

1 **The chromosomal-level genome assembly and comprehensive**
2 **transcriptomes of Chinese razor clam (*Sinonovacula constricta*) with**
3 **deep-burrowing life style and broad-range salinity adaptation**

4
5 Yinghui Dong^{1,†}, Qifan Zeng^{2,†}, Jianfeng Ren^{3,†}, Hanhan Yao¹, Wenbin Ruan¹, Liyuan
6 Lv¹, Lin He¹, Qinggang Xue¹, Zhenmin Bao^{2,4}, Shi Wang^{2,5*}, Zihua Lin^{1,*}

7
8
9 ¹Key Laboratory of Aquatic Germplasm Resource of Zhejiang, College of Biological
10 & Environmental Sciences, Zhejiang Wanli University, Ningbo 315100, China.

11
12 ²MOE Key Laboratory of Marine Genetics and Breeding, College of Marine Life
13 Sciences, Ocean University of China, Qingdao 266003, China

14
15 ³Key Laboratory of Exploration and Utilization of Aquatic Genetic Resources,
16 Ministry of Education, College of Fisheries and Life Science, Shanghai Ocean
17 University, Shanghai, 201306, China.

18
19 ⁴Laboratory for Marine Fisheries Science and Food Production Processes, Pilot
20 National Laboratory for Marine Science and Technology, Qingdao 266237, China

21
22 ⁵Laboratory for Marine Biology and Biotechnology, Pilot National Laboratory for
23 Marine Science and Technology, Qingdao 266237, China

24
25 [†]These authors contributed equally to this work.

26 Yinghui Dong, Email: dongyinghui118@126.com

27 Qifan Zeng, Email: zengqifan@ouc.edu.cn

28 Jianfeng Ren, Email: jfren@shou.edu.cn

29 Hanhan Yao, Email: yaohanhan1020@126.com

30 Wenbin Ruan, Email: wbruan@163.com

31 Liyuan Lv, Email: llyuan.2009@163.com

32 Lin He, Email: hlwithyou@qq.com

33 Qinggang Xue, Email: qxue@zwu.edu.cn

34 Zhenmin Bao, zmbao@ouc.edu.cn

35 *To whom correspondence should be addressed: Zhihua Lin: zhihua9988@126.com;

36 Shi Wang: swang@ouc.edu.cn.

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60 **Abstract**

61 **Background:** The Chinese razor clam, *Sinonovacula constricta*, is one of the
62 commercially important marine bivalves with deep-burrowing lifestyle and
63 remarkable adaptability of broad-range salinity. Despite its economic impact and
64 representative of the less-understood deep-burrowing bivalve lifestyle, there are few
65 genomic resources for exploring its unique biology and adaptive evolution. Herein,
66 we reported a high-quality chromosomal-level reference genome of *S. constricta*, the
67 first genome of the family Solenidae, along with a large amount of
68 short-read/full-length transcriptomic data of whole-ontogeny developmental stages,
69 all major adult tissues, and gill tissues under salinity challenge .

70 **Findings:** A total of 101.79 Gb and 129.73 Gb sequencing data were obtained with
71 the PacBio and Illumina platforms, which represented approximately 186.63X
72 genome coverage. In addition, a total of 160.90 Gb and 24.55 Gb clean data were also
73 obtained with the Illumina and PacBio platforms for transcriptomic investigation. A
74 *de novo* genome assembly of 1,340.13 Mb was generated, with a contig N50 of
75 689.18 kb. Hi-C scaffolding resulted in 19 chromosomes with a scaffold N50 of 57.99
76 Mb. The repeat sequences account for 50.71% of the assembled genome. A total of
77 26,273 protein-coding genes were predicted and 99.5% of them were annotated.
78 Phylogenetic analysis revealed that *S. constricta* diverged from the lineage of
79 Pteriomorphia at approximately 494 million years ago. Notably, cytoskeletal protein
80 tubulin and motor protein dynein gene families are rapidly expanded in the *S.*
81 *constricta* genome and are highly expressed in the mantle and gill, implicating
82 potential genomic bases for the well-developed ciliary system in the *S. constricta*.

83 **Conclusions:** The high-quality genome assembly and comprehensive transcriptomes
84 generated in this work not only provides highly valuable genomic resources for future
85 studies of *S. constricta*, but also lays a solid foundation for further investigation into
86 the adaptive mechanisms of benthic burrowing mollusks.

87 **Keywords:** *Sinonovacula constricta*, genome, transcriptome, gene family, ciliary
88 system

89

90 **Introduction**

91 The Chinese razor clam *Sinonovacula constricta* (Lamarck 1818) is a member of
92 the Solenidae family of bivalve molluscs, recognizing for its straight razor-like shape
93 and fragile shells (Figure 1A). It is widely distributed in the intertidal zone along the
94 west Pacific Ocean and engages in a pelago-benthic life cycle (Figure 1B). As
95 adaptation to a deep-burrowing lifestyle, the razor clam is characterized by smooth
96 shells, muscular foot, and elongated siphons (Figure 1). Benefit from its relatively
97 short production cycle and high productive efficiency, the razor clam has become one
98 of the four most important maricultured bivalve species (together with oyster, scallop,
99 and *Venerupis* spp.) in China, Japan, and Korea, with over 800,000 metric tons of
100 production in 2016 (FAO, 2018).

101 As living in estuarine and intertidal region, the razor clam faces tremendous
102 exposure to extreme environmental stresses such as drastic salinity fluctuation, highly
103 variable temperature, high concentration of ammonia nitrogen and hydrogen sulfide.
104 Unlike oysters, mussels and most clams with thick and sealed shells for protecting
105 their soft bodies, the razor clam with two thin and unclosed shells has chosen a
106 survival strategy of deep-burrowing lifestyle with high tolerance of a broad range of
107 salinities (5-45‰), making it an ideal model to investigate the adaptive mechanisms
108 of deep-burrowing lifestyle. Despite its economic impact and representative of the
109 less-understood deep-burrowing bivalve lifestyle, there are few genomic resources for
110 exploring its unique biology and adaptive evolution. Here, we generated the
111 high-quality chromosomal-level genome assembly and comprehensive transcriptomes
112 of *S. constricta* and investigated the transcriptomic changes under different
113 environmental stresses. These genomic resources will lay a prime foundation for
114 future studies of its lifestyle-related adaptive evolution and genetic improvement in
115 commercial breeding programs.

116

117 **Genomic DNA preparation, PacBio and Illumina sequencing**

118 An individual *S. constricta* was collected from the brood stock at the genetic
119 breeding research center of Zhejiang Wanli University. Genomic DNA was extracted

120 from muscle tissues using a phenol-chloroform method as described in the protocol
121 (Green and Sambrook, 2012). High molecular weight genomic DNA was sheared into
122 fragments of ~30 kb using a Covaris ultrasonicators (Covaris, Woburn, MA, USA).
123 The fragments were enzymatically repaired and converted into SMRTbell™ template
124 library following the manufacturer's instructions. Size-selection was performed to
125 enrich the DNA fragments longer than 10 kb for sequencing on a Pacific Biosciences
126 (PacBio) Sequel Single Molecule Real Time (SMRT) platform. The genomic library
127 was sequenced in 6 cells, generating 10,549,576 subreads with a N50 length of 13,619
128 bp, and accounting for a total of 101.79 Gb. A paired-end Illumina library with an
129 insert size of 300 bp were prepared with an Illumina Genomic DNA sample
130 Preparation kit and sequenced on an Illumina Xten system, yielding a total of 129.73
131 Gb reads with an insert size of 350 bp ([Supplementary Table S1](#)).

132

133 **Estimation of the genome size and sequencing coverage**

134 The Illumina short reads were first trimmed to remove adaptors and reads with
135 more than 10% ambiguous or more than 20% low-quality bases using Trimmomatic
136 (Bolger et al., 2014). The distribution of 17-mer frequency was estimated using the
137 clean reads. A total of 10^{10} k-mer was identified with the peak depth of coverage
138 being 80. Based on the formula: genome size = k-mer number/peak depth (Varshney
139 et al., 2011), the genome size of *S. constricta* was estimated to be 1,244.27 Mb, with a
140 heterozygous ratio of 1.53% and repeat rate of 53.12% ([Supplementary Figure S1](#)).

141

142 **De novo genome assembly and quality assessment**

143 PacBio long reads were corrected and assembled using the Falcon package (Chin et
144 al., 2016). Briefly, all the raw reads yielded by Pacbio platform were aligned to each
145 other to identify overlaps with DALIGNER. The overlap data and raw subreads were
146 then processed for consensus calling. After the error-correction, overlaps were
147 detected in the preassembled error-corrected reads and used to construct a directed
148 fragment assembly string graph. Contigs were constructed by finding the paths from
149 the string graph. The consensus calling of preceding assembly was performed with

150 Quiver. Subsequently, the paired-end clean reads yielded by Illumina platform were
151 aligned to polish the assembly using Pilon (Walker et al., 2014). The resulting
152 assembly consisted of 10,981 contigs, comprising 1,331.97 Mb with a contig N50 of
153 678,857 bp and GC contents of 35.46% (Table 1).

