

<b>Manuscript Number:</b>	GIGA-D-17-00315R1	
<b>Full Title:</b>	Draft genome of the Peruvian scallop <i>Argopecten purpuratus</i>	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	National Natural Science Foundation of China (31572618)	Dr. Chunde Wang
	National Natural Science Foundation of China (41676152)	Dr. Xiao Liu
	Fund for Shandong Modern Agro-Industry Technology Research System (SDAIT-14)	Dr. Chunde Wang
<b>Abstract:</b>	<p>Background: The Peruvian scallop, <i>Argopecten purpuratus</i>, is mainly cultured in Southern Chile and had been introduced into China in last century. Unlike other <i>Argopecten</i> scallops, the Peruvian scallop normally has a long life span of up to 7-10 years. Therefore, researchers have been employing it to develop hybrid vigor. Here, we performed whole genome sequencing, assembly, and gene annotation of the Peruvian scallop, with an important aim to develop genomic resources for genetic breeding in scallops. Findings: A total of 463.19-Gb (Gigabase) raw DNA reads were sequenced. The draft genome assembly of 724.78 Mb was generated (accounting for 81.87% of the estimated genome size 885.29 Mb), with a contig N50 size of 80.11 kb and scaffold N50 size of 1.02 Mb. Meanwhile, the repeat sequences were calculated to reach 33.74% of the whole genome, and a total of 26,256 protein-coding genes and 3,057 non-coding RNAs were predicted from the assembly. Conclusion: We generated a draft genome assembly of the Peruvian scallop, which will provide solid resource for further genetic breeding and evolutionary history analysis of this economically important scallop.</p>	
<b>Corresponding Author:</b>	Chunde Wang, Ph.D Qingdao Agricultural University CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Qingdao Agricultural University	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Chao Li, Ph.D	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Chao Li, Ph.D	
	Xiao Liu, D.Sc	
	Bo Liu, D.Sc	
	Bin Ma, MSc	
	Guilong Liu, MSc	
	Qiong Shi, Ph.D	
	Chunde Wang, Ph.D	
<b>Order of Authors Secondary Information:</b>		
<b>Response to Reviewers:</b>	Dear Editor, Thanks for your kindly consideration of publication of our manuscript in the journal. We	

	<p>appreciate the insightful comments from the anonymous reviewers very much. I have included the response letter in the supplementary file as the figures will not show up in this box. If you have more questions, please feel free to let us know in your earliest convenience.</p> <p>Best regards,</p> <p>Chao and Chunde</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p>	Yes

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

1 Draft genome of the Peruvian scallop *Argopecten purpuratus*

2 Chao Li<sup>1</sup>, Xiao Liu<sup>2</sup>, Bo Liu<sup>1</sup>, Bin Ma<sup>3</sup>, Guilong Liu<sup>1</sup>, Qiong Shi<sup>4</sup>, Chunde Wang<sup>1,\*</sup>

3  
4 <sup>1</sup>Marine Science and Engineering College, Qingdao Agricultural University, Qingdao  
5 266109, China

6 <sup>2</sup>Key Laboratory of Experimental Marine Biology, Institute of Oceanology, Chinese  
7 Academy of Sciences, Qingdao 266071, China

8 <sup>3</sup>Qingdao Oceanwide BioTech Co., Ltd., Qingdao 266101, China

9 <sup>4</sup>Shenzhen Key Lab of Marine Genomics, Guangdong Provincial Key Lab of  
10 Molecular Breeding in Marine Economic Animals, BGI Academy of Marine Sciences,  
11 BGI Marine, BGI, Shenzhen 518083, China

12  
13 \*Correspondence address. Professor Chunde Wang, Marine Science and Engineering  
14 College, Qingdao Agricultural University, Qingdao 266109, China (tel:  
15 +8613589227997; email: chundewang2007@163.com)

16  
17 Email addresses of all authors:

