1    **Genomic adaptations to chemosymbiosis in the deep-sea seep-dwelling**

2    **tubeworm *Lamellibrachia luymesi* (Siboglinidae, Annelida)**

3    Yuanning Li[1,2]*, Michael G. Tassia[1], Damien S. Waits[1], Viktoria E. Bogantes[1], Kyle T.

4    David[1], 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

## Abstract

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

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

**Significance**

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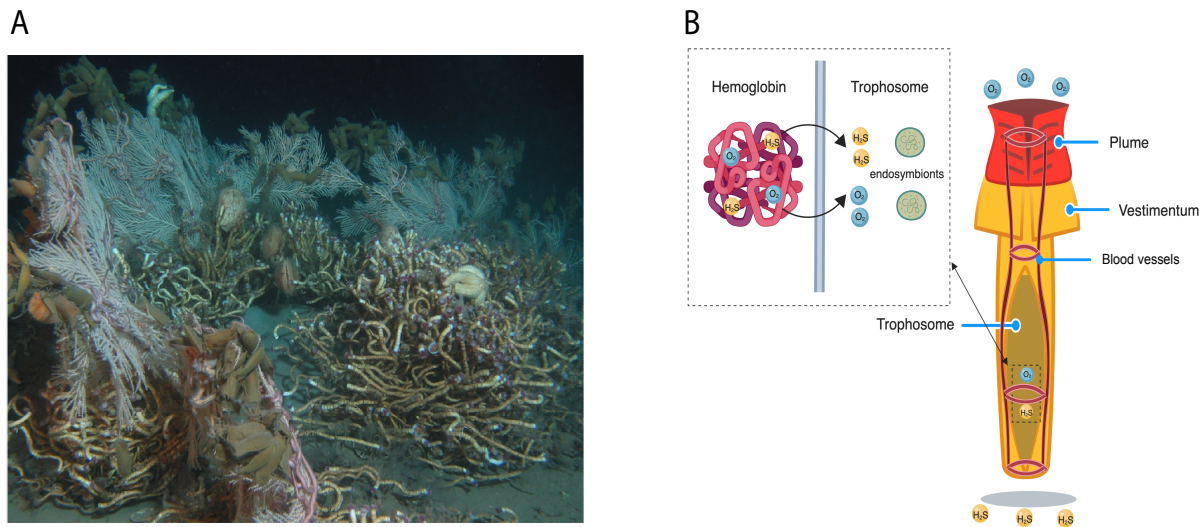
**Introduction**



**Fig. 1.** *Lamellibrachia luymesi*. (A). Seep habitat in the Gulf of Mexico. (B). Diagram of adult *L. luymesi* worm to model $O_2$ and $H_2S$ transport to symbionts in trophosome by hemoglobin molecules. The hemoglobin model was created with the help of Biorender (https://biorender.com/).

Recent advances in understanding the dominance of microbes on the planet has placed new emphasis on elucidating mechanisms that promote microbe-animal symbioses. Although considerable work has been undertaken on adaptations of microbial genomes to facilitate animal symbiosis (such as corals, termites, humans), examples of how animal host genomes have adapted to symbioses are limited (1). Vestimentiferan tubeworms inhabit some of Earth's most extreme environments, such as deep-sea hydrothermal vents and cold seeps, and are obligate dependents on symbiosis for survival. These animals lack a digestive tract and rely on sulfide-oxidizing bacterial symbionts for nutrition and growth. At some seeps, tubeworms, such as *Lamellibrachia luymesi* in the Gulf of Mexico, are so abundant that they transform the habitat (Fig. 1A) and thus facilitate biodiversity promoting adaptive radiations and evolutionary novelties (2). Given the obligate nature of the symbiosis between tubeworms and their gammaproteobacterial chemoautotrophic endosymbiont, one may 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. $H_2S$, $O_2$, $H_2$) 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
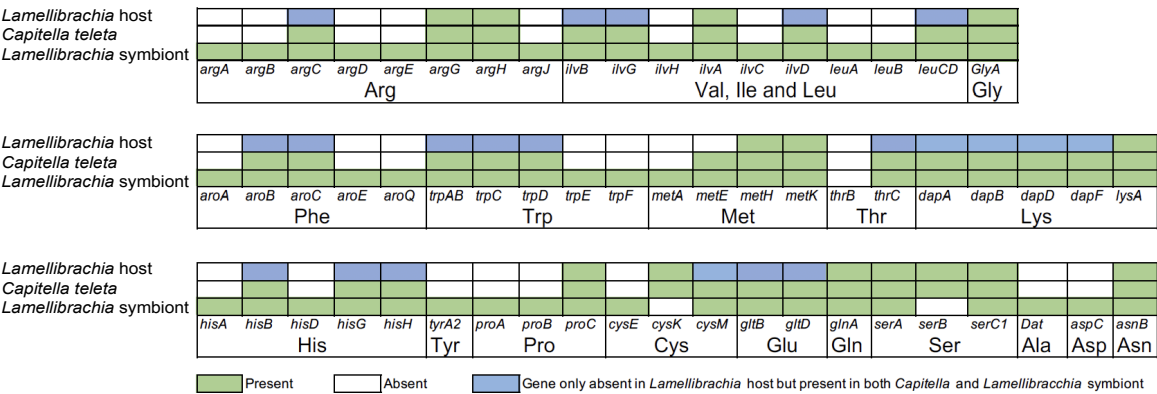
125    is lacking.