154 To assess the integrity of the genome assembly, Illumina short-insert library reads
155 were mapped to the contigs using BWA (version 0.6.2). In summary, 88.90% of the
156 assembled genome sequences were covered by 93.93% of the total reads
157 (Supplementary Table S2). The genome completeness was also evaluated using both
158 Core Eukaryotic Genes Mapping Approach (CEGMA) analysis (Parra et al., 2007)
159 and Benchmarking Universal Single-Copy Orthologs (BUSCO version 3) analysis
160 (Waterhouse et al., 2017). The CEGMA analysis identified 227 of the 248 core
161 eukaryotic genes (91.53%), and the BUSCO analysis unveiled 868 of the 978
162 near-universal single-copy metazoan orthologs (88.7%), indicating a high integrity of
163 the genome assembly (Supplementary Table S3 and S4).

164

165 **Illumina transcriptome sequencing and analysis**

166 Transcriptomic samples from different developmental stages and different adult
167 tissues were collected and sequenced for genome annotation. Embryos at four
168 developmental stages (eggs, four cells, blastulae, gastrulae), and larvae at four
169 developmental stages (trochophore larvae, D-shaped larvae, umbo larvae, and juvenile)
170 were collected at the hatchery of genetic breeding research center of Zhejiang Wanli
171 University. Artificial fertilization and larval culture were performed as previously
172 described (Dong et al., 2012). For each developmental stage, over 1,000 individuals
173 were collected for RNA extraction. Eight tissues (Figure 1A), including gill, digestive
174 gland, foot, mantle, adductor muscle, siphon, gonad (testis and ovary) were dissected
175 from one to three adult individuals and stored at -80°C after flash frozen in liquid
176 nitrogen.

177 Transcriptomic samples under salinities of 3‰, 25‰, and 38‰ were collected and
178 sequenced to identify genes and pathways involved in salt tolerance. The *S. constricta*
179 were subjected to salt stress for 16 hours under 22°C at extreme concentration of

180 low-salinity (3‰) and high-salinity (38‰) with the control concentration of
181 normal-salinity (25‰). Three replicate tanks for each group were set and each
182 replicate included 10 individuals. For the low-salinity group, the salinity was decreased
183 3‰ per hour though pouring into fresh water to target salinity of 3‰, and then
184 maintained for 16 hours. For the high-salinity groups, the salinity was raised 2‰ per
185 hour though pouring into artificial sea water to target salinity of 38‰, and then
186 maintained for 16 hours. Gills were dissected from three individuals of each replicate
187 and stored at -80°C.

188 Total mRNA was extracted from all the collected samples with TRIzol reagent
189 (OMEGA, America) following the manufacturer's instructions. A paired-end Illumina
190 library was constructed for each sample with an insert size of 300 bp and sequenced
191 on an Illumina X Ten system. Around 5-7 Gb of paired-end raw reads were yielded for
192 each sample. Clean reads were obtained by removing reads containing adapter, reads
193 containing ploy-N and low-quality reads by Trimmomatic (Bolger et al., 2014)
194 ([Supplementary Table S5](#)). The clean reads were aligned to the indexed *S. constricta*
195 reference genome using Hisat2 version 2.0.5 (Kim et al., 2015). The clean reads in the
196 samples of different adult tissues and salt stress were mapped onto the reference
197 genome with high proportion of around 70-80%, while samples of different
198 development stages with relative low proportion because of mixed thousand
199 individuals increasing the high SNP heterozygosity ([Supplementary Table S6](#)). The
200 featureCounts version 1.5.0 (Liao et al., 2014) was used to count the reads numbers
201 mapped to each gene and the gene expression level was calculated as FPKM
202 (Fragments Per Kilobase of transcript sequence per Millions base pairs) for each gene.

203 The identification of differentially expressed genes (DEGs) between different
204 salinity groups was performed using the DESeq2 R package version 1.16.1 with an
205 adjusted P-value <0.05 (Love et al., 2014). The numbers of up- and down-regulated
206 DEGs were 462 and 655 between the high-salinity group versus the normal-salinity
207 group, respectively while the numbers of up- and down-regulated DEGs were 898 and
208 826 between the low-salinity group versus the normal-salinity group, respectively
209 ([Supplementary Figure S2](#)). Gene Ontology (GO) enrichment analysis of DEGs was

210 implemented by the clusterProfiler R package with corrected P-value <0.05
211 considered significantly enriched GO terms (Yu et al., 2012). The clusterProfiler R
212 package is also used to test the statistical enrichment of DEGs in KEGG pathways
213 (Yu et al., 2012). The GO enrichment results demonstrated that the DEGs were
214 significantly enriched in the biological processes of transmembrane transport
215 (GO:0055085) and aminoglycan metabolic process (GO:0006022) and the molecular
216 functions of transmembrane transporter activity (GO:0022857) and ([Supplementary](#)
217 [Figure S3](#)). The KEGG pathway enrichment results indicated that the DEGs were
218 significantly enriched in amino acid metabolic pathways such as glycine, serine and
219 threonine metabolism, and arginine and proline metabolism, and the energy metabolic
220 pathways such as glycolysis/gluconeogenesis and citrate cycle ([Supplementary Figure](#)
221 [S4](#)).

222

223 **Full-length transcriptome sequencing and analysis**

224 Full-length RNA sequencing was also performed using the mixed RNAs from the
225 samples of eight development stages and eight adult tissues. Three libraries with
226 different insert lengths, e.g. 1-2k, 2-3k, and 3-6k, were constructed and sequenced on
227 a PacBio Sequel platform. A total of 1,064,194 post-filter polymerase reads were
228 obtained from 7 SMRT cells, including 688,944 full-length non-chimeric reads
229 ([Supplementary Table S7](#)). The full-length RNA transcriptomic analysis was
230 performed with the SMRT Link v4.0.0 software suite
231 (<https://www.pacb.com/support/software-downloads>). After redundant sequences
232 clustering using the ICE (Iterative Clustering and Error correction) algorithm,
233 consensus sequences building using the pbdagcon tool with DAGCon (Directed
234 Acyclic Graph Consensus) algorithm, and consensus sequence polishing with Quiver ,
235 a total of 61,620 high-quality (>0.99) and 358,297 low-quality (<0.99) transcript
236 sequences were obtained. Then, the transcript sequences were polished and corrected
237 using Illumina reads with LoRDEC (Salmela and Rivals, 2014). Finally, the corrected
238 transcripts were further clustered with CD-HIT (version 4.6) (Li and Godzik, 2006),
239 resulting in 75,225 Unigenes and 276,484 transcript isoforms. The full-length

240 transcripts were further used to annotate the protein-coding genes in the genome as
241 the direct evidences. The statistical information for full-length transcriptome analysis
242 is listed in [Supplementary Table S7](#).

243

244 **Repetitive sequence annotation**

245 Repetitive sequences in the genome assembly were identified through *ab initio*
246 prediction and homology-based searches. RepeatScout (version 1.0.5) and Repeat
247 Modeler version 1.0.11 (<http://www.repeatmasker.org>) were used for *de novo*
248 identification of repeat families in the *S. constricta* genome. Full length long terminal
249 repeat (LTR) retrotransposons were also identified using the LTR-finder (version 1.0.2)
250 (Xu and Wang, 2007) with the parameters “-E -C”. Tandem Repeats Finder (TRF
251 version 4.09) (Benson, 1999) was used to screen tandem repeats with the parameters
252 “match=2, mismatching penalty=7, indel penalty=7, match probability=80, indel
253 probability=10, minimum alignment score=50, maximum period size=500”. The
254 predicted repetitive sequences along with the RepBase database (Bao et al., 2015)
255 were used for homology-based searches using Repeatmasker (version 4.5.0) with the
256 parameters “-a -nolow -no_is -norna -parallel 32 -small -xsmall -poly -e ncbi -pvalue
257 0.0001” (Tarailo-Graovac and Chen, 2009).

258 Finally, a total of 675,404,889 bp repetitive sequences were identified, accounting
259 for 50.71% of the assembled genome ([Table 2](#)), which is consistent with our genome
260 survey result of 53.12%. Repetitive sequences were dominated by tandem repeats
261 (15.39%) and followed by DNA transposons (14.38%) and LTR retrotransposons
262 (10.84%) ([Table 2](#)).

263

264 **Protein-coding gene prediction and annotation**

265 Gene annotation was performed based on *de novo* prediction, homology-based
266 searches, and transcriptome assisted methods. Protein sequences of Yesso scallop
267 (*Patinopecte yessoensis*), Pacific oyster (*Crassostrea gigas*), owl limpet (*Lottia*
268 *gigantea*), octopus (*Octopus bimaculoides*), leech (*Helobdella robusta*), nematode
269 (*Caenorhabditis elegans*), fruit fly (*Drosophila melanogaster*), sea urchin

270 (*Strongylocentrotus purpuratus*), ascidian (*Ciona intestinalis*), Florida lancelet
271 (*Branchiostoma floridae*), and human (*Homo sapiens*) were downloaded from NCBI
272 and aligned to the genome assembly using TBLASTN with the parameters “-evalue
273 1e-5”. The gene structures were predicted with GeneWise (version 2.4.1) (Doerks et
274 al., 2002). The Illumina RNA-seq reads of the eight tissues and eight developmental
275 stages were aligned to the genome assembly using Tophat (version 2.1.1) (Trapnell et
276 al., 2009). Cufflinks (version 2.1.0) (Trapnell et al., 2010) was used to generate gene
277 models with the parameter “-multi-read-correct”. The resulting GTF file along with
278 the PacBio Iso-seq transcripts was utilized to model gene structures with the PASA
279 pipeline (version 2.0.2) (Haas et al., 2008).