18 Chao Li, leochao@163.com

19 Xiao Liu, liuxiao@qdio.ac.cn

20 Bo Liu, liubomusic@126.com

21 Bin Ma, hereiammabin@163.com

22 Guilong Liu, 969613442@qq.com

23 Qiong Shi, shiqiong@genomics.cn

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 24 **Abstract**

2  
3  
4 25 **Background:** The Peruvian scallop, *Argopecten purpuratus*, is mainly cultured in  
5  
6  
7 26 Southern Chile and had been introduced into China in last century. Unlike other  
8  
9  
10 27 *Argopecten* scallops, the Peruvian scallop normally has a long life span of up to 7-10  
11  
12  
13 28 years. Therefore, researchers have been employing it to develop hybrid vigor. Here,  
14  
15  
16 29 we performed whole genome sequencing, assembly, and gene annotation of the  
17  
18  
19 30 Peruvian scallop, with an important aim to develop genomic resources for genetic  
20  
21  
22 31 breeding in scallops. **Findings:** A total of 463.19-Gb (Gigabase) raw DNA reads were  
23  
24  
25 32 sequenced. The draft genome assembly of 724.78 Mb was generated (accounting for  
26  
27  
28 33 81.87% of the estimated genome size 885.29 Mb), with a contig N50 size of 80.11 kb  
29  
30  
31 34 and scaffold N50 size of 1.02 Mb. Meanwhile, the repeat sequences were calculated  
32  
33  
34 35 to reach 33.74% of the whole genome, and a total of 26,256 protein-coding genes and  
35  
36  
37 36 3,057 non-coding RNAs were predicted from the assembly. **Conclusion:** We  
38  
39  
40 37 generated a draft genome assembly of the Peruvian scallop, which will provide solid  
41  
42  
43 38 resource for further genetic breeding and evolutionary history analysis of this  
44  
45  
46 39 economically important scallop.  
47  
48

49  
50 41 **Keywords:** *Argopecten purpuratus*; Peruvian scallop; genome assembly; annotation;  
51  
52  
53 42 gene prediction; phylogenetic analysis  
54  
55  
56 43  
57  
58 44  
59  
60  
61  
62  
63  
64  
65

1 45

2  
3 46 **Data description**

4  
5 47

6  
7 48 Introduction

8  
9  
10 49 The Peruvian scallop (*Argopecten purpuratus*), also known as the Chilean  
11  
12 50 scallop, is a medium-sized bivalve with a wide distribution in Peru and Chile [1]. In  
13  
14 51 Chile, the cultured scallops reach a commercial size of around 9 cm in shell height  
15  
16 52 within 14-16 months [2]. It is a relatively stenothermic species as its natural habitat is  
17  
18 53 largely under the influence of upwelling currents from Antarctica [3]. Unlike other  
19  
20 54 *Argopecten* scallops, the Peruvian scallop normally has a long life span of up to 7-10  
21  
22 55 years [4, 5]. This scallop was introduced into China in the late 2000's and had played  
23  
24 56 an important role in stock improvement of *Argopecten* scallops via inter-specific  
25  
26 57 hybridization [6, 7] with bay scallops.  
27  
28  
29  
30

31 58

32  
33  
34 59 Whole genome sequencing

35  
36  
37 60 Genomic DNA was extracted from adductor muscle sample of a single *A.*  
38  
39 61 *purpuratus* (Figure 1), which was obtained from a local scallop farm in Laizhou,  
40  
41 62 Shandong Province, China. A whole genome shotgun sequencing strategy was applied.  
42  
43 63 Six libraries with different insert length (250 bp, 450 bp, 2 kb, 5 kb, 10 kb, and 20 kb)  
44  
45 64 were constructed according to the standard protocol provided by Illumina (San Diego,  
46  
47 65 CA, USA). Briefly, the DNA sample was randomly broken into fragments by Covaris  
48  
49 66 ultrasonic fragmentation apparatus. The library was prepared following end repair,  
50  
51 67 adding sequence adaptor, purification, and PCR amplification. The mate-pair libraries  
52  
53 68 (2 kb, 5 kb, 10 kb, and 20 kb) and paired-end libraries (250 bp, 450 bp) were all  
54  
55 69 sequenced on Illumina HiSeq4000 platform with paired-end 150 bp. In addition,  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

