parsimony.

64

65 **Manual annotation of gene families with potential interest.**

66 We manually annotated some gene families of interest including hemoglobin gene families,
67 genes related to amino acid synthesize, immunity function and longevity. Hbs and linker
68 sequences (Fig. S4) of interest were obtained from *L. luymesigenome* and assembled siboglinid
69 transcriptomes derived from previous studies (6, 7) via diamond BLASTP (evalue cutoff 1E-5) by
70 using bait Hb and linker bait sequences acquired from *RiftiaHbs* (downloaded from SwissProt
71 Database). Sequences with best hits to target proteins were annotated for protein domain
72 architecture using the Pfam databases included in InterProScan. After manual removal of
73 redundant and incorrect sequences (e.g., sequences are too short or lack of globin domain), we
74 used MAFFT 7.2.15 (8) to align Hb amino acid sequences. Maximum likelihood analyses were
75 performed in IQTree v1.5 (9) under the best-fitting models for associated partition schemes

76 determined by ModelFinder implemented in IQTree with ultrafast bootstrapping of 1000
77 replicates.

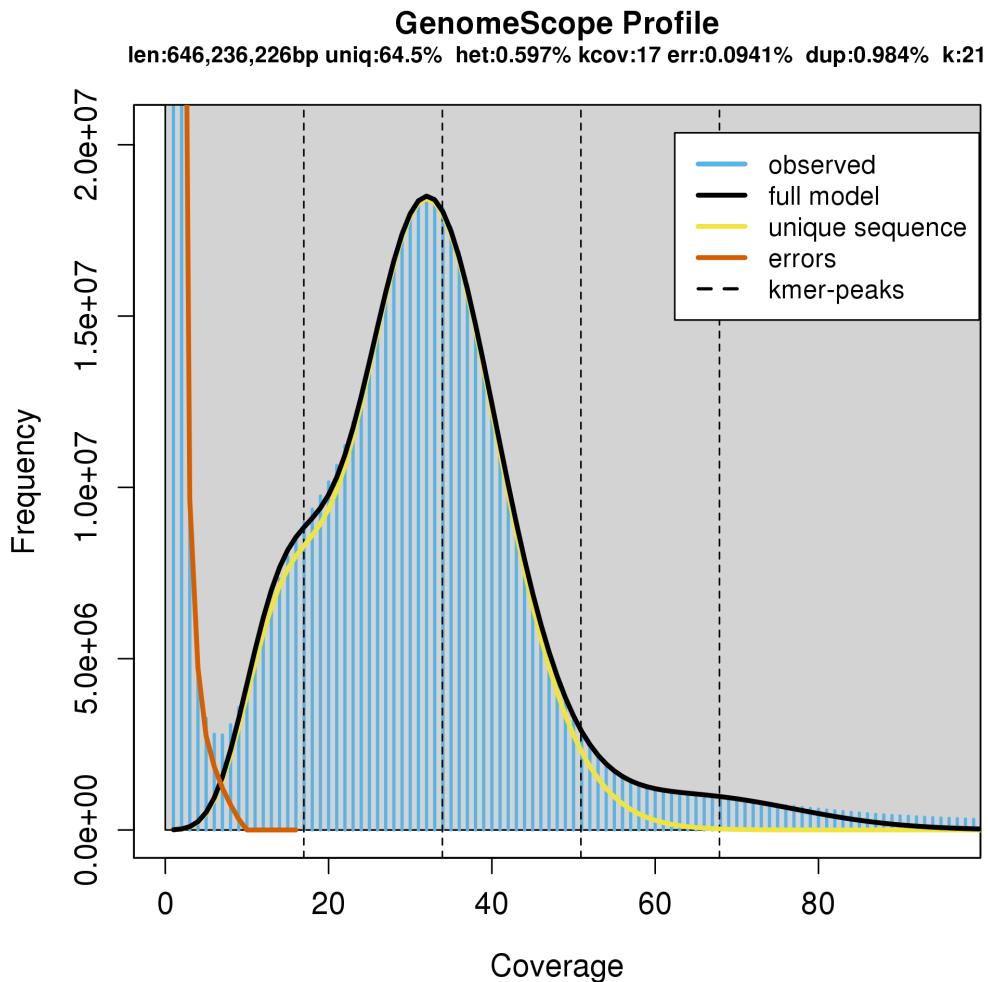
78 Discovery of SODs and immunity-related genes largely follows the same workflow as used for
79 Hbs. For immunity genes, targeted genes were additionally processed through the
80 Extract_Homologs2 script used in (10) (available at
81 https://github.com/mtassia/Homolog_identification). We examined major signaling components
82 of the TLR signaling pathway, as well as RLRs, NFkB-associated proteins and interferon-
83 regulatory factors. We only included identification of TLR and RIGs signaling components in the
84 manuscript as other immunity related genes did not clearly reveal any evolutionary patterns of
85 interest across lophotrochozoans (Table S9). Importantly, the Extract_Homologs2 script
86 identifies unique protein sequences within an amino acid dataset that fall within user-defined
87 domain architecture criteria (Table S10). Due to this stringency, the pipeline only identifies the
88 complement of unique proteins for any target family encoded in a genome. Full amino acid
89 sequences for TLRs were placed in a phylogenetic context using the bioinformatic workflow
90 delineated above for Hbs.

91 Searches for genes related to amino acids synthesis from *Lamellibrachia*, *Lamellibrachia*
92 symbionts, and *Capitella teleta* genomes were performed by using the KEGG2 KAAS genome
93 annotation web server and then visualized by the KEGG Mapper Reconstruct Pathway. A
94 BLASTp search of protein sequences from the genome annotation were queried against the
95 Swiss-Prot database was used to search and supplement for proteins that were missing in the
96 visualized KEGG pathway.
97 The results of gene alignments, tree files mentioned above were available at the Github
98 Repository (<https://github.com/yzl0084/Lamellibrachia-genome>).
99

100 References

- 101 1. Schauer KL, Freund DM, Prenni JE, Curthoys NP (2013) Proteomic profiling and pathway
102 analysis of the response of rat renal proximal convoluted tubules to metabolic acidosis.
103 *American Journal of Physiology-Renal Physiology* 305(5):F628–F640.
- 104 2. Li Y, Liles MR, Halanych KM (2018) Endosymbiont genomes yield clues of tubeworm success.
105 *ISME J* 12(11):2785.
- 106 3. Keller A, Nesvizhskii AI, Kolker E, Aebersold R (2002) Empirical statistical model to estimate the
107 accuracy of peptide identifications made by MS/MS and database search. *Anal Chem*
108 74(20):5383–5392.
- 109 4. Searle BC, Turner M, Nesvizhskii AI (2008) Improving sensitivity by probabilistically combining
110 results from multiple MS/MS search methodologies. *J Proteome Res* 7(1):245–253.
- 111 5. Nesvizhskii AI, Keller A, Kolker E, Aebersold R (2003) A statistical model for identifying proteins
112 by tandem mass spectrometry. *Anal Chem* 75(17):4646–4658.
- 113 6. Li Y, et al. (2017) Phylogenomics of tubeworms (Siboglinidae, Annelida) and comparative
114 performance of different reconstruction methods. *Zool Scr* 46(2):200–213.
- 115 7. Waits DS, Santos SR, Thornhill DJ, Li Y, Halanych KM (2016) Evolution of Sulfur Binding by
116 Hemoglobin in Siboglinidae (Annelida) with Special Reference to Bone-Eating Worms, Osedax.
117 *J Mol Evol* 82(4-5):219–229.