280 Five *de novo* gene prediction packages, including Augustus (version 2.5.5) (Stanke
281 et al., 2006), glimmerHMM (version 3.01) (Majoros et al., 2004), SNAP (version
282 2006-07-28) (Leskovec and Susic, 2016), Geneid (version 1.4) (Parra et al., 2000),
283 and Genscan (version 3.1) (Burge and Karlin, 1997) were used to predict genes with
284 the repeat-masked genome sequences by default settings. All the gene model
285 evidences were integrated using EVidenceModeler (version 1.1.1) (Haas et al., 2008).
286 Finally, 26,273 protein-coding genes were identified in the *S. constricta* genome
287 ([Supplementary Table S8](#)).

288 The functional annotations were performed by aligning the predicted protein
289 sequences to public databases including KEGG, SwissProt and NCBI-NR databases
290 using BLASTP with the E-value threshold of 1e-5. InterProScan (v.4.8) (Jones et al.,
291 2014) was also used to identify motifs and domains by searching the Pfam, InterPro
292 and Gene Ontology (GO) databases. Taken together, 26,140 (99.5%) of the 26,273
293 genes could be annotated by at least one database ([Table 3](#)).

294

295 **Noncoding RNA prediction and annotation**

296 The noncoding RNA genes, including miRNAs, rRNAs, snRNAs, and tRNAs, were
297 annotated in the *S. constricta* genome. tRNAs were predicted by tRNAscan-SE 2.0
298 (Lowe and Chan, 2016) with eukaryote parameters. The miRNAs and snRNAs were
299 screened using INFERNAL 1.1.2 against the Rfam database (version 14.1) (Kalvari et

300 al., 2018) with default parameters. Finally, 968 miRNAs, 3,354 tRNAs, 822 rRNAs,
301 and 298 snRNAs were identified ([Supplementary Table S9](#)).

302

303 **Gene family and phylogenetic analysis**

304 Fifteen Eumetazoa species were selected for gene family analysis, including *H.*
305 *sapiens*, *B. floridae*, *D. melanogaster*, European honey bee (*Apis mellifera*),
306 Californian leech (*Helobdella robusta*), ocean-dwelling worm (*Capitella teleta*), *O.*
307 *bimaculoides*, *L. gigantea*, California sea hare (*Aplysia californica*), *C. gigas*,
308 American oyster (*Crassostrea virginica*), Sydney rock oyster (*Saccostrea glomerata*),
309 *P. yessoensis*, Chinese scallop (*Chlamys farreri*), and Starlet sea anemone
310 (*Nematostella vectensis*). All data were downloaded from either NCBI or Ensembl.
311 The longest protein sequence was selected to represent the gene with multiple
312 alternative splicing isoforms. Gene family clusters from all the 16 species were first
313 assigned using OrthoMCL (version 2.0.9) (Li et al., 2003) with an inflation value of
314 1.5. CAFE (version 3) (De Bie et al., 2006) was used to analyze gene family
315 expansion and contraction under maximum likelihood framework. The protein-coding
316 genes from all the 16 species were assigned into 39,058 families with 337 strict
317 single-copy orthologs. In the *S. constricta* genome, a total of 12,945 gene families
318 were identified, 803 of which were specifically possessed by *S. constricta*. Compared
319 with the other 15 species, *S. constricta* has 193 expanded and 31 contracted gene
320 families ([Figure 2](#)). Notably, cytoskeletal protein alpha tubulin (*TUA*) family and
321 motor protein dynein heavy chain (*DYH*) family are rapidly expanded in the *S.*
322 *constricta* genome ([Figure 3A](#)). They play vital roles in the microtubule architecture
323 and the bending movement of cilia (Mohri et al., 2012). The razor clam has a
324 well-developed ciliary system for gill filtering, food-particles retaining, and water
325 pumping (Morton, 1984). The adjoining cilia generate effective beat through
326 coordinated wavelike movements. The pumping rate of the ciliary system in the gill
327 and mantle cavity can be adjusted to generate powerful currents to facilitate the
328 principal sorting and retaining of suspended particles in the labial palps. Effluxes can
329 also be ejected from the pedal gape to flush away sources of irritation detected by the

330 sensory tentacles (Morton, 1984). The transcriptomic data revealed that the *TUA* and
331 *DYH* genes are highly expressed in the gills (Figure 3B and 3C), suggesting that the
332 expansion of these genes could be an adaptation to the deep-burrowing lifestyle.

333 Phylogenetic inference of the 16 species was performed with the 337 single-copy
334 orthologs. Multiple sequences alignment was conducted for the protein sequences of
335 each ortholog gene using MUSCLE (version 3.8.31) (Edgar, 2004) separately. The
336 alignments for all the orthologs were then concatenated into a super alignment matrix
337 with 241,349 amino acids. RAxML (version 8.2.12) (Stamatakis, 2014) was used to
338 infer the alignment matrix by a maximum likelihood method with the substitution
339 model PROTGAMMAAUTO. Bootstrapping with 100 replicates was used for node
340 support. Divergence time between species was estimated using MCMCTree in PAML
341 package (version 4.7a) (Yang, 1997) with the parameters of “burn-in = 1,000,
342 sample-number = 1,000,000, sample-frequency = 2”. The constructed maximum
343 likelihood phylogenetic tree revealed that *S. constricta* clustered with other bivalve
344 species and diverged ~494 million years ago (Mya) from the lineage leading to
345 oysters and scallops (Figure 2).

346

347 **Hi-C scaffolding and macro-synteny analysis**

348 Adductor muscle tissue of a razor clam from the same population was collected for
349 Hi-C library construction. The tissue specimen was fixed with 1% formaldehyde and
350 the genomic DNA was cross-linked, digested by restriction enzymes HindIII, labeled
351 with biotinylated residue, and end repaired. The library was sequenced on an Illumina
352 NovaSeq platform, generating 156.73 G of raw reads. The raw reads were truncated at
353 the junctions and aligned to the polished genome using BWA (version 0.7.17) with
354 default parameters. Only the unique aligned reads with a mapping quality over 20
355 were further processed. After filtering invalid interaction pairs by HiC-Pro (v.2.8.0)
356 (Servant et al., 2015), 30.32% of the clean reads were valid pairs and utilized to
357 evaluate the interaction strength among whole genome contigs. Lachesis (version
358 2e27abb) was used to cluster and anchor the contigs to the chromosomes using an
359 agglomerative hierarchical clustering method (Burton et al., 2013). Finally, 3,068

360 contigs, accounting for 87.82% of the total bases, were clustered into 19 linkage
361 groups (Figure 4A), which was consistent with the karyotype revealed by previous
362 studies (Wang et al., 1998). The ancient ortholog genes exhibited remarkable
363 preservation of ancestral bilaterian linkage groups (Simakov et al., 2013; Wang et al.,
364 2017) with a conservation index (CI) of 0.71, indicating the considerable accuracy of
365 the Hi-C clustering (Figure 4B).

366

367 **Conclusions**

368 We assembled a high-quality chromosomal-level reference genome of *S. constricta*,
369 the first genome of the family Solenidae, along with comprehensive transcriptomic
370 data of whole-ontogeny developmental stages and all major tissues (under normal and
371 stressed conditions). The genomic and transcriptomic resources reported here would
372 lay a prime foundation for future studies to elucidate the razor clam' adaptive traits
373 relating to deep-burrowing lifestyle (e.g., thin shells, advanced ciliary and siphon
374 system, , extraordinary adaption to broad-range salinity and high concentration of
375 ammonia nitrogen and hydrogen sulfide) and genetic improvement in commercial
376 breeding programs.

377

378 **Ethics approval and consent to participate**

379 All experimental procedures were approved by the Institutional Animal Care and
380 Use Committee (IACUC) of Zhejiang Wanli University, China. All participants
381 consent to publish the work under the "Consent to publish" heading.

382

383 **Data availability**

384 The *S. constricta* genome assembly is available at NCBI (BioProject:
385 PRJNA559038). RNA sequencing data files are available through the NCBI Sequence
386 Read Archive (BioProject: PRJNA559056). The *S. constricta* genome assembly and
387 annotation files also could be downloaded from the website
388 <http://202.121.66.128/clam-genome/zwu.htm>.

389

390 **Acknowledgements**

391 This work was financially supported by the National Key Research and
392 Development Program of China (No.2018YFD0901405), Zhejiang Major Program of
393 Science and Technology (No.2016C02055-9), Modern Agro-industry Technology
394 Research System (No. CARS-49).

395

396 **Author contributions**

397 Y. D., Z. L. and Z. B. conceived the project. H. Y., W. R. and L.H. conducted the
398 environmental stress and collected the samples. Q. Z., S. W., J. R. and L. L. performed
399 the genome assembly, annotation, transcriptome analysis and other bioinformatics
400 analysis. Y. D., J. R., Q. Z., S. W. and Q. X. wrote and revised the manuscript. All
401 authors read and approved the final manuscript.