70 SMRTbell libraries were prepared using either 10-kb or 20-kb preparation protocols.  
71 Briefly, the DNA sample was sheared by Diagenode Megaruptor2 (the Kingdom Of  
72 Belgium), the SMRTbell library was produced by ligating universal hairpin adapters  
73 onto double-stranded DNA fragments. Adapter dimers were efficiently removed using  
74 PacBio's MagBead kit. The final step of the protocol was to remove failed ligation  
75 products through the use of exonucleases. After the exonuclease and AMPure PB  
76 purification steps, sequencing primer was annealed to the SMRTbell templates,  
77 followed by binding of the sequence polymerase to the annealed templates. Subsequent  
78 sequencing was performed on PacBio Sequel instrument with Sequel™ Sequencing Kit  
79 1.2.1 ( Pacific Biosciences of California, USA ). Finally, the 10X Genomics library  
80 was constructed and sequenced with paired-end 150 bp on the Hiseq platform. The  
81 Chromium™ Genome Solution ( 10X Genomics, USA ) massively partitions and  
82 molecularly barcodes DNA using microfluidics, producing sequencing-ready libraries  
83 with >1,000,000 unique barcodes. In total, 463.19 gigabases (Gb) of raw reads were  
84 generated, including 75.72, 70.22, 19.21, 45.71, 28.34, 11.78, 18.01 and 194.20 Gb  
85 from the 250-bp, 450-bp, 2-kb, 5-kb, 10-kb, 20-kb libraries, Pacbio sequencing, and  
86 10X Genomics library respectively. The raw reads were trimmed before used for  
87 subsequent genome assembling. For Illumina HiSeq sequencing, the adaptor sequences,  
88 the reads contain more than 10% ambiguous nucleotides, as well as the reads contains  
89 more than 20% of the low quality nucleotides (quality score less than 5) were all  
90 removed. For PacBio sequencing, the generated polymerase reads were firstly broke at  
91 the adaptor position, the subreads were generated following removing the adaptor  
92 sequences. Then the subreads were filtered by minimum length = 50.

94 Estimation of the genome size and sequencing coverage

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

95 The 17-mer frequency distribution analysis [8] was performed on all the  
96 remaining clean reads to estimate the genome size of the Peruvian scallop using the  
97 following formula: genome size = k-mer number / peak depth. A total number of 6.22  
98  $\times 10^{10}$  k-mers and the peak k-mer depth of 69 was employed to obtain the estimated  
99 genome size at 885.29 Mb (Table 1), and the estimated repeat sequencing ratio  
100 reaches 33.74%.

101

#### 102 *De novo* genome assembly and quality assessment of *A. purpuratus* genome

103 All the pair-end Illumina reads were first assembled into scaffolds using  
104 Platanus\_v1.2.4 [9], and then were applied to fill the gaps by GapCloser\_v1.12-r6  
105 [10]. Subsequently, the Pacbio data were used for additional gap filling by  
106 PBJelly\_v14.1 with default parameters [11], and then all the Illumina reads were  
107 employed for two rounds to correct the genome assembly by Pilon\_v1.18 [12]. After  
108 that, the 10X linked-reads were used to link scaffolds by fragScaff\_140324.1 [13].  
109 Particularly, in order to solve the issue of heterozygosity, in our assembly process, we  
110 chose 19-kmer to draw k-mer distribution histogram, and classified all the kmers into  
111 homozygous kmer and heterozygous kmer according to the coverage depth at first.  
112 Secondly, we utilized 45-kmer to construct the de bruijn figure, combine the bubbles  
113 for heterozygous sites, according to the sequences with longer length and deeper  
114 coverage depth. Then, the reads information and pair-end information were used to  
115 determine the connection between the heterozygous parts, and construct the contigs,  
116 filtered the contigs lacking support. Finally, the heterozygous contigs and  
117 homozygous contigs were distinguished through contigs coverage depth information.  
118 After assembly, the reads from short insert length libraries were mapped to the  
119 assembled genome. And only one peak was observed in the sequencing depth



120 distribution analysis with the average sequencing depth 148.2 X, which is consistent  
121 with the sequencing depth, indicated high quality of the assembled scallop genome.  
122 Finally, a draft genome of 724.78 Mb was assembled (accounting for 81.87% of the  
123 estimated genome size at 885.29 Mb), with a contig N50 size of 80.11 kb and a  
124 scaffold N50 size of 1.02 Mb (Table 1).