- 118 8. Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7:
119 improvements in performance and usability. *Mol Biol Evol* 30(4):772–780.
- 120 9. Chernomor O, von Haeseler A, Minh BQ (2016) Terrace Aware Data Structure for
121 Phylogenomic Inference from Supermatrices. *Syst Biol* 65(6):997–1008.
- 122 10. Tassia MG, Whelan NV, Halanych KM (2017) Toll-like receptor pathway evolution in
123 deuterostomes. *Proceedings of the National Academy of Sciences* 114(27):7055–7060.
124
125
126
127
128
129
130
131
132
133
134
135
136
137



138
139 **Figure S1.** Estimation of genome size, repetitive content and level of heterozygosity
140 from 100 million Illumina paired-end reads.
141

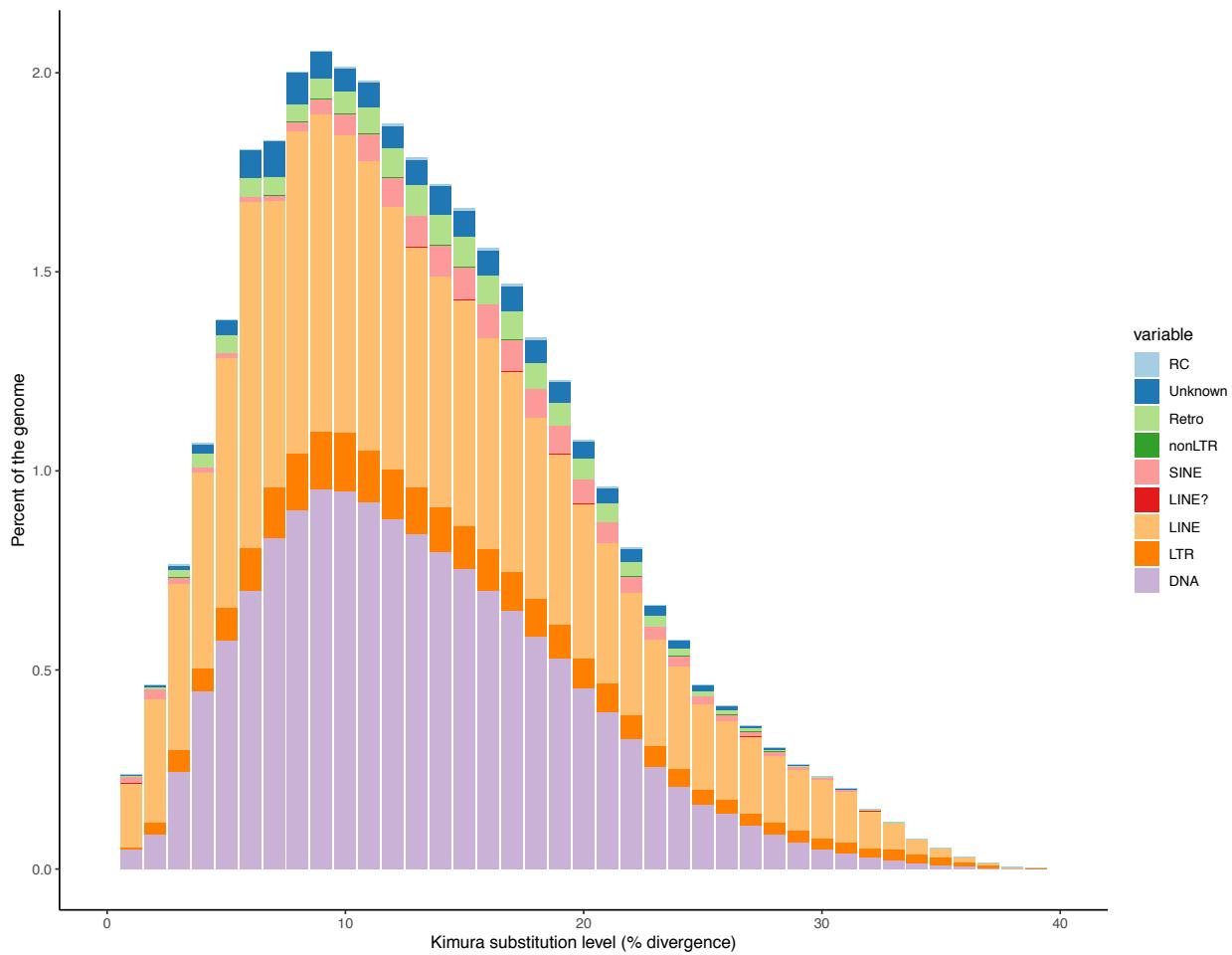


Figure S2. History from major superfamilies of transposable elements in the *Lamellibrachia* genome. Kiruma distances are arranged from value 0 representing recent TE copies to 40 for ancient TE insertions.

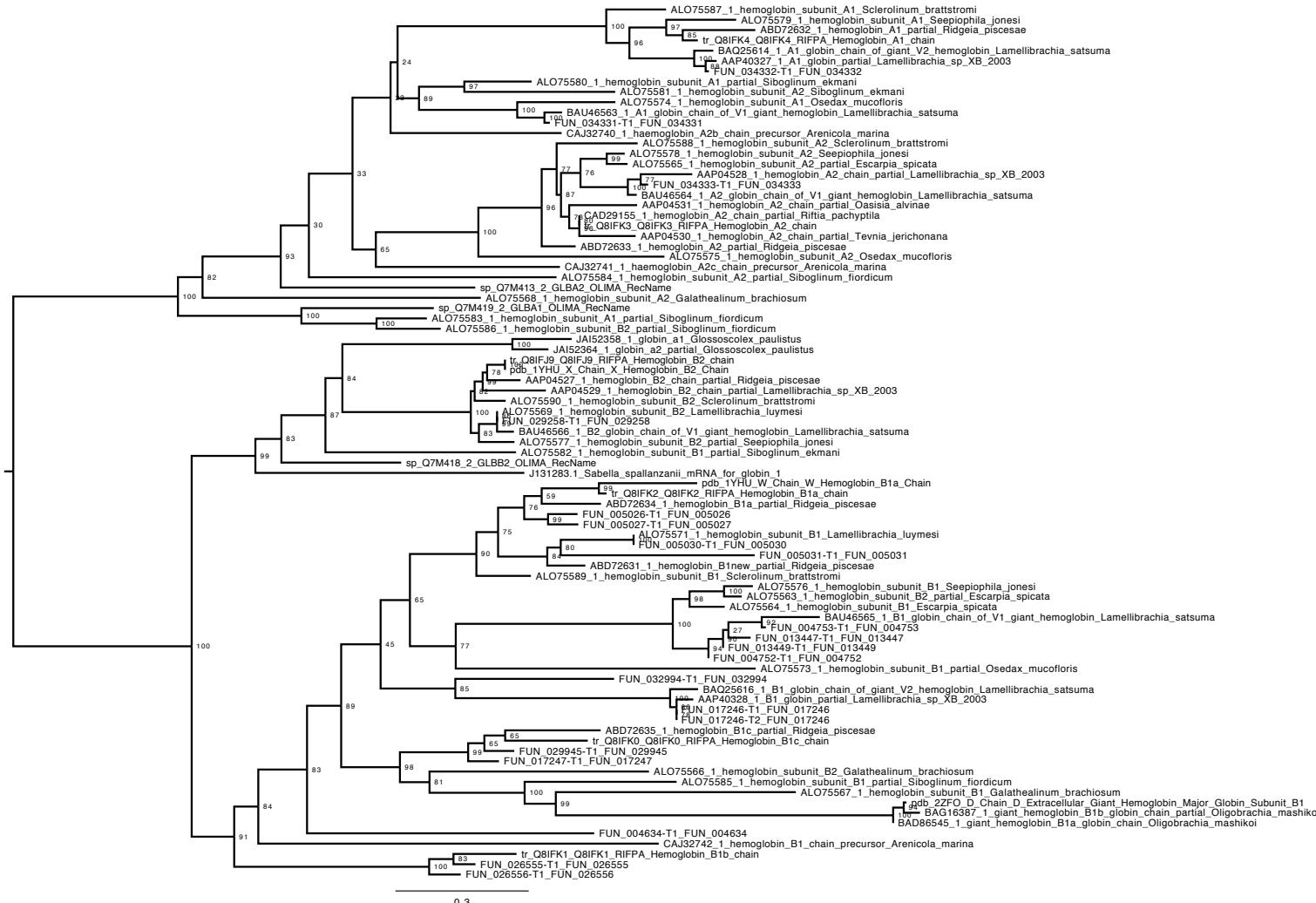


Figure S3. Siboglinid hemoglobin maximum-likelihood tree reconstructed with IQtree with midpoint rooting with 1000 ultrafast bootstraps using LG model. Bootstrap support values are shown at the relevant node. GenBank accession numbers are listed on the terminal nodes.