402

403

404 **Competing interests**

405 The authors declare that they have no competing interests.

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420 **Reference**

- 421 Bao, W., Kojima, K.K., and Kohany, O. (2015). Repbase Update, a database of repetitive elements
422 in eukaryotic genomes. *Mobile DNA* 6, 11.
- 423 Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids*
424 *research* 27, 573-580.
- 425 Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
426 sequence data. *Bioinformatics* 30, 2114-2120.
- 427 Burge, C., and Karlin, S. (1997). Prediction of complete gene structures in human genomic DNA.
428 *Journal of molecular biology* 268, 78-94.
- 429 Burton, J.N., Adey, A., Patwardhan, R.P., Qiu, R., Kitzman, J.O., and Shendure, J. (2013).
430 Chromosome-scale scaffolding of de novo genome assemblies based on chromatin interactions.
431 *Nature biotechnology* 31, 1119-1125.
- 432 Chin, C.S., Peluso, P., Sedlazeck, F.J., Nattestad, M., Concepcion, G.T., Clum, A., Dunn, C.,
433 O'Malley, R., Figueroa-Balderas, R., Morales-Cruz, A., *et al.* (2016). Phased diploid genome
434 assembly with single-molecule real-time sequencing. *Nature methods* 13, 1050-1054.
- 435 De Bie, T., Cristianini, N., Demuth, J.P., and Hahn, M.W. (2006). CAFE: a computational tool for
436 the study of gene family evolution. *Bioinformatics* 22, 1269-1271.
- 437 Doerks, T., Copley, R.R., Schultz, J., Ponting, C.P., and Bork, P. (2002). Systematic identification
438 of novel protein domain families associated with nuclear functions. *Genome research* 12, 47-56.
- 439 Dong, Y.H., Yao, H.H., Zhang, P.Y., Shen, P.Y., Liiu, H.M., and Lin, Z.H. (2012). Cytological
440 observation on fertilization and early cleavage in *Sinonovaula constricta*. *Journal of fisheries of*
441 *China* 36, 1400-1409.
- 442 Edgar, R.C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high
443 throughput. *Nucleic acids research* 32, 1792-1797.
- 444 FAO (2018). FAO yearbook. Fishery and Aquaculture Statistics 2016.
445 <http://www.fao.org/fishery/publications/yearbooks/en>.
- 446 Green, M., and Sambrook, J. (2012). *Molecular Cloning: A Laboratory Manual*. 4th Edition, Vol.
447 II, Cold Spring Harbor Laboratory Press, New York.
- 448 Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J., White, O., Buell, C.R., and
449 Wortman, J.R. (2008). Automated eukaryotic gene structure annotation using EVIDENCEModeler
450 and the Program to Assemble Spliced Alignments. *Genome biology* 9, R7.
- 451 Jones, P., Binns, D., Chang, H.Y., Fraser, M., Li, W., McAnulla, C., McWilliam, H., Maslen, J.,
452 Mitchell, A., Nuka, G., *et al.* (2014). InterProScan 5: genome-scale protein function classification.
453 *Bioinformatics* 30, 1236-1240.
- 454 Kalvari, I., Nawrocki, E.P., Argasinska, J., Quinones-Olvera, N., Finn, R.D., Bateman, A., and
455 Petrov, A.I. (2018). Non-Coding RNA Analysis Using the Rfam Database. *Current protocols in*
456 *bioinformatics* 62, e51.
- 457 Kim, D., Langmead, B., and Salzberg, S.L. (2015). HISAT: a fast spliced aligner with low memory
458 requirements. *Nature methods* 12, 357-360.
- 459 Leskovec, J., and Soscic, R. (2016). SNAP: A General Purpose Network Analysis and Graph
460 Mining Library. *ACM transactions on intelligent systems and technology* 8.
- 461 Li, L., Stoeckert, C.J., Jr., and Roos, D.S. (2003). OrthoMCL: identification of ortholog groups for
462 eukaryotic genomes. *Genome research* 13, 2178-2189.
- 463 Li, W., and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of

- 464 protein or nucleotide sequences. *Bioinformatics* 22, 1658-1659.
- 465 Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose program for
466 assigning sequence reads to genomic features. *Bioinformatics* 30, 923-930.
- 467 Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion
468 for RNA-seq data with DESeq2. *Genome biology* 15, 550.
- 469 Lowe, T.M., and Chan, P.P. (2016). tRNAscan-SE On-line: integrating search and context for
470 analysis of transfer RNA genes. *Nucleic acids research* 44, W54-57.
- 471 Majoros, W.H., Pertea, M., and Salzberg, S.L. (2004). TigrScan and GlimmerHMM: two open
472 source ab initio eukaryotic gene-finders. *Bioinformatics* 20, 2878-2879.
- 473 Mohri, H., Inaba, K., Ishijima, S., and Baba, S.A. (2012). Tubulin-dynein system in flagellar and
474 ciliary movement. *Proceedings of the Japan Academy Series B, Physical and biological sciences*
475 88, 397-415.
- 476 Morton, B. (1984). The functional morphology of *Sinonovucufu constvictu* with a discussion on
477 the taxonomic status of the Novaculininae (Bivalvia). *J Zool, Lond* 202, 299-325.
- 478 Parra, G., Blanco, E., and Guigo, R. (2000). GeneID in *Drosophila*. *Genome research* 10, 511-515.
- 479 Parra, G., Bradnam, K., and Korf, I. (2007). CEGMA: a pipeline to accurately annotate core genes
480 in eukaryotic genomes. *Bioinformatics* 23, 1061-1067.
- 481 Salmela, L., and Rivals, E. (2014). LoRDEC: accurate and efficient long read error correction.
482 *Bioinformatics* 30, 3506-3514.
- 483 Servant, N., Varoquaux, N., Lajoie, B.R., Viara, E., Chen, C.J., Vert, J.P., Heard, E., Dekker, J.,
484 and Barillot, E. (2015). HiC-Pro: an optimized and flexible pipeline for Hi-C data processing.
485 *Genome biology* 16, 259.
- 486 Simakov, O., Marletaz, F., Cho, S.J., Edsinger-Gonzales, E., Havlak, P., Hellsten, U., Kuo, D.H.,
487 Larsson, T., Lv, J., Arendt, D., *et al.* (2013). Insights into bilaterian evolution from three spiralian
488 genomes. *Nature* 493, 526-531.
- 489 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of
490 large phylogenies. *Bioinformatics* 30, 1312-1313.
- 491 Stanke, M., Keller, O., Gunduz, I., Hayes, A., Waack, S., and Morgenstern, B. (2006).
492 AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic acids research* 34, W435-439.
- 493 Trapnell, C., Pachter, L., and Salzberg, S.L. (2009). TopHat: discovering splice junctions with
494 RNA-Seq. *Bioinformatics* 25, 1105-1111.
- 495 Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg, S.L.,
496 Wold, B.J., and Pachter, L. (2010). Transcript assembly and quantification by RNA-Seq reveals
497 unannotated transcripts and isoform switching during cell differentiation. *Nature biotechnology*
498 28, 511-515.
- 499 Varshney, R.K., Chen, W., Li, Y., Bharti, A.K., Saxena, R.K., Schlueter, J.A., Donoghue, M.T.,
500 Azam, S., Fan, G., Whaley, A.M., *et al.* (2011). Draft genome sequence of pigeonpea (*Cajanus*
501 *cajan*), an orphan legume crop of resource-poor farmers. *Nature biotechnology* 30, 83-89.
- 502 Walker, B.J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C.A., Zeng,
503 Q., Wortman, J., Young, S.K., *et al.* (2014). Pilon: an integrated tool for comprehensive microbial
504 variant detection and genome assembly improvement. *PloS one* 9, e112963.
- 505 Wang, J.X., Zhao, X.F., Zhou, L.H., and Xiang, J.H. (1998). Chromosome study of *Sinonovacula*
506 *constricta* (Bivalvia) *Oceanologia Et Limnologia Sinica* 29, 191-196.
- 507 Wang, S., Zhang, J., Jiao, W., Li, J., Xun, X., Sun, Y., Guo, X., Huan, P., Dong, B., Zhang, L., *et*

508 *al.* (2017). Scallop genome provides insights into evolution of bilaterian karyotype and
509 development. *Nature ecology & evolution* *1*, 120.
510 Waterhouse, R.M., Seppey, M., Simao, F.A., Manni, M., Ioannidis, P., Klioutchnikov, G.,
511 Kriventseva, E.V., and Zdobnov, E.M. (2017). BUSCO applications from quality assessments to
512 gene prediction and phylogenomics. *Molecular biology and evolution*.
513 Xu, Z., and Wang, H. (2007). LTR_FINDER: an efficient tool for the prediction of full-length LTR
514 retrotransposons. *Nucleic acids research* *35*, W265-268.
515 Yang, Z. (1997). PAML: a program package for phylogenetic analysis by maximum likelihood.
516 *Computer applications in the biosciences : CABIOS* *13*, 555-556.
517 Yu, G., Wang, L.G., Han, Y., and He, Q.Y. (2012). clusterProfiler: an R package for comparing
518 biological themes among gene clusters. *Omics : a journal of integrative biology* *16*, 284-287.
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545

546 **Tables**

547

548 **Table 1.** Statistics of the genome assembly of *Sinonovacula constricta*

Types	Number	Length (bp)	N50 (bp)	Maximum (bp)	GC content (%)
Contig	10,981	1,331,972,725	678,857	5,402,231	35.46
Scaffold	7,932	1,332,277,427	57,991,182	93,300,556	35.46

549

550

551

552

553 **Table 2.** Statistics of the repetitive sequences

Type	Repeat size (bp)	% of genome
DNA	191,499,094	14.38
LINE	71,938,692	5.4
LTR	144,451,530	10.84
SINE	5,528,172	0.42
Tandem repeat	204,889,587	15.38
Other	157,232	0.01
Unknown	56,940,582	4.27
Total	675,404,889	50.71

554

555

556

557

558

559

560

561 **Table 3.** Statistics of gene annotation to different databases

Annotation database	Number of annotated genes	Percentage (%)
NR	23,844	90.88
Swiss-Prot	18,131	69.1
KEGG	18,928	72.14
InterProScan	25,475	97.1
Pfam	15,391	58.66
GO	22,956	87.49
Total	26,140	99.50