125 With this initial assembly, we applied the short insert library reads to map with  
126 the assembled genome using BWA\_0.6.2 software [14] to calculate the mapping ratio  
127 and assess the assembly integrity. In summary, 91.05% short reads were mapped onto  
128 the assembled genome with a coverage of 89.40%, indicating a good reliability of our  
129 genome assembly. CEGMA\_v2.5 (Core Eukaryotic Genes Mapping Approach)  
130 defines a set of conserved protein families that occur in a wide range of eukaryotes,  
131 and presents a mapping procedure to accurately identify their exon-intron structures in  
132 a novel genomic sequence [15]. A protein is classified as complete if the alignment of  
133 the predicted protein to the HMM profile represents at least 70% of the original KOG  
134 domain, otherwise it is classified as partial. Through mapping to the 248 core  
135 eukaryotic genes, a total of 222 genes (89.52%) were identified. BUSCO\_v3  
136 (Benchmarking Universal Single-Copy Orthologs) provides quantitative measures for  
137 the assessment of genome assembly completeness, based on evolutionarily-informed  
138 expectations of gene content from near-universal single-copy orthologs [16]. We  
139 confirmed that 89% of the 843 single-copy genes were identified, indicating a good  
140 integrity of the genome assembly.

141  
142 Repeat sequence analysis of the genome assembly

143 We searched transposable elements (TEs) in the assembled genome through  
144 *ab-initio* and homology based methods. For the former method, we applied

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

145 RepeatModeler\_1.0.4 [17] (the parameter set as '--engine\_db wublast') to build a  
146 specific repeat database. For the latter method, we employed known repeat library  
147 (Rebase) [18] to identify repeats with RepeatMasker\_open-4.0 [19] (the parameter  
148 set as '-a -nolow -no\_is -norna -parallel 3 -e wublast --pvalue 0.0001') and  
149 RepeatProteinMask (the parameter set as '-noLowSimple -pvalue 0.0001 -engine  
150 wublast') [19]. Tandem repeats finder\_4.04 (TRF) was used to find tandem repeats  
151 with the parameters setting as 'Match = 2, Mismatching penalty = 7, Delta = 7, PM =  
152 80, PI = 10, Minscore = 50, MaxPeriod = 2,000' [20]. Finally, we summarized that the  
153 total repeat sequences are 294,496,811 bp, accounting for 40.63% of the assembled  
154 genome, and including 11.46% of tandem repeats, which is consistent with our  
155 above-mentioned estimation (Table 2).

## 157 Gene annotation

### 158 (1) Annotation of protein coding genes

159 The annotation strategy for protein-coding genes integrated *de novo* prediction  
160 with homology and transcriptome data based evidence. Homology sequences from  
161 Pacific oyster (*Crassostrea gigas*), Molluscs (*Lottia gigantea*), Mosquito (*Anopheles*  
162 *gambiae*), Amphioxus (*Branchiostoma floridae*), Nematode (*Caenorhabditis elegans*),  
163 Ascidian (*Ciona intestinalis*), Fruit fly (*Drosophila melanogaster*), Leech (*Helobdella*  
164 *robusta*), Human (*Homo sapiens*), Octopuses (*Octopus bimaculoides*), Sea urchin  
165 (*Strongylocentrotus purpuratus*) were downloaded from Ensemble [21]. The protein  
166 sequences of homology species were aligned to the assembled genome with  
167 TBLASTn (e-value  $\leq 10^{-5}$ ) [22] and predicted gene structures with  
168 GeneWise\_2.4.1 (the parameter set as '-genesf') [23]. The transcriptome data were  
169 generated from adductor muscle, gill, hepatopancreas and mantle on Illumina

170 HiSeq4000 platform. And Tophat\_2.1.1 (the parameter set as '--max-intron-length  
171 500000 -m 2 --library-type fr-unstranded') [24] was utilized to map the transcriptome  
172 data onto our genome assembly and then Cufflinks\_2.1.0 (the parameter set as  
173 '--multi-read-correct') [25] was employed to generated gene model according to the  
174 pair-end relationships and the overlap between aligned reads. The *de novo* prediction  
175 of genes was carried out with four programs: Augustus\_3.0.3 (the parameter set as  
176 '-uniqueGeneId true --noInFrameStop=true --gff3 on --genemodel complete --strand  
177 both') [26], Genscan (the default parameter) [27], GlimmerHMM\_3.0.2 (the  
178 parameter set as '-f -g') [28] and SNAP (the default parameter) [29]. All evidences of  
179 gene model were integrated using EvidenceModeler\_r2012-06-25 (EVM) [29].  
180 Finally, we identified 26,256 protein-coding genes in the Peruvian scallop genome. A  
181 total of 26,513 genes were predicted through the *de novo* method, 19,394 genes were  
182 annotated by RNA transcripts or raw RNA reads, and 15,608 genes were supported by  
183 homolog evidence. The average transcript length, CDS length and intron length were  
184 calculated to be 10,534 bp, 1,418 bp and 1,505 bp respectively (Table 1).