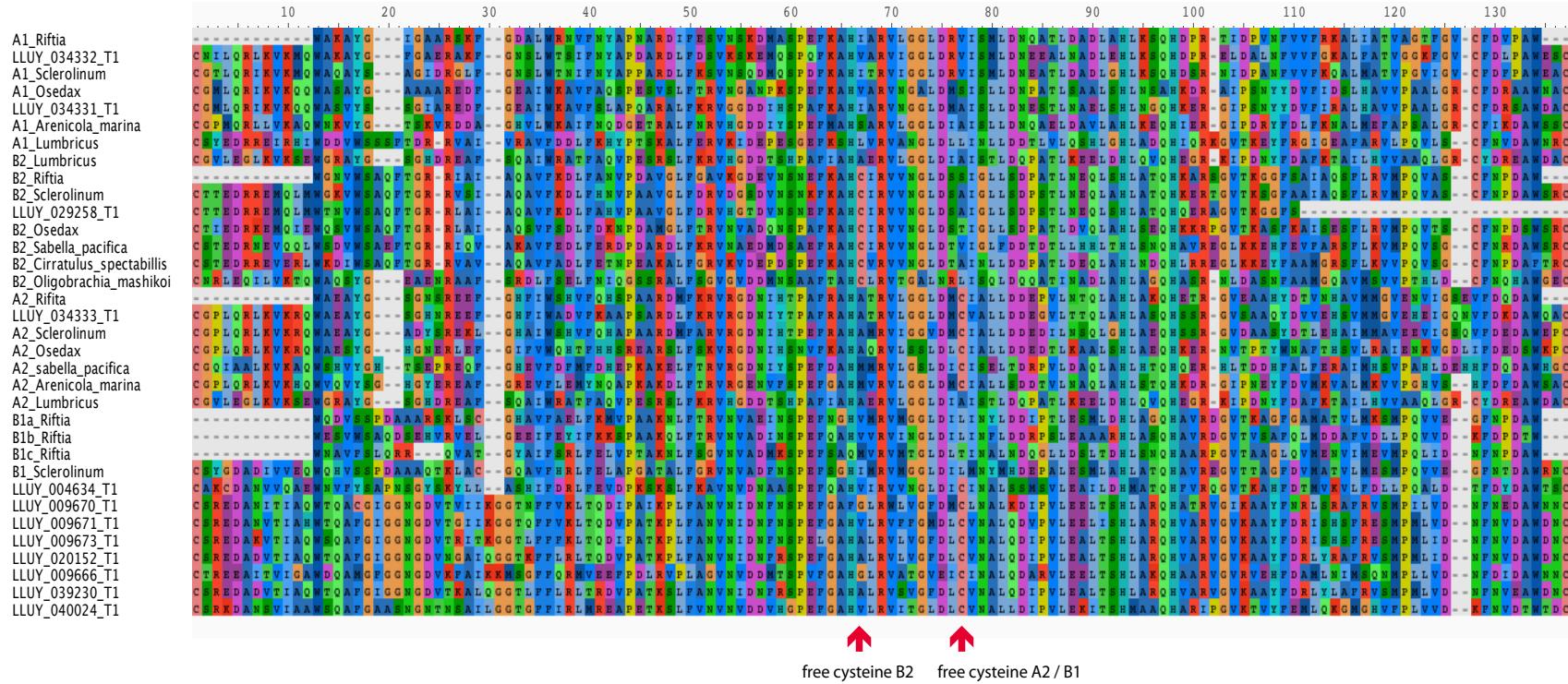


Figure S4. Hemoglobin gene diversity in *Lamellibrachia luymesi*. (Partial alignment of sampled siboglinid Hb subunit A1, A2, B1, B2 sequences. Red arrows indicate positions contain free cysteines or cysteine residues in HB B2 chains, and B1/A2 chains, respectively.

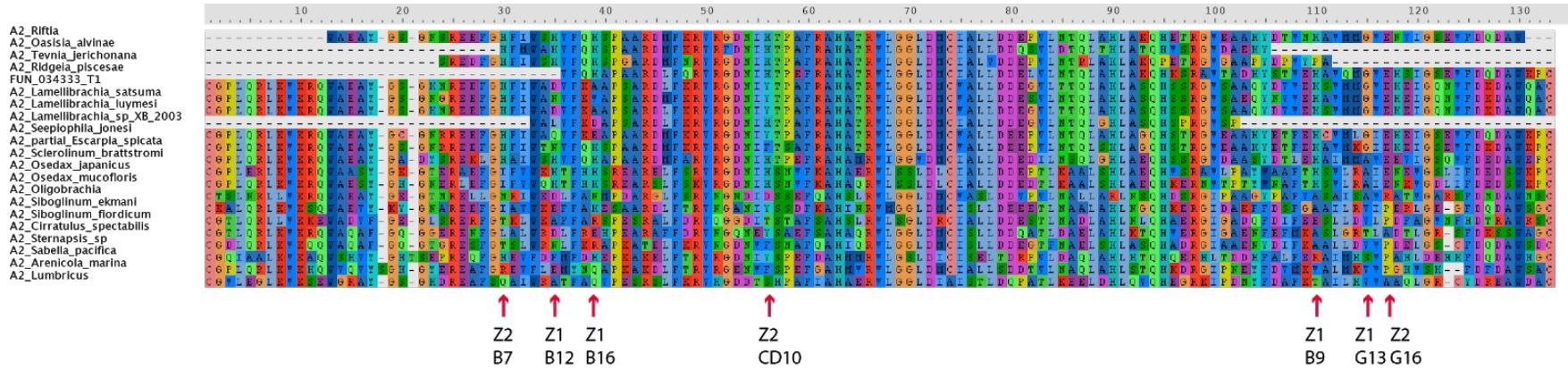


Figure S5. Partial alignment of sampled siboglinid HB subunit A2 sequences. Red arrows indicate amino acid residues at the interface between pairs of A2 chains with zinc moieties for H₂S binding.