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

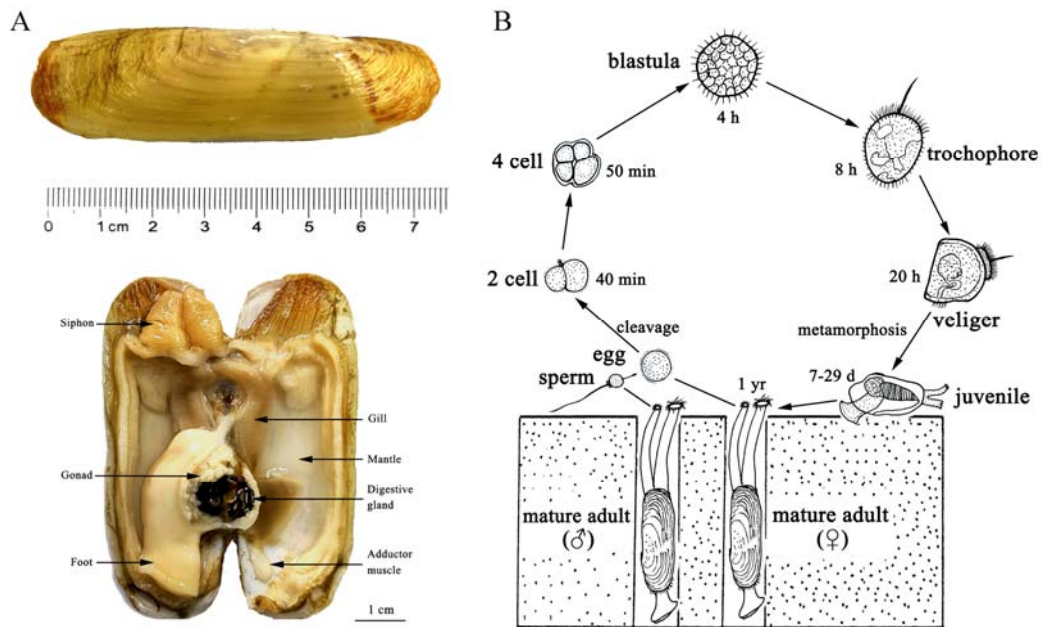
577

578

579

580 **Figures**

581



582

583 **Figure 1.** A. The appearance and anatomic structures of an adult razor clam. B. A

584 pelago-benthic life cycle of the razor clam.

585

586

587

588

589

590

591

592

593

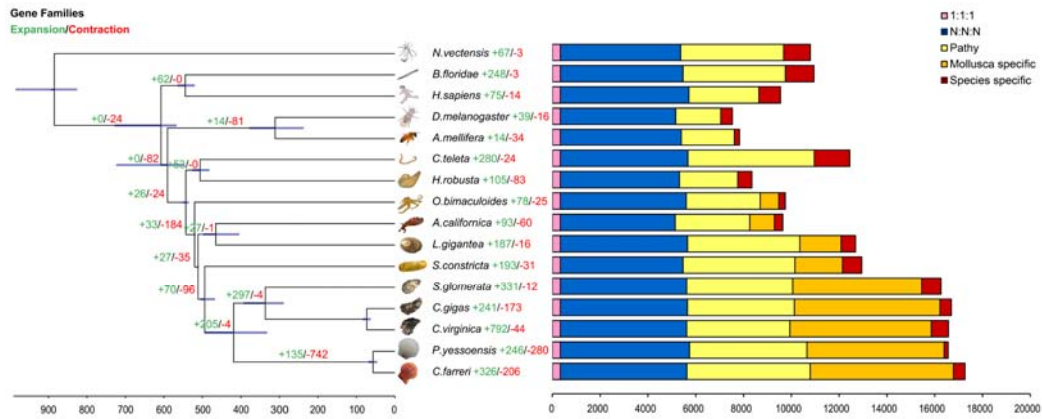
594

595

596

597

598



599

600 **Figure 2.** Phylogenetic tree and number of shared orthologs among *S. constricta* and

601 other animal species. Numbers of gene families undergoing expansion and contraction

602 for each lineage are exhibited as red and green, respectively.

603

604

605

606

607

608

609

610

611

612

613

614

615

616

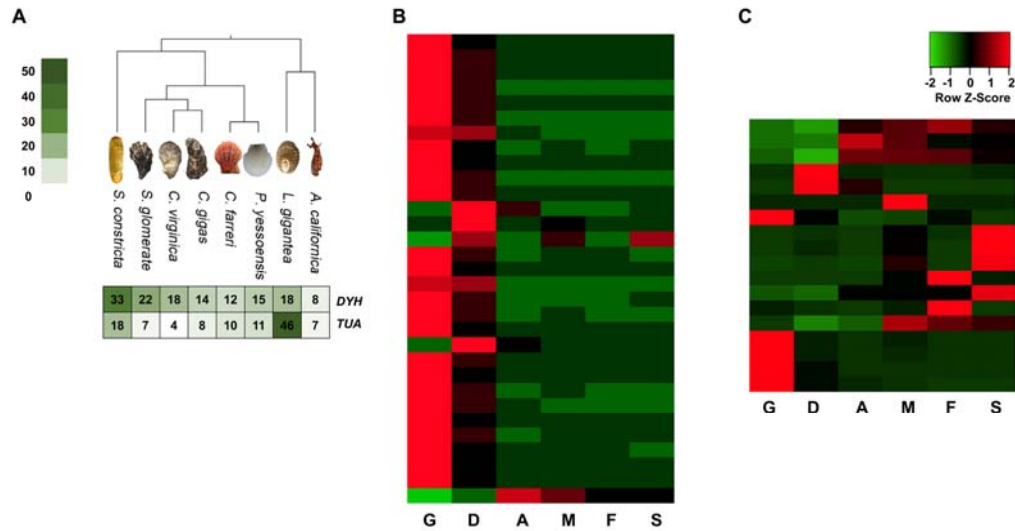
617

618

619

620

621



622

623 **Figure 3. A.** The comparison of the copy numbers of dynein heavy chain (*DYH*) and
 624 alpha tubulin (*TUA*) genes in 8 molluscan species. **B & C.** The tissue-wide expression
 625 patterns of *DYH* genes and *TUA* genes. Abbreviations: G, gill; D, digestive gland; A,
 626 adductor muscle; M, mantle; F, foot; S, siphon.

627

628

629

630

631

632

633

634

635

636

637

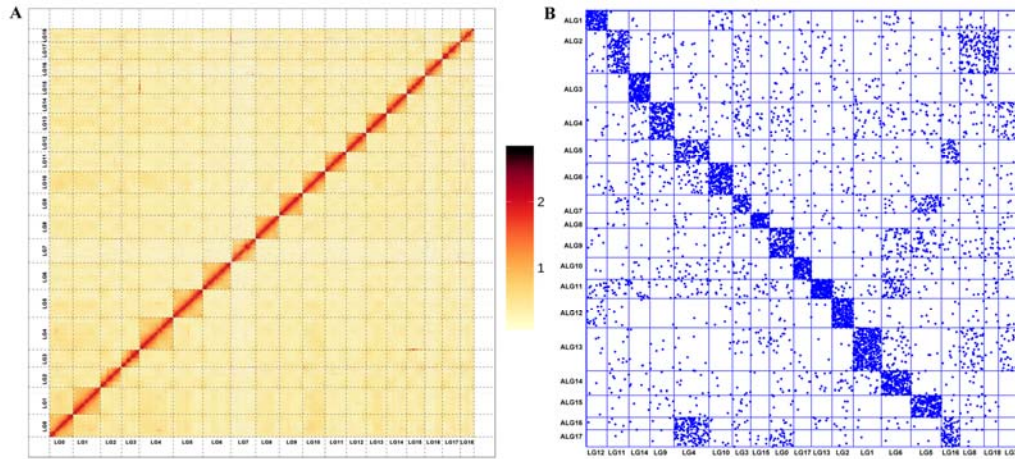
638

639

640

641

642



643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

Figure 4. A. Hi-C interaction heat map of of *S. constricta*. **B.** Chromosome-based macro-synteny between *S. constricta* and the 17 presumed bilaterian ALGs retrieved from Simakov et al. (2013).