185

## 186 (2) Gene functional annotation

187 Gene functions were predicted from the best BLASTP (e-value  $\leq 10^{-5}$ ) hits in  
188 SwissProt databases [30]. Gene domain annotation was performed by searching the  
189 InterPro database [31]. All genes were aligned against Kyoto Encyclopedia of Genes  
190 and Genomes (KEGG) [32] to identify the best hits for pathways. Gene Ontology (GO)  
191 terms for genes were obtained from the corresponding InterPro entry [33]. Finally,  
192 among these annotated genes, 70.3% encoded proteins showed homology to proteins  
193 in the SwissProt database, 91.1% were identified in the non-redundant (Nr) database,

194 70.4% were identified in the KEGG database, 72.1% were identified in the InterPro,  
195 and a total of 92.1% could be mapped to the functional databases.

196

### 197 (3) *Non-coding RNA annotation*

198 The non-coding RNA genes, including miRNAs, rRNAs, snRNAs and tRNAs,  
199 were identified. The tRNAscan-SE 2.0 software with eukaryote parameters [34] was  
200 employed to predict tRNA genes. The miRNA and snRNA genes in the assembled  
201 genome were extracted by INFERNAL\_1.1.2 software [35] against the Rfam database  
202 [36] with default parameters. Finally, 1,132 miRNAs, 1,664 tRNAs, 41 rRNAs and  
203 220 snRNAs were discovered from the Peruvian scallop genome.

204

### 205 Global gene family classification

206 Protein-coding genes from the Peruvian scallop and other sequenced species,  
207 including Human (*H. sapiens*), Amphioxus (*B. floridae*), Fruit fly (*D. melanogaster*),  
208 Red flour beetle (*T. castaneum*), Nematode (*Caenorhabditis elegans*), brachiopod  
209 (*Lingula anatine*), *Helobdella robusta*, *Capitella teleta*, *Octopus bimaculoides*, *Lottia*  
210 *gigantea*, mollusk (*Aplysia californica*), Pacific Abalone (*Haliotis discus*), Pacific  
211 oyster (*C. gigas*), pearl oyster (*Pinctada fucata*), Yesso scallop (*Patinopecten*  
212 *yessoensis*), and cold seep mussel (*Bathymodiolus platifrons*), Brown mussel  
213 (*Modiolus philippinarum*) were analyzed. All data were downloaded from Ensemble  
214 [21] or NCBI [37]. For each protein-coding gene with alternative splicing isoforms,  
215 we only kept the longest protein sequence as the representative.

216 Gene family analysis was based on the homolog of gene sequences in related  
217 species, which was initially implemented by the alignment of an "all against all"  
218 BLASTP (with a cutoff of 1e-7), and subsequently the alignments with high-scoring

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

219 segment pairs were conjoined for each gene pair by TreeFam\_3.0 [38]. To identify  
220 homologous gene pairs, we required more than 30% coverage of the aligned regions  
221 in both homologous genes. Finally, homologous genes were clustered into gene  
222 families by OrthoMCL-5 [39] with the optimized parameter of '-inflation 1.5'. All  
223 protein-coding genes from the examined 18 genomes were employed to assign gene  
224 families. In total, the protein-coding genes were classified into 45,268 families and  
225 108 strict single-copy orthologs (Figure 2).

226

## 227 Phylogenetic analysis

228 Evolutionary analysis was performed using these single-copy protein-coding  
229 genes from the examined 18 species. Amino acid and nucleotide sequences of the  
230 ortholog genes were aligned by the multiple alignment software MUSCLE with  
231 default parameters [40]. A total number of 108 single-copy ortholog alignments were  
232 concatenated into a super alignment matrix of 242,085 nucleotides. A Maximum  
233 Likelihood method (ML) deduced tree was inferred based on the matrix of nucleotide  
234 sequences using RAxML\_v8.0.19 with default nucleotide substitution  
235 model-PROTGAMMAAUTO [41]. Clade support was assessed using bootstrapping  
236 algorithm in the RAxML package with 100 alignment replicates (Figure 3) [42]. The  
237 constructed phylogenetic tree (Figure 3) indicated that the Peruvian scallop and Yesso  
238 scallop are clustered closely first and then clustered with Oysters and Mussels, which  
239 is in consistent with their putative evolution relationships [43-46].