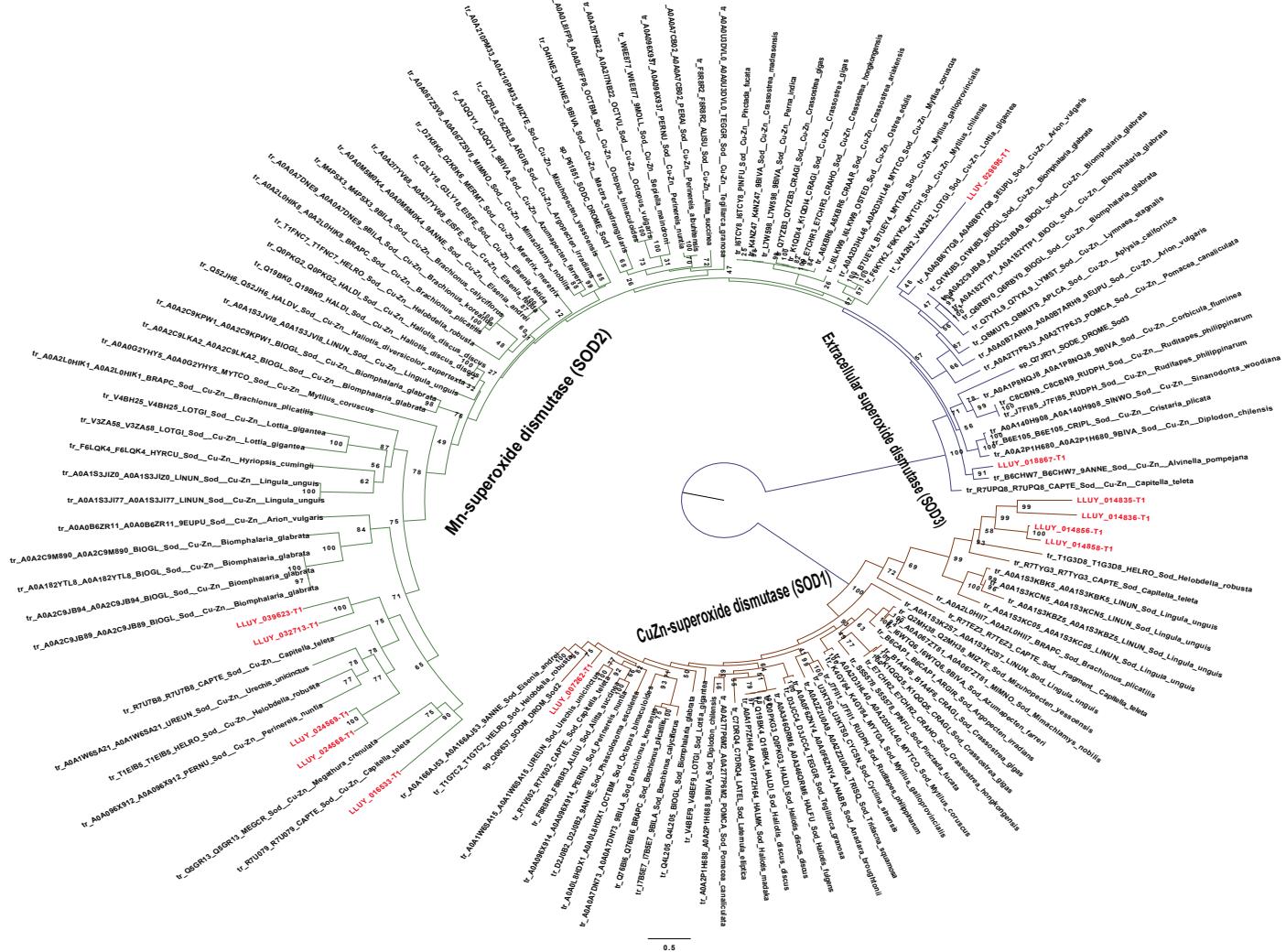


Figure S6. Lophotrochozoan SOD maximum-likelihood tree reconstructed with IQtree with midpoint rooting with 1000 ultrafast bootstraps using LG model. Bootstrap support values are shown at the relevant node. GenBank accession numbers are listed on the terminal nodes. GenBank accession numbers are next to the tip names.

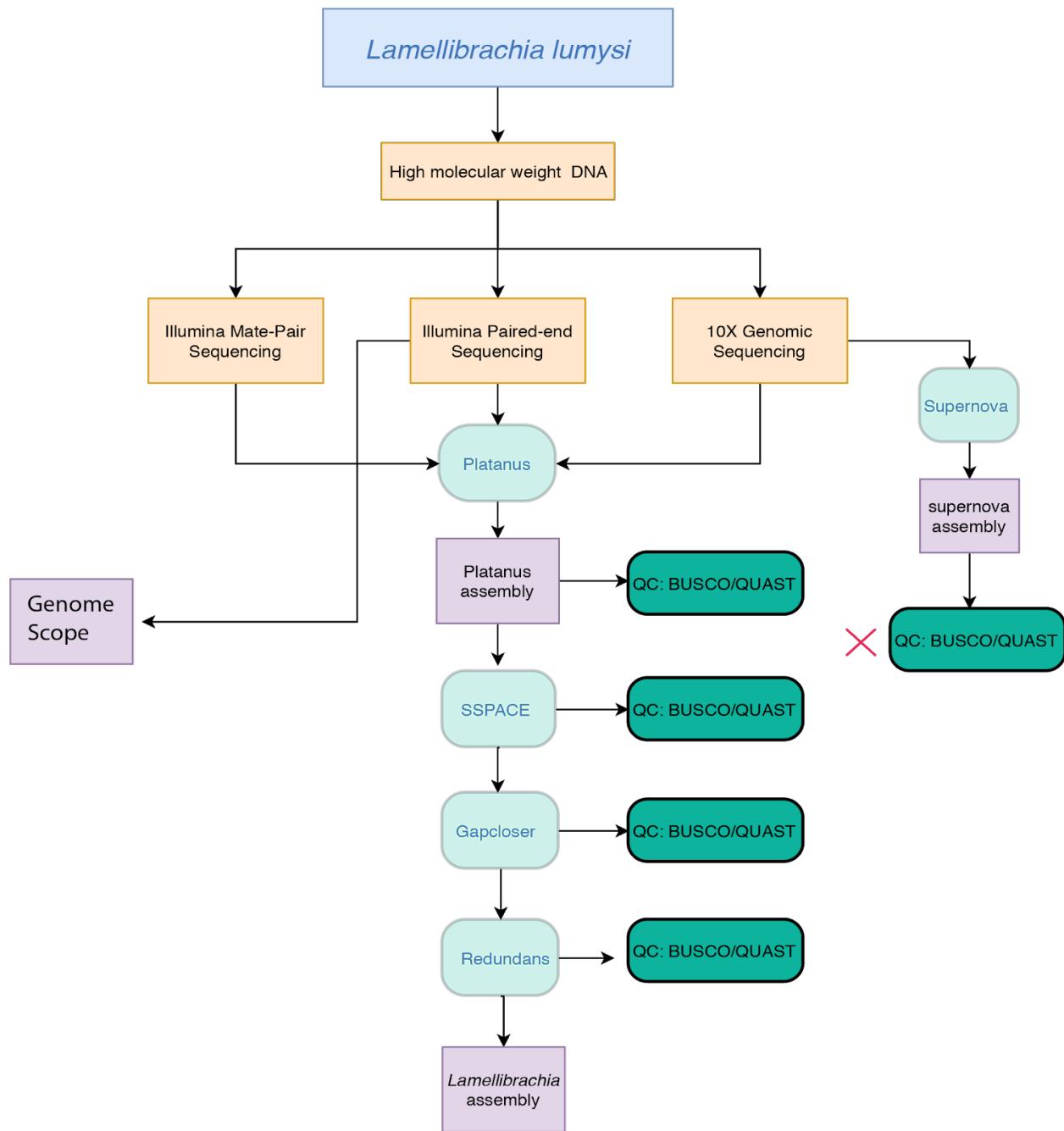


Figure S7. Workflow of *Lamellibrachia luymesi* genome assembly. 10X genomics assembly alone provide worse assembly and failed QC compared to *Platanus* indicated by red X.