666 **Supplementary materials:**

667 **Table S1.** Summary of the genomic sequencing reads

668 **Table S2.** Statistics of Illumina short reads coverage

669 **Table S3.** Summary genomic completeness by CEGMA

670 **Table S4.** Summary genomic completeness by BUSCO

671 **Table S5.** Summary of Illumina transcriptome sequencing data

672 **Table S6.** Summary of clean reads mapping

673 **Table S7.** Summary of PacBio full-length transcriptome sequencing

674 **Table S8.** Summary of the gene prediction results

675 **Table S9.** Summary of the non-coding RNA annotation

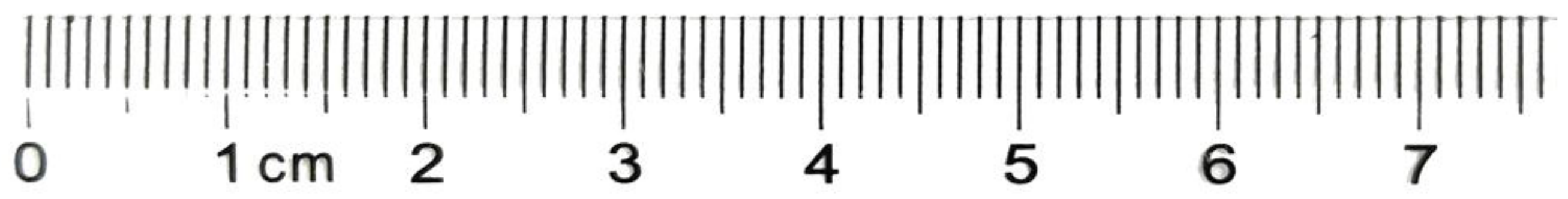
676 **Figure S1.** Genome survey of *Sinonovacula constricta* using 17-mer analysis

677 **Figure S2.** Volcano map of differentially expressed genes

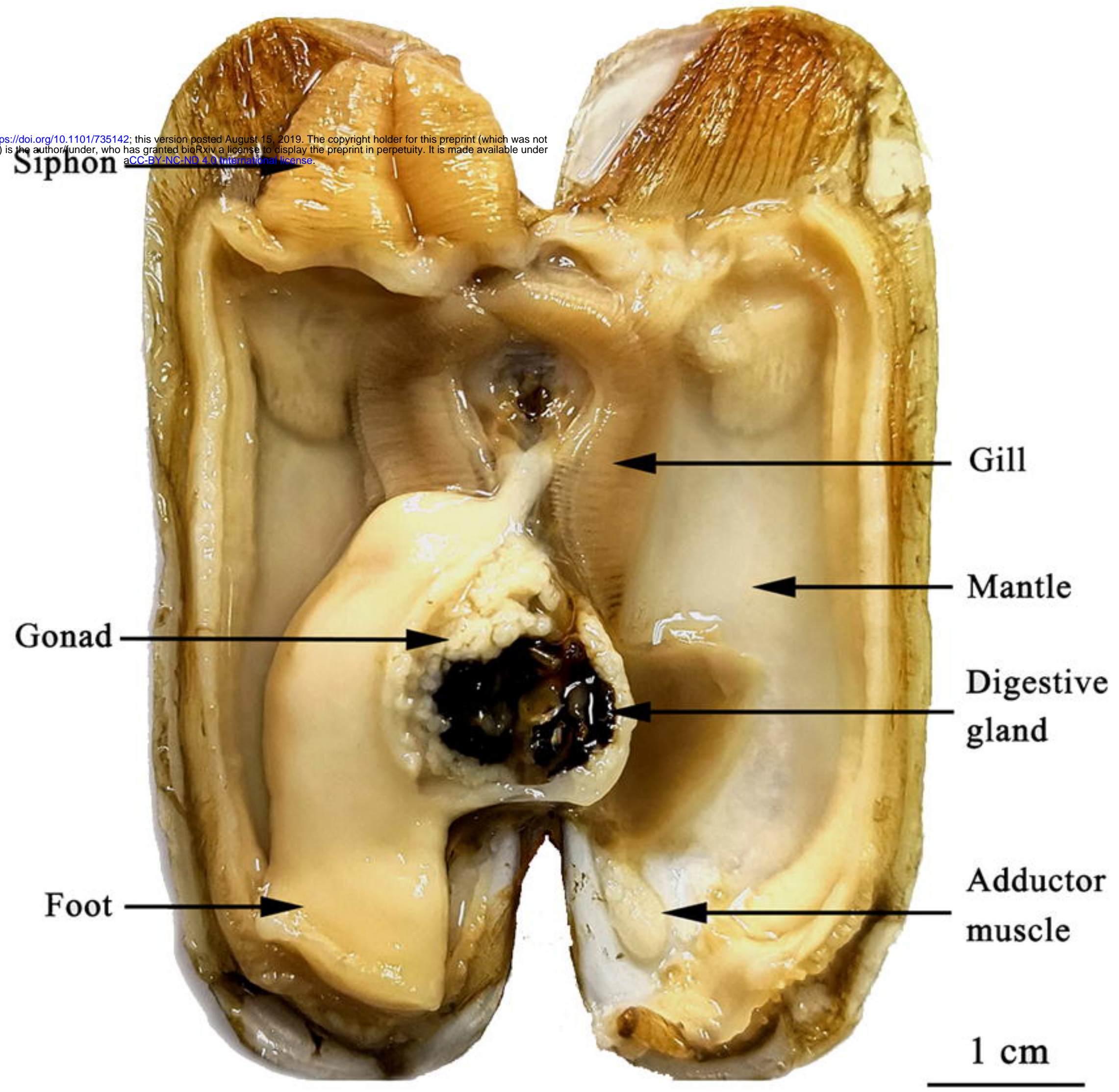
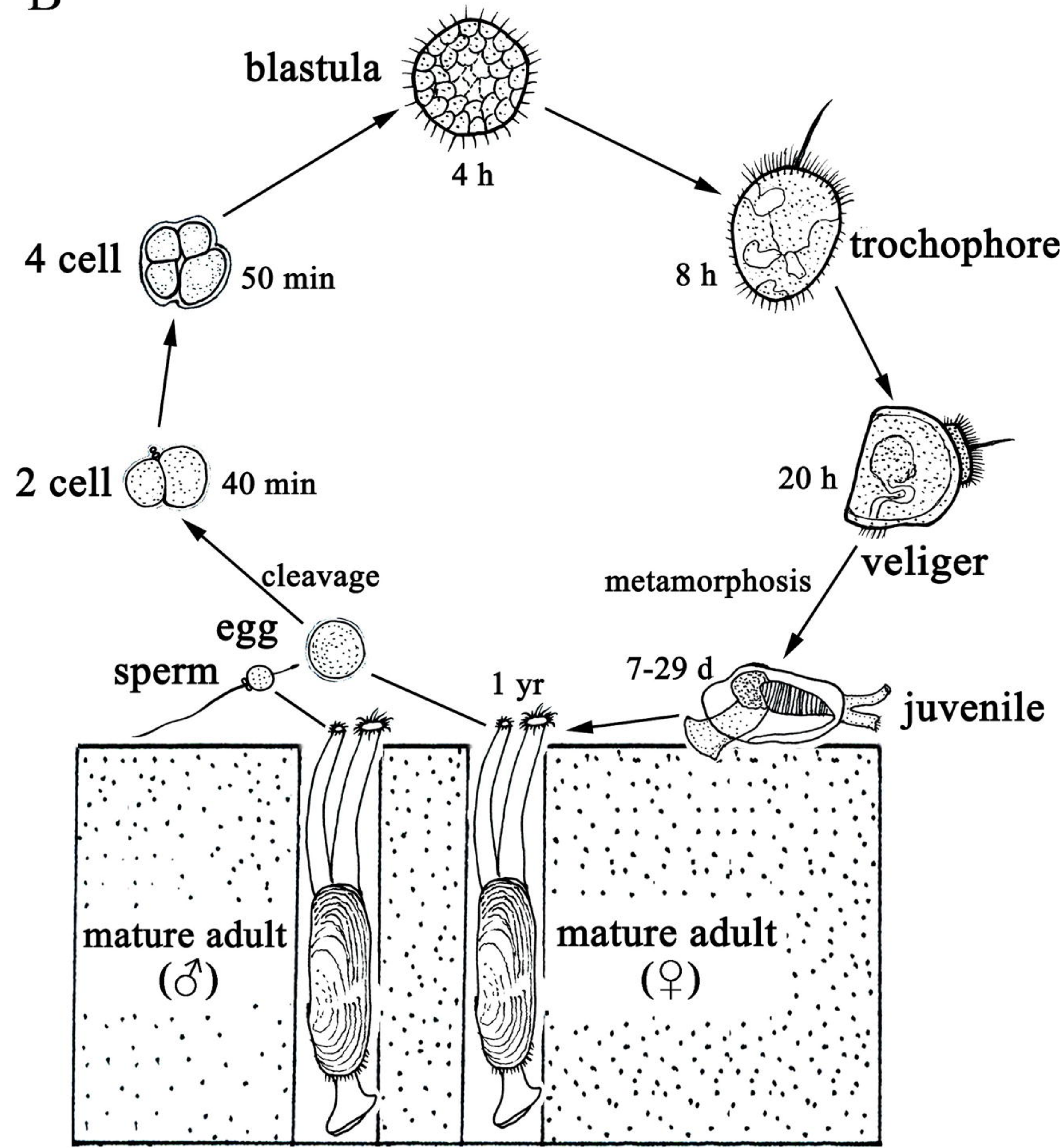
678 **Figure S3.** Dot plot of GO enrichment of differentially expressed genes

679 **Figure S4.** Dot plot of KEGG pathway enrichment of differentially expressed gen

680

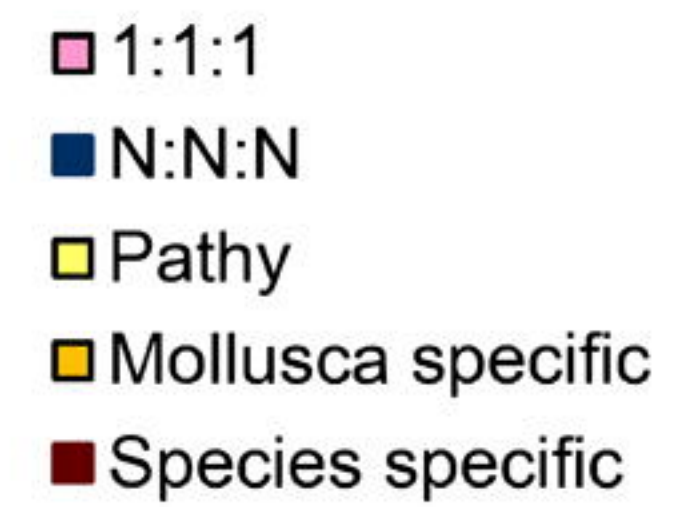
A

bioRxiv preprint doi: <https://doi.org/10.1101/735142>; this version posted August 16, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