240

## 241 The estimation of divergence time

242 The species divergence times were inferred with MCMCTree included in PAML  
243 v4.7a [47] with the parameter set as 'burn-in=1,000, sample-number=1000,000,

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

244 sample-frequency=2', and evolutionary analysis was performed using single-copy  
245 protein-coding genes from the 18 examined species. Based on the phylogenetic tree  
246 (Figure 3), the molecular clock was calibrated based on the fossil records according to  
247 previous studies [48-51]. Finally, we estimated that the divergence between the  
248 Peruvian scallop and Yesso scallops happened at 113.6 Mya ago.

249

## 250 **Conclusion**

251 In the present study, we reported the first whole genome sequencing, assembly and  
252 annotation of Peruvian scallop (*A. purpuratus*), an economically important bivalve in  
253 China. The assembled draft genome of 724.78 Mb accounts for 81.87% of the  
254 estimated genome size (885.29 Mb). A total of 26,256 protein-coding genes and 3,057  
255 non-coding RNAs were predicted from the assembly. In the coming future, this  
256 generated genome assembly will provide solid support for deep biological studies.  
257 With availability of these genomic data, subsequent development of genetic markers  
258 for further genetic selection and molecular breeding of scallops could be realized. Our  
259 current genome data will also definitely facilitate the genetic evolutionary history  
260 analysis for the abundant scallops in the world.

261

## 262 **Availability of Data**

263 Supporting raw data have been deposited in NCBI with the project accession  
264 PRJNA418203.

265

## 266 **Acknowledgements**

267 This work was supported by National Natural Science Foundation of China  
268 (31572618 granted to C. Wang and 41676152 granted to X. Liu) and the earmarked

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

269 fund for Shandong Modern Agro-Industry Technology Research System (SDAIT-14)

270 granted to C. Wang.

271

272 **Conflicts of interest**

273 The authors declare that they have no competing interests.

274

275 **Author's contributions**

276 C.W., X.L. and C.L. designed the project. B.M. and G.L. collected the samples and

277 prepared the quality control. C.L., C.W. and X.L. were involved in the data analysis.

278 C.W., X.L., C.L. and Q.S. wrote the manuscript. All authors read and approved the

279 final manuscript.

280

281 **Reference**

282 1. Dall WH. The mollusca and branchiopoda. Report of dredging operation,  
283 Albatros' 1891. Bulletin Mollusca Comparative Zoology, 1909; 37: 147-294.

284 2. Gonzalez ML, Lopez DA, Perez MC, Riquelme VA, Uribe JM and Le PM.  
285 Growth and the scallop, *Argopecten purpuratus* (Lamarck, 1819), in southern  
286 Chile. Aquaculture. 1999; 175 3–4:307-16.

287 3. Genética D, Morfológica Y, Dos E, Del P, Argopecten P, De P, et al. Genetic  
288 and morphological differentiation between two pectinid populations of  
289 *Argopecten purpuratus* from the northern Chilean coast. Estudios  
290 Oceanológicos. 2001; 1:51-60.

291 4. Disalvo LH, Alarcon E, Martinez E and Uribe E. Progress in mass culture of  
292 *Chlamys (Argopecten) purpurata* Lamarck (1819) with notes on its natural  
293 history. Revista Chilena de Historia Natural. 1984, 57:35-45.

294 5. Estabrooks SL. The possible role of telomeres in the short life span of the bay  
295 scallop, *Argopecten irradians irradians* (Lamarck 1819). Journal of Shellfish  
296 Research. 2007; 26 2:307-13.

297 6. Wang C, Liu B, Li J and Liu S. Inter-specific hybridization between  
298 *Argopecten purpuratus* and *Argopecten irradians irradians*. Marine Sciences.  
299 2009; 33 10:84-75.