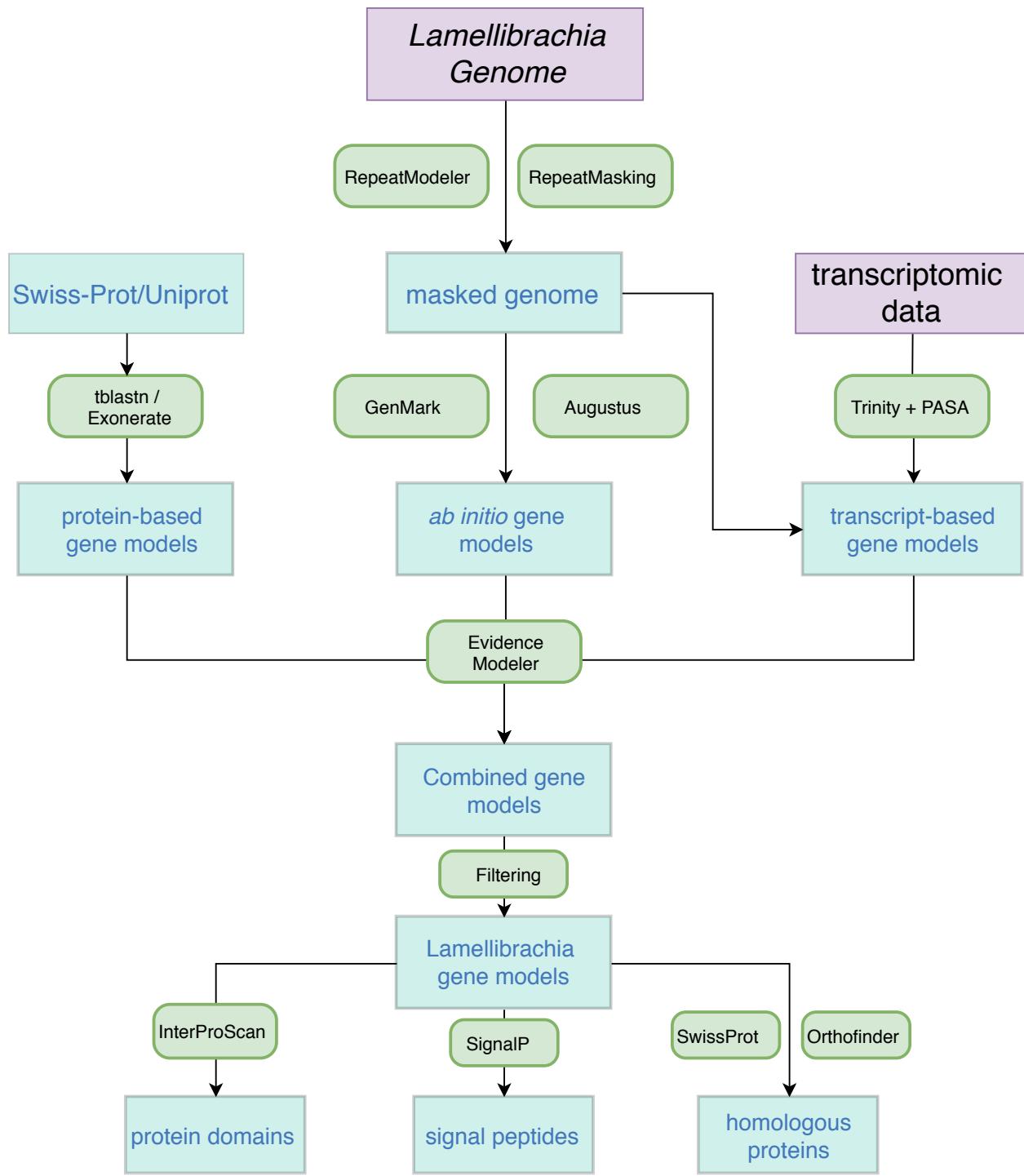


Figure S8. Workflow of *Lamellibrachia* genome annotation pipeline using Funannotate.

Table S1. Sequencing information of *Lamellibrachia luymesi* genome.

Tissue	Data Type	Sequencing Chemistry	Total Read Number	Lab Accession	Accession Number	Coverage (X)
Vestimentum	Genomics	10X Genomics	648,546,716	KH-4260-0006	SRR8519110	141.48
Vestimentum	Genomics	180 bp Paired-end	530,601,282	SL84794	SRR8519115	96.43
Vestimentum	Genomics	180 bp Paired-end	318,356,186	SL115013	SRR8519114	57.86
Vestimentum	Genomics	400 bp Paired-end	237,879,494	SL84795	SRR8519113	43.12
Vestimentum	Genomics	750 bp Paired-end	118,008,914	SL84796	SRR8519112	21.40
Vestimentum	Genomics	3-5 kbp Mate-pair	344,803,888	SL85812	SRR8519119	60.77
Vestimentum	Genomics	5-7 kbp Mate-pair	352,639,094	SL85813	SRR8519118	64.04
Plume	Transcriptome	Paired-end	58,660,044	SL85796	SRR8519117	
Trophosome	Transcriptome	Paired-end	75,640,660	SL85798	SRR8519111	
Vestimentum	Transcriptome	Paired-end	50,537,812	SL85797	SRR8519116	

Table S2.Genome assembly and BUSCO statistics of *Lamellibrachia luymesi* compared to other lophotrochozoan genomes.

Taxon	Species Name	# contigs	Total length	Largest contig	GC (%)	N50	BUSCO (%)				
							Complete	Single	Duplicate	Fragment	Missing
Annelida	<i>Lamellibrachia lumysi</i>	11,871	687,711,696	2,117,112	40.16	372,990	95.80	93.00	2.80	2.90	1.30
	<i>Capitella teleta</i>	20,803	333,283,208	1,620,044	40.36	188,402	92.30	87.60	4.70	1.10	6.60
	<i>Hydroides elegans</i>	188,407	1,026,046,400	244,066	35.43	17,725	79.90	53.00	26.90	8.80	11.30
	<i>Helobdella robusta</i>	1,991	235,376,169	13,640,604	32.82	3,060,193	85.30	83.80	1.50	3.70	11.00
Phoronida	<i>Phoronis australis</i>	3,983	498,443,662	4,871,659	39.34	655,058	91.90	89.40	2.50	1.30	6.80
Nemertea	<i>Notospermus geniculatus</i>	11,108	858,599,399	1,576,180	42.85	239,235	91.90	89.40	2.50	1.30	6.80
Mollusca	<i>Crassostrea virginica</i>	10	684,723,884	104,168,038	34.83	75,944,018	90.70	88.10	2.60	0.80	8.50
Mollusca	<i>Crassostrea gigas</i>	7,658	557,717,710	1,964,558	33.42	402,213	90.80	85.70	5.10	0.90	8.30
Mollusca	<i>Bathymodiolus platifrons</i>	65,662	1,658,191,953	2,790,175	34.17	345,477	89.30	88.00	1.30	2.20	8.50
Mollusca	<i>Mytilus galloprovincialis</i>	1,002,334	1,500,149,602	67,529	31.77	3,239	91.20	63.00	28.20	0.90	7.90
Mollusca	<i>Octopus bimaculoides</i>	151,674	2,338,188,782	4,064,693	36.06	485,615	85.80	85.30	0.50	3.40	10.80
Mollusca	<i>Modiolus philippinarum</i>	74,573	2,629,556,424	715,382	33.96	100,386	85.20	82.70	2.50	4.90	9.90
Mollusca	<i>Mizuhopecten yessoensis</i>	82,658	987,568,220	7,498,238	36.52	827,226	89.80	87.80	2.00	1.20	9.00
Mollusca	<i>Lottia gigantea</i>	4,469	359,505,668	9,386,848	33.28	1,870,055	91.30	90.20	1.10	0.90	7.80
Brachiopoda	<i>Lingula anatina</i>	2,677	406,282,338	2,166,018	36.42	460,090	90.20	70.20	20.00	0.90	8.90
Rotifer	<i>Aplysia californica</i>	4,331	927,296,314	6,102,535	40.35	917,541	88.30	87.80	0.50	1.60	10.10

Table S3.