126         Using comparative genomics, transcriptomic and proteomic analyses on the

127    tubeworm *Lamellibrachia luymesi*, 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.

## Results and Discussion

### Genome features

133         Using Illumina paired-end, mate-pair and 10X genomic sequencing (Table S1),

134    we the assembled genome of a single *Lamellibrachia luymesi* 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    *luymesi* 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. luymesi* 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**
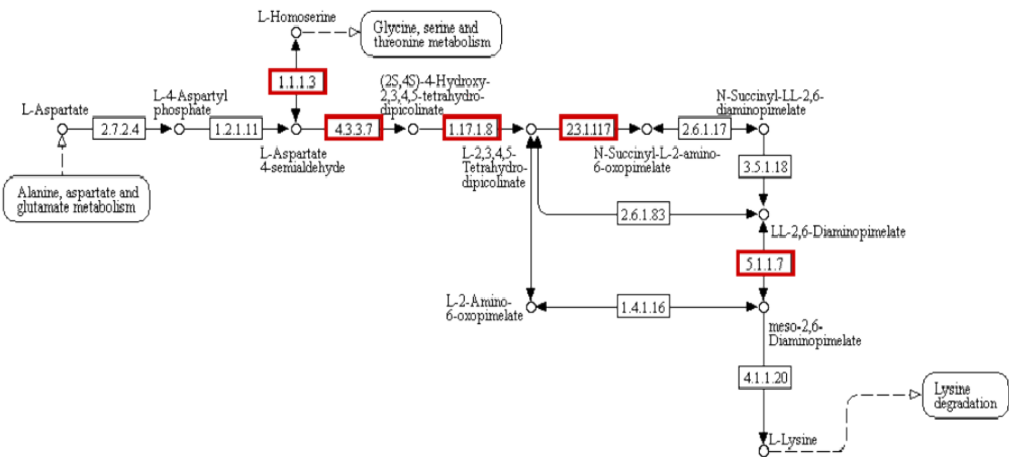


146

147    **Fig. 2.** *Lamellibrachia luymesi* 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. luymesi* and *L. luymesi* symbionts. Blue boxes

150    represents genes present in *C. teleta* and *L. luymesi* gammaproteobacterial symbionts

151    but absent in *L. luymesi*. (B) Example of Lysine biosynthesis pathway. Red boxes

152    indicate genes missing in *L. luymesi*. 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    *luymesi* 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. luymesi* 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. luymesi* 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. luymesi* 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 magnifia* (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. luymesi* (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

**Hemoglobin evolution**



**Fig. 3.** Hemoglobin gene diversity in *Lamellibrachia luymesi*. Gene tree of siboglinid Hb subunits A1, A2, B1 and B2 reconstructed using IQtree with 1000 ultrafast bootstrap. Only Siboglinid Hb sequences (from SwissProt database or this study) were labeled. *L. luymesi* sequences labeled red. Accession numbers associated with each sequence was shown in the full tree (Fig. S3).

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

10

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 $H_2S$ detoxification process (21).