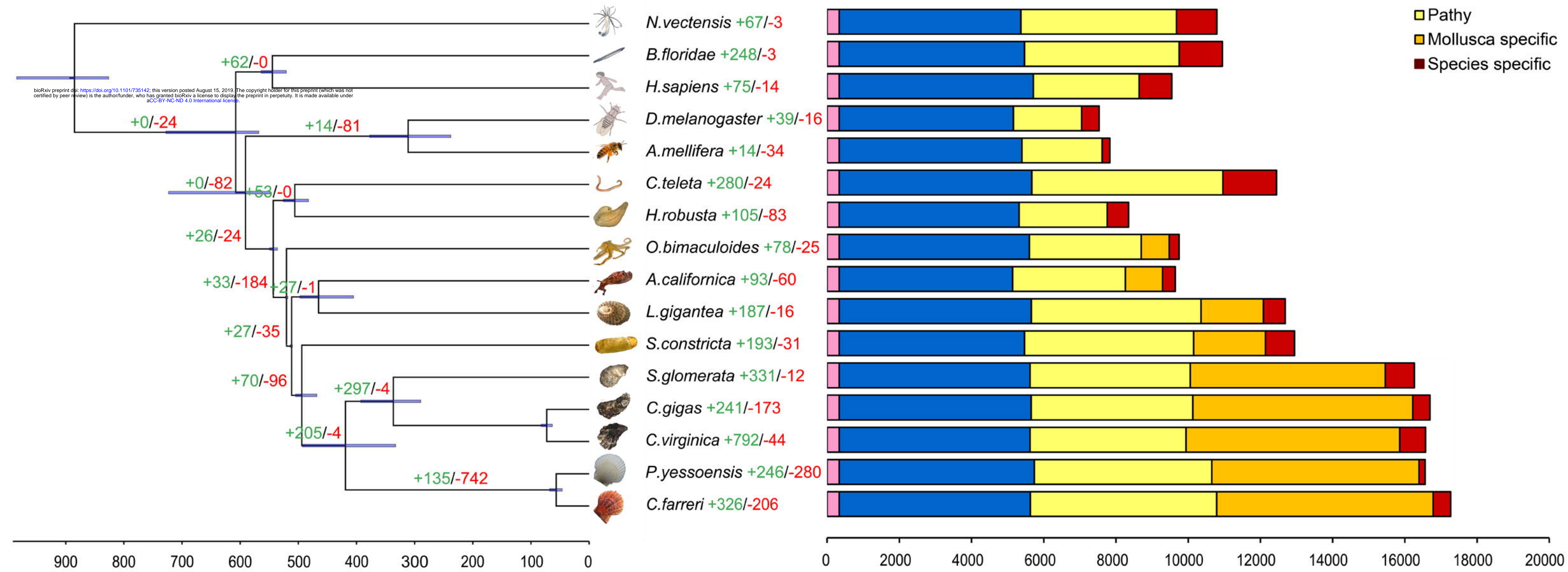
**B**

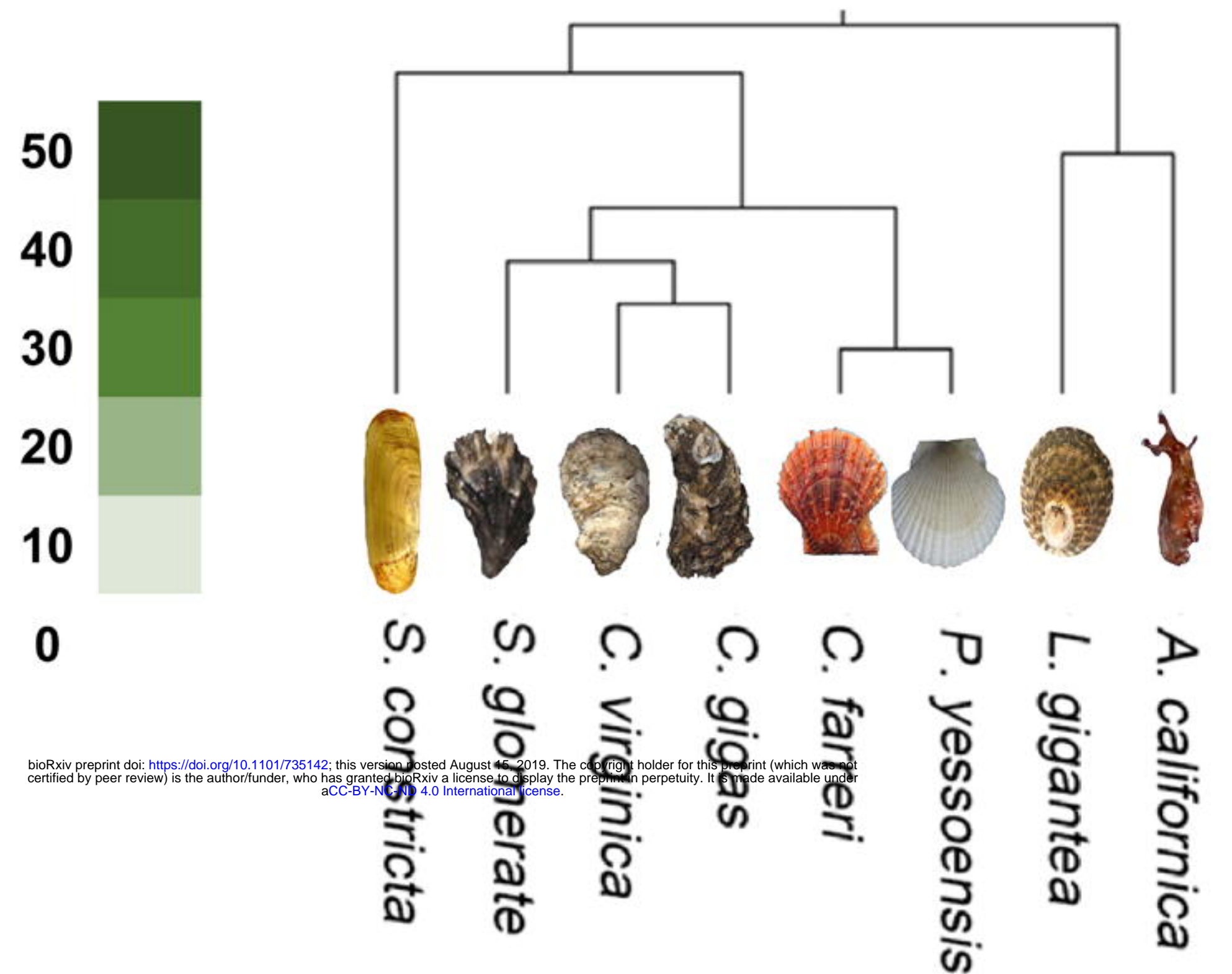
Gene Families

Expansion/Contraction



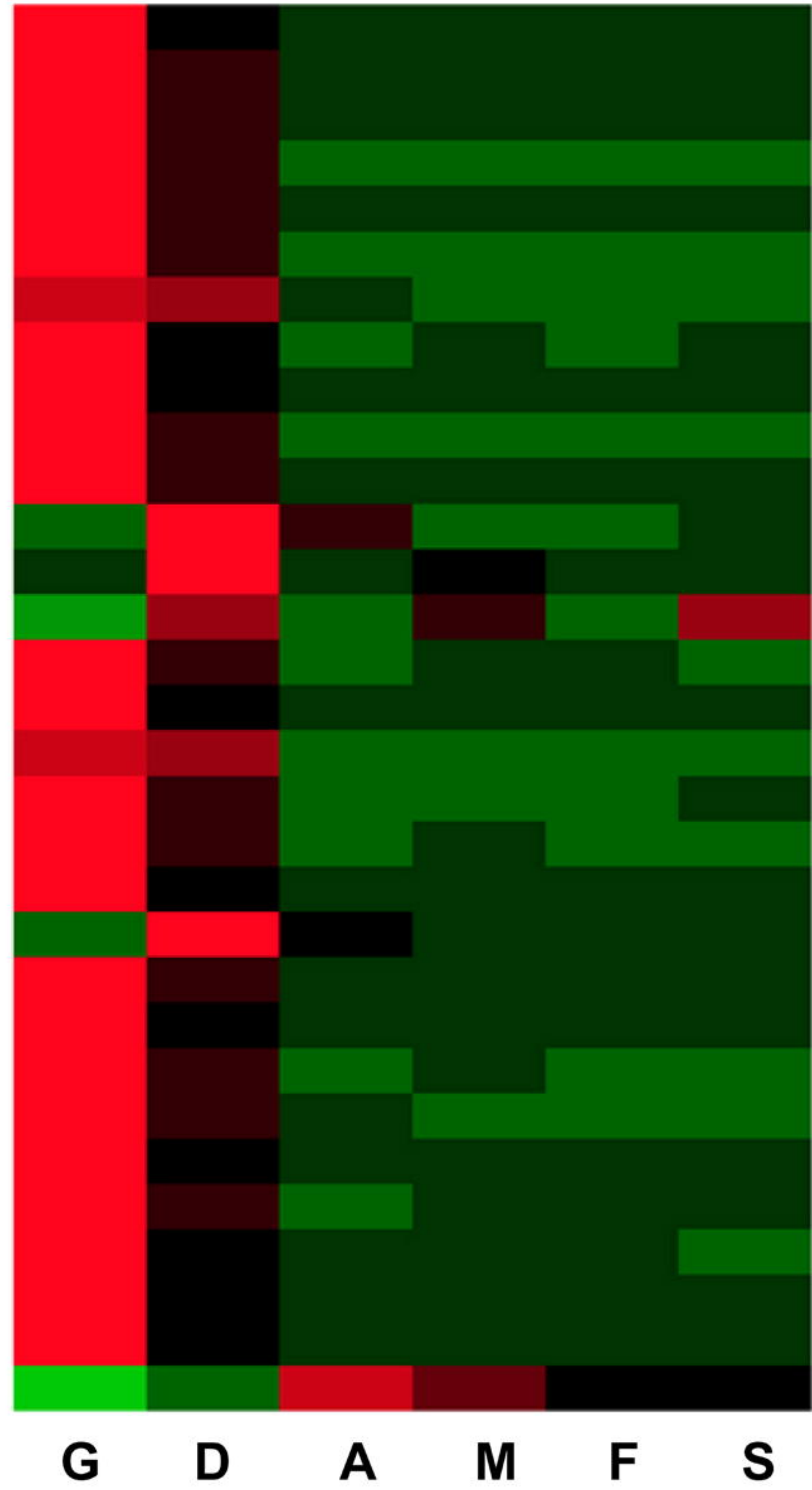
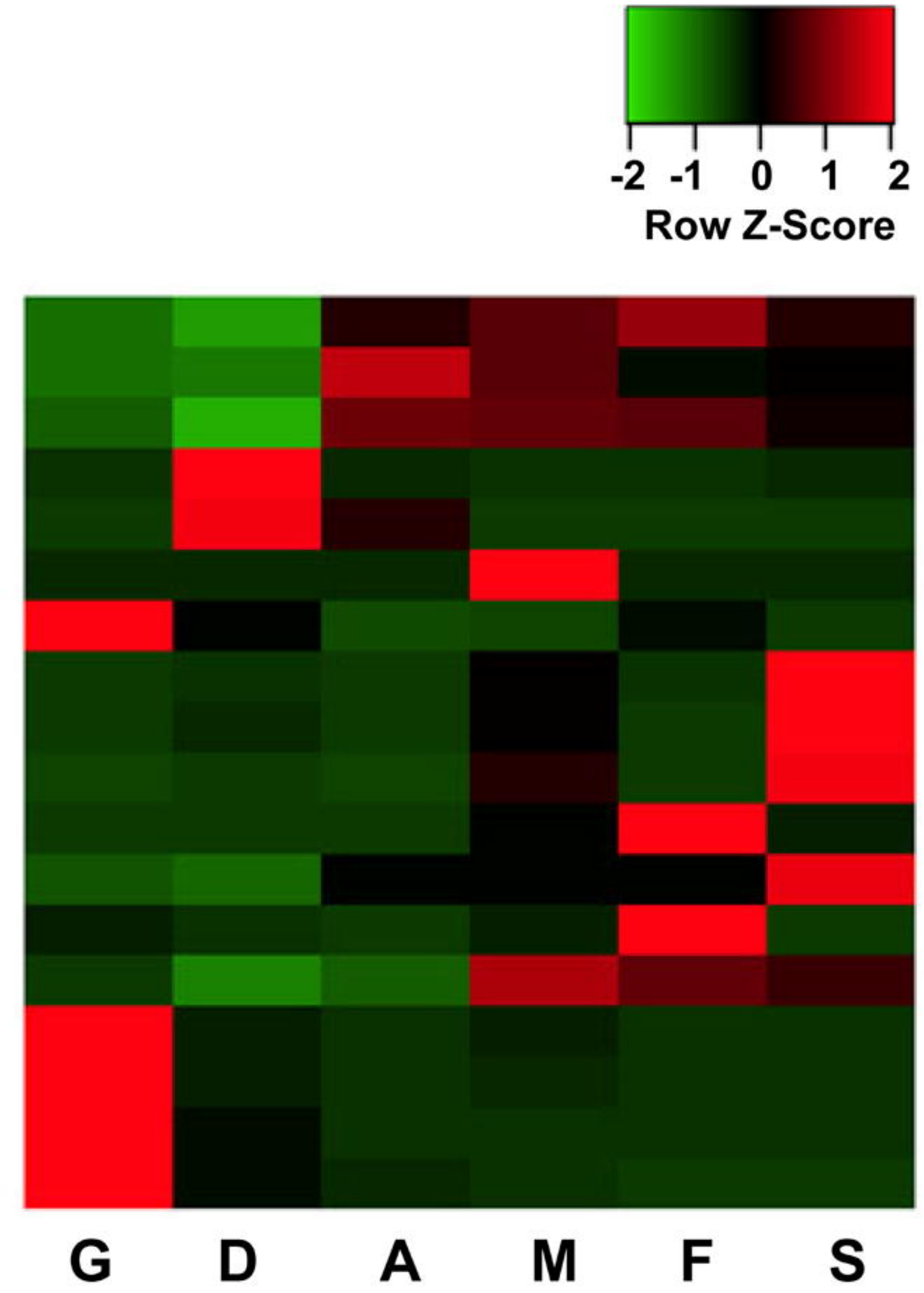
bioRxiv preprint doi: <https://doi.org/10.1101/735142>; this version posted August 15, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