Proteomics and genome assemblies used in comparative analyses.

Taxon	Species	Genome source	RefSeq assembly accession
Annelida	<i>Lamellibrachia luymsi</i>	This study	SDWI00000000
	<i>Capitella teleta</i>	NCBI	GCA_000328365.1
	<i>Helobdella robusta</i>	NCBI	GCA_000326865.1
Mollusca	<i>Lottia gigantea</i>	NCBI	GCA_000327385.1
	<i>Octopus bimaculoides</i>	NCBI	GCA_001194135.1
	<i>Chlamys farreri</i>	NCBI	
	<i>Bathymodiolus platifrons</i>	NCBI	<u>GCA_002080005.1</u>
	<i>Biomphalaria glabrata</i>	NCBI	
	<i>Mizuhopecten yessoensis</i>	NCBI	<u>GCA_002113885.2</u>
	<i>Modiolus philippinarum</i>	NCBI	GCA_000457365.1
	<i>Patinopecten yessoensis</i>	NCBI	
	<i>Crassostrea gigas</i>	NCBI	GCF_000297895.1
	<i>Crassostrea virginica</i>	NCBI	<u>GCA_002022765.4</u>
Nemertea	<i>Notospermus geniculatus</i>	NCBI	<u>GCA_002633025.1</u>
Phoronida	<i>Phoronis australis</i>	NCBI	<u>GCA_002633005.1</u>
Brachiopoda	<i>Lingula anatina</i>	NCBI	<u>GCA_001039355.2</u>
Flatworm	<i>Schistosoma mansoni</i>	NCBI	<u>GCA_000237925.2</u>
	<i>Schmidtea mediterranea</i>	NCBI	<u>GCA_002600895.1</u>
	<i>Macrostomum lignano</i>	NCBI	<u>GCA_002269645.1</u>
	<i>Echinococcus multilocularis</i>	NCBI	<u>GCA_000469725.3</u>
Rotifera	<i>Aplysia californica</i>	NCBI	<u>GCA_000002075.2</u>
Ecdysozoa	<i>Diphania pulex</i>	NCBI	<u>GCA_000187875.1</u>
	<i>Drosophila melanogaster</i>	NCBI	<u>GCA_000001215.4</u>

Table S4.Repetitive element contained in the *Lamellibrachia luymesi* genome

	Subclass	Number of elements	length occupied (bp)	percentage of sequence
SINEs		56,783	8,181,892	1.19
LINEs		307,497	96,235,311	13.99
	LINE1	1,333	337,473	0.05
	LINE2	51,791	18,267,908	2.66
	L3/CR1	83,528	33,202,610	4.83
LTR elements		70,108	17,360,126	2.52
DNA elements		442,010	101,825,130	14.81
	hAT-Charlie	6,246	1,710,664	0.25
	TcMar-Tigger	4,078	1,496,397	0.22
Unclassified		95,876	17,185,149	2.5
Tot. interspersed repeats			240,787,608	35.01
Small RNA		101	27,655	0
Satellites		813	376,172	0.05
Simple repeats		211,949	13,781,625	2
Low complexity		8,125	682,495	0.1

Table S5.

PANTHER gene family annotation of gene families that are under expansion or contraction as identified by CAFE.

Orthology Group	Number of gain and loss	PANTHER gene family	PANTHER annotation
OG0000040	+23*	PTHR44025	LOW-DENSITY LIPOPROTEIN RECEPTOR-RELATED PROTEIN
OG0000044	+16*	PTHR12011	G-PROTEIN COUPLED RECEPTOR
OG0000059	+25*	PTHR11177	CHITINASE
OG0000062	+20*	PTHR28576	PIGGYBAC TRANSPOSABLE ELEMENT-DERIVED PROTEIN
OG0000075	+20*	PTHR19325	COMPLEMENT COMPONENT-RELATED SUSHI DOMAIN-CONTAINING
OG0000091	+27*	PTHR11119	XANTHINE-URACIL / VITAMIN C PERMEASE FAMILY MEMBER
OG0000118	+30*	PTHR10877	POLYCYSTIN-RELATED
OG0000124	+9*	PTHR13802	MUCIN 4-RELATED
OG0000128	+23*	PTHR10131	TNF RECEPTOR ASSOCIATED FACTOR
OG0000137	+14*	PTHR11908	XANTHINE DEHYDROGENASE
OG0000146	+18*	PTHR11709	MULTI-COPPER OXIDASE
OG0000150	+20*	PTHR14647	GALACTOSE-3-O-SULFOTRANSFERASE
OG0000155	+13*	PTHR23033	BETA1,3-GALACTOSYLTRANSFERASE
OG0000171	+17*	PTHR10283	SOLUTE CARRIER FAMILY 13 MEMBER
OG0000180	+16*	PTHR24039	FIBRILLIN
OG0000183	+21*	PTHR14453	PARP/ZINC FINGER CCCH TYPE DOMAIN CONTAINING PROTEIN
OG0000184	+35*	PTHR24221	ABC TRANSPORTER
OG0000198	+27*	PTHR24033	FAMILY NOT NAMED
OG0000219	+24*	PTHR10579	CALCIUM-ACTIVATED CHLORIDE CHANNEL REGULATOR
OG0000238	+26*	PTHR11039	NEBULIN
OG0000242	+11*	PTHR34415	FAMILY NOT NAMED
OG0000247	+15*	PTHR14453	PARP/ZINC FINGER CCCH TYPE DOMAIN CONTAINING PROTEIN
OG0000250	+29*	PTHR13800	TRANSIENT RECEPTOR POTENTIAL CATION CHANNEL, SUBFAMILY M, MEMBER 6
OG0000262	+16*	PTHR43645	UPF0214 PROTEIN YFEW
OG0000268	+13*	PTHR43998	FILAMIN
OG0000277	+13*	PTHR44131	CUB DOMAIN-CONTAINING PROTEIN
OG0000287	+18*	PTHR44097	PROTEIN SERRATE
OG0000290	+7*	PTHR10166	VOLTAGE-DEPENDENT CALCIUM CHANNEL SUBUNIT ALPHA-2/DELTA-RELATED
OG0000293	+26*	PTHR22605	AAA+ ATPASE, CORE DOMAIN-CONTAINING PROTEIN
OG0000304	+12*	PTHR23097	TUMOR NECROSIS FACTOR RECEPTOR SUPERFAMILY MEMBER
OG0000305	+17*	PTHR23232	KRAB DOMAIN C2H2 ZINC FINGER
OG0000307	+10*	PTHR18966	IONOTROPIC GLUTAMATE RECEPTOR
OG0000318	+11*	PTHR12622	DELTEX-RELATED
OG0000322	+17*	PTHR23024	MEMBER OF 'GDXG' FAMILY OF LIPOLYTIC ENZYMES
OG0000324	+8*	PTHR11106	GANGLIOSIDE INDUCED DIFFERENTIATION ASSOCIATED PROTEIN 2-RELATED