236    Surprisingly, we found a significant expansion of B1 Hbs, 25 copies, in *L. luymesi*

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. luymesi* 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 $H_2S$ binding, another hypothesis suggested

244    that zinc moieties bound to amino acid residues at the interface between pairs of A2

245    chains influence $H_2S$ binding (9). The $Zn^{2+}$-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 $H_2S$ transport.
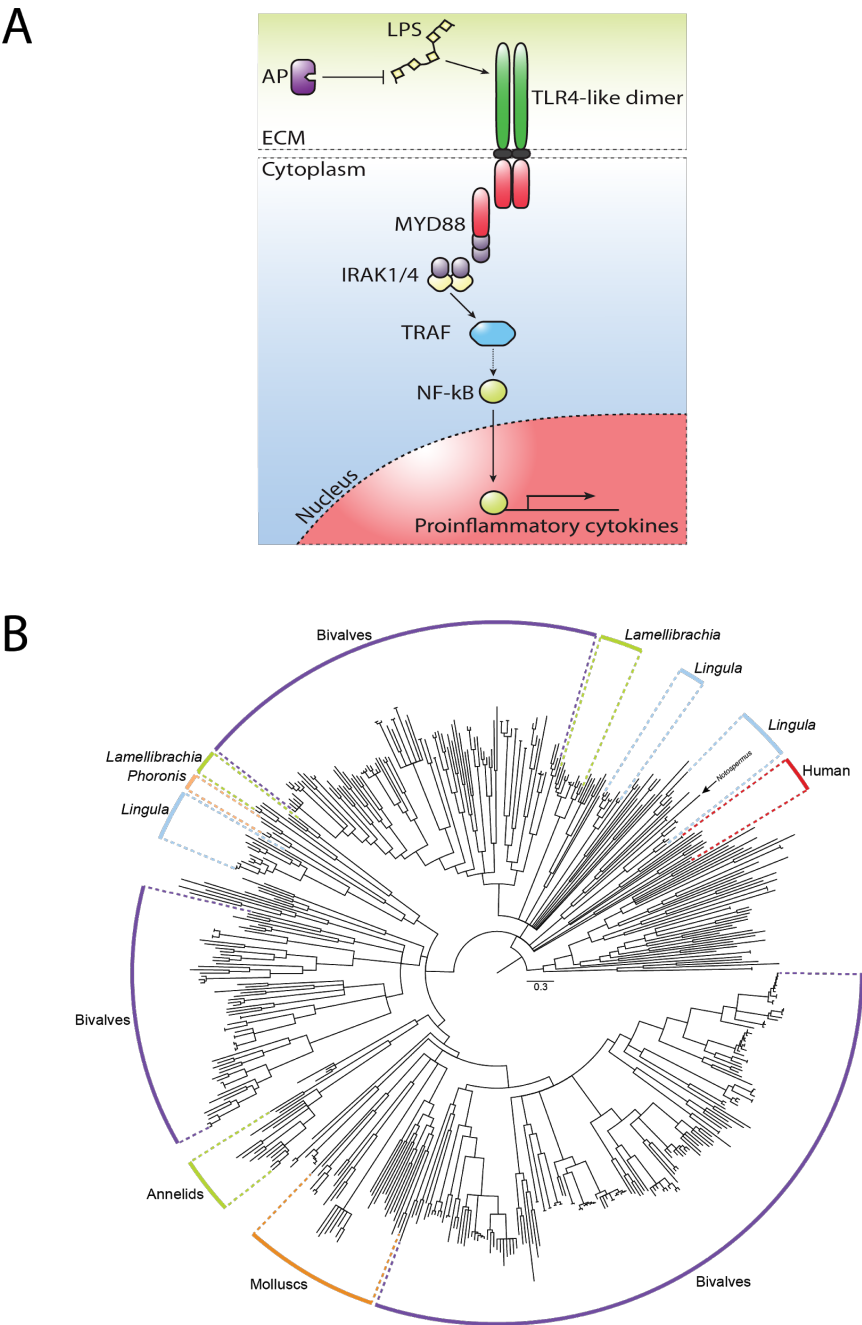
## Immunity function



**Fig. 4.** Toll-Like Receptors (TLRs) in *Lamellibrachia luymesi*. (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 ≥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 form *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-kB, 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-kB, 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-kB 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    Seep-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. luymesi*'s genome compared to other

311    lophotrochozoans (Fig. S6). Most lophotrochozoan genomes contain one or two copies

312    of SOD1 and SOD2, but *L. luymesi* 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 (≥ 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.


**Acknowledgments**

**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

**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 Seep 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. Vurture 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. Pryszcz 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.*

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