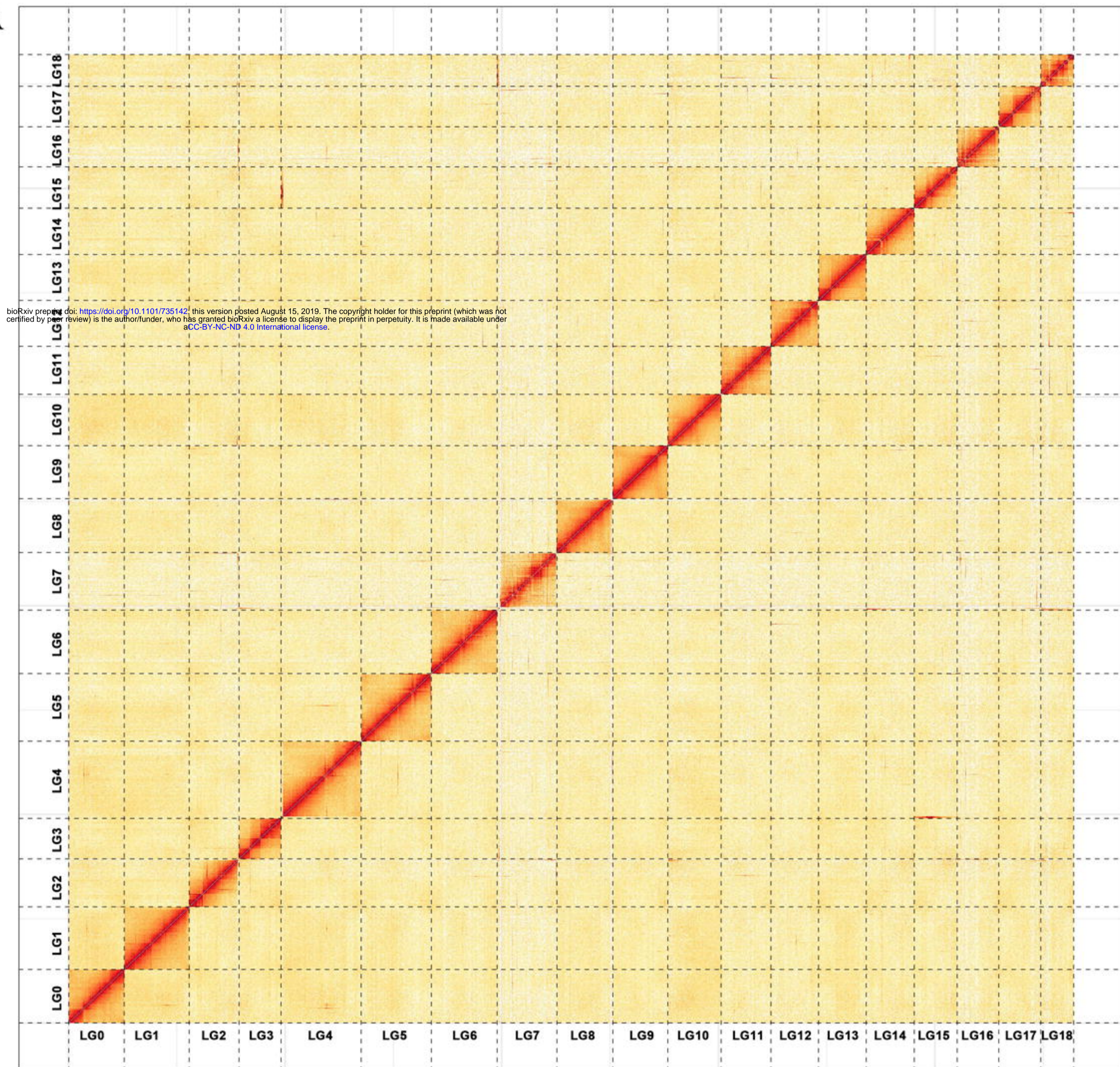


A

bioRxiv preprint doi: <https://doi.org/10.1101/735142>; this version posted August 4, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

33	22	18	14	12	15	18	8	<i>DYH</i>
18	7	4	8	10	11	46	7	<i>TUA</i>

B**C**

A**B**