OG0000337	+24*	PTHR10887	DNA2/NAM7 HELICASE FAMILY
OG0000342	+25*	PTHR16897	CARNOSINE N-METHYLTRANSFERASE
OG0000347	+11*	PTHR10796	PATCHED-RELATED
OG0000357	+10*	PTHR23302	TRANSMEMBRANE CHANNEL-RELATED
OG0000365	+7*	PTHR12042	LACTOSYLCERAMIDE 4-ALPHA-GALACTOSYLTRANSFERASE ALPHA- 1,4-GALACTOSYLTRANSFERASE
OG0000373	+23*	PTHR31649	LD46221P-RELATED
OG0000378	+39*	PTHR13715	RYANODINE RECEPTOR AND IP3 RECEPTOR
OG0000388	+16*	PTHR15600	SACSin
OG0000398	+13*	PTHR11697	GENERAL TRANSCRIPTION FACTOR 2-RELATED ZINC FINGER PROTEIN
OG0000423	+8*	PTHR44014	FAMILY NOT NAMED
OG0000432	+14*	PTHR14454	GRB2-ASSOCIATED AND REGULATOR OF MAPK PROTEIN
OG0000460	+10*	PTHR43905	CONTACTIN
OG0000465	+11*	PTHR23130	FERRIC-CHELATE REDUCTASE
OG0000468	+15*	PTHR16897	CARNOSINE N-METHYLTRANSFERASE
OG0000477	+10*	PTHR11616	SODIUM/CHLORIDE DEPENDENT TRANSPORTER
OG0000491	+20*	PTHR11046	OLIGORIBONUCLEASE, MITOCHONDRIAL
OG0000500	+12*	PTHR24106	CASPASE RECRUITMENT DOMAIN-CONTAINING PROTEIN 8/NACHT, LRR AND PYD DOMAINS-CONTAINING PROTEIN
OG0000524	+50*	PTHR19277	PENTRAXIN
OG0000555	+9*	PTHR31009	S-ADENOSYL-L-METHIONINE:CARBOXYL METHYLTRANSFERASE FAMILY PROTEIN
OG0000562	+8*	PTHR15698	PHYTANOYL-COA HYDROXYLASE-INTERACTING PROTEIN
OG0000570	+27*	PTHR10697	MAMMALIAN EPENDYMIN-RELATED PROTEIN 1
OG0000578	+21*	PTHR33748	FAMILY NOT NAMED
OG0000586	+14*	PTHR34153	SI:CH211-262H13.3
OG0000592	+41*	PTHR19325	COMPLEMENT COMPONENT-RELATED SUSHI DOMAIN-CONTAINING
OG0000613	+11*	PTHR44854	FIBROCYSTIN-
OG0000617	+20*	PTHR19325	COMPLEMENT COMPONENT-RELATED SUSHI DOMAIN-CONTAINING
OG0000640	+28*	PTHR23145	NUCLEOSOMAL BINDING PROTEIN 1
OG0000656	+9*	PTHR19297	GLYCOSYLTRANSFERASE 14 FAMILY MEMBER
OG0000687	+13*	PTHR24243	G-PROTEIN COUPLED RECEPTOR
OG0000715	+10*	PTHR10796	PATCHED-RELATED
OG0000728	+22*	PTHR45240	ZINC METALLOPROTEINASE NAS
OG0000789	+14*	PTHR10773	DNA-DIRECTED RNA POLYMERASES I, II, AND III SUBUNIT RPABC2
OG0000849	+37*	PTHR12673	FACIOGENITAL DYSPLASIA PROTEIN
OG0000862	+23*	PTHR11514	MYC
OG0000917	+10*	PTHR31569	ZINC FINGER SWIM DOMAIN-CONTAINING PROTEIN
OG0000927	+21*	PTHR44252	CARBONYL REDUCTASE NADPH
OG0000938	+17*	PTHR11360	MONOCARBOXYLATE TRANSPORTER
OG0000943	+39*	PTHR13954	IRE1-RELATED
OG0000987	+10*	PTHR24280	CYTOCHROME P450 20A1

OG0001012	+19*	PTHR16897	CARNOSINE N-METHYLTRANSFERASE
OG0001096	+11*	PTHR11799	PARAOXONASE
OG0001166	+10*	PTHR16897	HISTAMINE N-METHYLTRANSFERASE
OG0001243	+12*	PTHR11903	PROSTAGLANDIN G/H SYNTHASE
OG0001293	+52*	PTHR23194	PYGOPUS
OG0001377	+10*	PTHR31513	GLYCINE-RICH PROTEIN
OG0001383	+9*	PTHR20956	HEH2P
OG0001422	+41*	PTHR23259	RIDDLE
OG0001533	+24*	PTHR37558	FAMILY NOT NAMED
OG0001594	+10*	PTHR37445	FAMILY NOT NAMED
OG0001671	+22*	PTHR10199	THROMBOSPONDIN
OG0001730	+7*	PTHR24106	NACHT, LRR AND PYD DOMAINS-CONTAINING PROTEIN
OG0001853	+6*	PTHR13627	FUKUTIN RELATED PROTEIN
OG0001862	+10*	PTHR11697	GENERAL TRANSCRIPTION FACTOR 2-RELATED ZINC FINGER PROTEIN
OG0001954	+10*	PTHR15031	CARTILAGE INTERMEDIATE LAYER PROTEIN CLIP
OG0002014	+17*	PTHR12419	OTU DOMAIN CONTAINING PROTEIN
OG0002286	+6*	NONE	NA
OG0002380	+6*	NONE	NA
OG0002381	+22*	PTHR11046	OLIGORIBONUCLEASE, MITOCHONDRIAL
OG0002455	+11*	PTHR35558	FAMILY NOT NAMED
OG0000077	-10*	PTHR14918	PROTEIN SZT2

Table S6.

Key genes of host genes identified as proteins from proteomic analysis.

Function System	Feature ID	Function
Proteosome	LLUY_028786-T1	PSMA6; 20S proteasome subunit alpha 1 [EC:3.4.25.1]
	LLUY_012964-T1	PSMA2; 20S proteasome subunit alpha 2 [EC:3.4.25.1]
	LLUY_008627-T1	PSMA4; 20S proteasome subunit alpha 3 [EC:3.4.25.1]
	LLUY_027290-T1	PSMA7; 20S proteasome subunit alpha 4 [EC:3.4.25.1]
	LLUY_003537-T1	PSMA5; 20S proteasome subunit alpha 5 [EC:3.4.25.1]
	LLUY_040090-T1	PSMA1; 20S proteasome subunit alpha 6 [EC:3.4.25.1]
	LLUY_012936-T1	PSMA3; 20S proteasome subunit alpha 7 [EC:3.4.25.1]
	LLUY_023805-T1	PSMB3; 20S proteasome subunit beta 1
	LLUY_029937-T1	PSMB3; 20S proteasome subunit beta 3
	LLUY_002855-T1	PSMB2; 20S proteasome subunit beta 4
	LLUY_015520-T1	PSMB5; 20S proteasome subunit beta 5
	LLUY_022140-T1	PSMB1; 20S proteasome subunit beta 6
	LLUY_022807-T1	PSMB4; 20S proteasome subunit beta 7
	LLUY_003098-T1	SMD2, RPN1; 26S proteasome regulatory subunit N1
	LLUY_026382-T1	PSMD7, RPN8; 26S proteasome regulatory subunit N8
	LLUY_017796-T1	EIF3E, INT6; translation initiation factor 3 subunit E
	LLUY_014232-T1	HSPA1s; heat shock 70kDa protein 1/2/6/8
	LLUY_037621-T1	HSP90A, htpG; molecular chaperone HtpG
	LLUY_000099-T1	HSP90A, htpG; molecular chaperone HtpG
	LLUY_000099-T1	RAD23, HR23; UV excision repair protein RAD23
Lysosome	LLUY_033571-T1	cathepsin C [EC:3.4.14.1]
	LLUY_007735-T1, LUY_007737-T1, LUY_036642-T1, LUY_036643-T1	cathepsin B [EC:3.4.22.1]
	LLUY_009810-T1, LUY_026908-T1	cathepsin L [EC:3.4.22.15]
	LLUY_012982-T1	legumain [EC:3.4.22.34]
	LLUY_007693-T1	cathepsin F [EC:3.4.22.41]
	LLUY_016715-T1	cathepsin D [EC:3.4.23.5]
	LLUY_023349-T2, LUY_023349-T1	clathrin heavy chain
	LLUY_017989-T1	lysosomal-associated membrane protein 1/2
	LLUY_024617-T1	cathepsin X [EC:3.4.18.1]
	LLUY_002319-T1	lysosomal alpha-mannosidase [EC:3.2.1.24]
	LLUY_004297-T1	lysosomal alpha-glucosidase [EC:3.2.1.20]
	LLUY_005359-T1	saposin
	LLUY_016480-T1	lysosome membrane protein 2
	LLUY_003890-T1	cathepsin A (carboxypeptidase C) [EC:3.4.16.5]
Longevity regulating pathway	LLUY_007262-T1, LUY_014836-T1, LUY_014856-T1	SOD2; superoxide dismutase, Fe-Mn family [EC:1.15.1.1]
	LLUY_018867-T1	SOD1; superoxide dismutase, Cu-Zn family [EC:1.15.1.1]

Table S7.

Key genes of symbiont genes identified as proteins from proteomic analysis.

Function System	Feature ID	Length (bp)	Function
rTCA Cycle	Lamellibrachia_symbiont.peg.150	1524	Fumarate hydratase class I, aerobic (EC 4.2.1.2)
	Lamellibrachia_symbiont.peg.2185	1911	2-oxoglutarate oxidoreductase, alpha subunit (EC 1.2.7.3)
	Lamellibrachia_symbiont.peg.2186	954	2-oxoglutarate oxidoreductase, beta subunit (EC 1.2.7.3)
	Lamellibrachia_symbiont.peg.2943	987	Malate dehydrogenase (EC 1.1.1.37)
	Lamellibrachia_symbiont.peg.986	1161	ATP citrate lyase beta chain (EC 4.3.1.8)
	Lamellibrachia_symbiont.peg.987	870	ATP citrate lyase alpha chain (EC 4.3.1.8)
	Lamellibrachia_symbiont.peg.2924	1167	2-oxoglutarate oxidoreductase, alpha subunit (EC 1.2.7.3)
	Lamellibrachia_symbiont.peg.2925	999	2-oxoglutarate oxidoreductase, beta subunit (EC 1.2.7.3)
	Lamellibrachia_symbiont.peg.2926	570	2-oxoglutarate oxidoreductase, gamma subunit (EC 1.2.7.3)
Calvin Cycle	Lamellibrachia_symbiont.peg.2754	1389	Ribulose bisphosphate carboxylase (EC 4.1.1.39)
	Lamellibrachia_symbiont.peg.2757	804	Rubisco activation protein CbbQ
	Lamellibrachia_symbiont.peg.2758	2256	Rubisco activation protein CbbO
Sulfer Oxidataion	Lamellibrachia_symbiont.peg.303	942	Dissimilatory sulfite reductase, beta subunit (EC 1.8.99.3)
Nitrogen Mitobolism	Lamellibrachia_symbiont.peg.724	1602	Respiratory nitrate reductase beta chain (EC 1.7.99.4) Denitrifying reductase gene clusters; Nitrate and nitrite ammonification
	Lamellibrachia_symbiont.peg.725	3762	Respiratory nitrate reductase alpha chain (EC 1.7.99.4)
Adhesion-related proteins	Lamellibrachia_symbiont.peg.2820	978	Ankyrin
	Lamellibrachia_symbiont.peg.2856	3813	Fibronectin type III domain protein
Oxidative stress	Lamellibrachia_symbiont.peg.1649	582	Superoxide dismutase [Fe] (EC 1.15.1.1) (FeSOD)
	Lamellibrachia_symbiont.peg.2936	489	Rubrerythrin

Table S8.

Lamellibrachia luymesi Hb sequences identified that are highly expressed in the trophosome tissue or as proteins from proteomic data.

HBs	<i>Lamellibrachia luymsei</i> Hbs	Differential expressed	Mass spectrum
A1 Chain	LLUY_034331-T1	*	*
	LLUY_034332-T1	*	*
A2 Chain	LLUY_034333-T1	*	*
B2 Chain	LLUY_029258-T1	X	*
B1 Chain	LLUY_004752-T1	*	*
	LLUY_004753-T1	*	X
	LLUY_005026-T1	X	X
	LLUY_005027-T1	X	X
	LLUY_005030-T1	*	X
	LLUY_005031-T1	X	X
	LLUY_009666-T1	X	X
	LLUY_009670-T1	X	X
	LLUY_009671-T1	*	X
	LLUY_009673-T1	X	X
	LLUY_013447-T1	*	*
	LLUY_013449-T1	*	X
	LLUY_017246-T1	*	*
	LLUY_017247-T1	*	*
	LLUY_020152-T1	*	X
	LLUY_026555-T1	X	X
	LLUY_026556-T1	*	*
	LLUY_029945-T1	X	X
	LLUY_032994-T1	X	X
	LLUY_038743-T1	X	X
	LLUY_039230-T1	X	X
	LLUY_040024-T1	X	X
	LLUY_004634-T1	*	X
Linkers	LLUY_002344-T1	*	*
	LLUY_024479-T1	*	X
	LLUY_026441-T1	*	*
	LLUY_026443-T1	*	*

*: Positive; X: negative.

Sequences with putative free-cysteine were colored as red.

Table S9.

Number of unique TLR proteins encoded in lophotrochozoan genomes.

Taxon	Species	Number of TLR identified	Number of RLR identified
Annelida	<i>Lamellibrachia luymesi</i>	33	2
	<i>Capitella telata</i>	5	3
	<i>Helobdella robusta</i>	4	3
Mollusca	<i>Bathymodiolus platifrons</i>	61	7
	<i>Crassostrea gigas</i>	61	11
	<i>Modiolus philippinarum</i>	90	2
	<i>Mizuhopecten yessoensis</i>	30	4
	<i>Octopus bimaculoides</i>	5	5
	<i>Patinopecten yessoensis</i>	22	4
	<i>Lottia gigantea</i>	7	3
	<i>Crassostrea virginica</i>	109	8
	<i>Notospermus geniculatus</i>	5	4
Nemertea	<i>Phoronis australis</i>	23	3
Brachiopoda	<i>Lingula anatina</i>	46	13
Flatworm	<i>Biomphalaria glabrata</i>	17	6
Rotifera	<i>Aplysia californica</i>	15	2
Vertebrata	<i>Homo sapiens</i>	11	3

Table S10.

Domain requirements for identifying components of TLR pathway.

Protein	Domain Requirements
<i>TLR/TOLL</i>	TIR+LRR(≥ 3)
<i>MYD88</i>	TIR+DEATH
<i>SARM1</i>	TIR+SAM(2)
<i>DDX58</i>	CARD(2)+Helicase_ATP_binding+Helicase_C_Terminal+Rig1_Regulatory_Domain
<i>DHX58</i>	Helicase_ATP_binding+Helicase_C_Terminal+Rig1_Regulatory_Domain
<i>IFIH1</i>	CARD(2)+Helicase_ATP_binding+Helicase_C_Terminal+Rig1_Regulatory_Domain
<i>TRAF</i>	Zn_Finger+MATH/TRAF or Zn_Finger+WD40_repeats(≥ 3)
<i>NLR</i>	(x)+NACHT+LRR(≥ 3)
<i>IKK</i>	Kinase
<i>IKB</i>	ANK(≥ 3)
<i>NFKB</i>	RHD+ANK(≥ 3)+DEATH
<i>NEMO</i>	NEMO/Coiled_coil
<i>IRF</i>	Interferon_regulatory_factor_DNA_binding_domain