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 300 7. Wang C, Liu B, Li J, Liu S, Hu L, Fan X, et al. Introduction of the Peruvian  
301 scallop and its hybridization with the bay scallop in China. *Aquaculture*. 2011;  
302 310 3–4:380-7.
  - 303 8. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel  
304 counting of occurrences of k-mers. *Bioinformatics*. 2011; 27 6:764.
  - 305 9. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al.  
306 Efficient *de novo* assembly of highly heterozygous genomes from  
307 whole-genome shotgun short reads. *Genome Research*. 2014; 24 8:1384-95.
  - 308 10. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an  
309 empirically improved memory-efficient short-read *de novo* assembler.  
310 *Gigascience*. 2012; 1 1:18.
  - 311 11. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the gap:  
312 upgrading genomes with Pacific Biosciences RS long-read sequencing  
313 technology. *Plos One*. 2012; 7 11:e47768.
  - 314 12. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al.  
315 Pilon: an integrated tool for comprehensive microbial variant detection and  
316 genome assembly improvement. *Plos One*. 2014; 9 11:e112963.
  - 317 13. Adey A, Kitzman JO, Burton JN, Daza R, Kumar A, Christiansen L, et al. In  
318 vitro, long-range sequence information for *de novo* genome assembly via  
319 transposase contiguity. *Genome Research*. 2014; 24 12:2041.
  - 320 14. Li H and Durbin R. Fast and accurate short read alignment with  
321 Burrows-Wheeler transform. 2009; 25 14:1754-1760..
  - 322 15. Parra G, Bradnam K and Korf I. CEGMA: a pipeline to accurately annotate  
323 core genes in eukaryotic genomes. *Bioinformatics*. 2007; 23 9:1061.
  - 324 16. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM.  
325 BUSCO: assessing genome assembly and annotation completeness with  
326 single-copy orthologs. *Bioinformatics*. 2015; 31 19:3210.
  - 327 17. Grundmann N, Demester L and Makalowski W. TEclass-a tool for automated  
328 classification of unknown eukaryotic transposable elements. *Bioinformatics*.  
329 2009; 25 10:1329.
  - 330 18. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, et al. Repbase Update, a  
331 database of eukaryotic repetitive elements. *Cytogenetic & Genome Research*.  
332 2005 110:462-7.
  - 333 19. Tarailograovac M and Chen N. Using RepeatMasker to identify repetitive  
334 elements in genomic sequences. *Current protocols in bioinformatics*. 2009;  
335 Chapter 4 Unit 4:Unit 4.10.
  - 336 20. Benson G. Tandem repeats finder: a program to analyze DNA sequences.  
337 *Nucleic Acids Research*. 1999; 27 2:573-80.
  - 338 21. Kersey PJ, Allen JE, Allot A, Barba M, Boddu S, Bolt BJ, et al. Ensembl  
339 Genomes 2018: an integrated omics infrastructure for non-vertebrate species.  
340 *Nucleic Acids Res*. 2017; doi:10.1093/nar/gkx1011.



- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- 341 22. Kent WJ. BLAT--the BLAST-like alignment tool. *Genome Research*. 2002;  
342 12 4:656-64.
  - 343 23. Birney E, Clamp M and Durbin R. GeneWise and Genomewise. *Genome*  
344 *Research*. 2004; 14 5:988.
  - 345 24. Trapnell C, Pachter L and Salzberg SL. TopHat: discovering splice junctions  
346 with RNA-Seq. *Bioinformatics*. 2009; 25 9:1105-11.
  - 347 25. Pertea G. Transcript assembly and quantification by RNA-Seq reveals  
348 unannotated transcripts and isoform switching during cell differentiation.  
349 2010.
  - 350 26. Stanke M and Waack S. Gene prediction with a hidden Markov model and a  
351 new intron submodel. *Bioinformatics*. 2003; 19 suppl\_2:215-25.
  - 352 27. Salamov AA and Solovyev VV. Ab initio gene finding in *Drosophila* genomic  
353 DNA. *Genome Research*. 2000; 10 4:516.
  - 354 28. Majoros WH, Pertea M and Salzberg SL. TigrScan and GlimmerHMM: two  
355 open source ab initio eukaryotic gene-finders. *Bioinformatics*. 2004; 20  
356 16:2878-9.
  - 357 29. Korf I. Gene finding in novel genomes. *Bmc Bioinformatics*. 2004; 5 1:59.
  - 358 30. Bairoch A and Apweiler R. The SWISS-PROT protein sequence database and  
359 its supplement TrEMBL in 2000. *Nucleic Acids Research*. 2000; 28 1:45-8.
  - 360 31. Mulder N and Apweiler R. InterPro and InterProScan: tools for protein  
361 sequence classification and comparison. *Methods in Molecular Biology*. 2007;  
362 396:59.
  - 363 32. Kanehisa M and Goto S. KEGG: Kyoto Encyclopedia of genes and genomes.  
364 *Nucleic Acids Research*. 2000; 27 1:29-34.
  - 365 33. Sherlock G. Gene Ontology: tool for the unification of biology. *Canadian*  
366 *Institute of Food Science & Technology Journal*. 2009; 22 4:415.
  - 367 34. Lowe TM and Eddy SR. tRNAscan-SE: a program for improved detection of  
368 transfer RNA genes in genomic sequence. *Nucleic Acids Research*. 1997; 25  
369 5:955-64.
  - 370 35. Nawrocki EP, Kolbe DL and Eddy SR. Infernal 1.0: inference of RNA  
371 alignments. *Bioinformatics*. 2009; 25 10:1335.
  - 372 36. Griffiths-Jones S, Moxon S, Marshall M, Khanna A, Eddy SR and Bateman A.  
373 Rfam: annotating non-coding RNAs in complete genomes. *Nucleic Acids*  
374 *Research*. 2005; 33 Database issue:121-4.
  - 375 37. Sayers EW, Barrett T, Benson DA, Bryant SH, Canese K, Chetvernin V, et al.  
376 Database resources of the national center for biotechnology information. In:  
377 *Haptics Symposium 2010*, pp.199-205.
  - 378 38. Ruan J, Li H, Chen Z, Coghlan A, Coin LJ, Guo Y, et al. TreeFam: 2008  
379 Update. *Nucleic Acids Research*. 2008; 36 Database issue:D735-D40.
  - 380 39. Li L, Stoeckert CJ and Roos DS. OrthoMCL: Identification of ortholog groups  
381 for eukaryotic genomes. *Genome Research*. 2003; 13 9:2178.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

382 40. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and  
383 high throughput. *Nucleic Acids Research*. 2004; 32 5:1792-7.

384 41. Stamatakis A. RAxML Version 8: A tool for Phylogenetic Analysis and  
385 post-analysis of large phylogenies. *Bioinformatics*. 2014; 30 9:1312.

386 42. Stamatakis A, Ott M and Ludwig T. RAxML-OMP: An Efficient Program for  
387 Phylogenetic Inference on SMPs. In: *International Conference on Parallel*  
388 *Computing Technologies 2005*, pp.288-302.

389 43. Sun J, Zhang Y, Xu T, Mu H, Lan Y, Fields CJ, et al. Adaptation to deep-sea  
390 chemosynthetic environments as revealed by mussel genomes. *Nature Ecology*  
391 *& Evolution*. 2017;1 5:121.

392 44. Shi W, Zhang J, Jiao W, Ji L, Xun X, Yan S, et al. Scallop genome provides  
393 insights into evolution of bilaterian karyotype and development. *Nature*  
394 *Ecology & Evolution*. 2017; 1 5:120.

395 45. Kocot KM, Cannon JT, Todt C, Citarella MR, Kohn AB, Meyer A, et al.  
396 Phylogenomics reveals deep molluscan relationships. *Nature*. 2011; 477  
397 7365:452.

398 46. Smith SA, Wilson NG, Goetz FE, Feehery C, Andrade SC, Rouse GW, et al.  
399 Resolving the evolutionary relationships of molluscs with phylogenomic tools.  
400 *Nature*. 2011; 480 7377:364-7.

401 47. Yang Z. PAML: a program package for phylogenetic analysis by maximum  
402 likelihood. *Computer Applications in the Biosciences Cabios*. 1997; 13 5:555.

403 48. Simakov O, Kawashima T, Marlétaz F, Jenkins J, Koyanagi R, Mitros T, et al.  
404 Hemichordate genomes and deuterostome origins. *Nature*. 2015; 527  
405 7579:459-65.

406 49. Benton, M, Donoghue, P & Asher, R. Calibrating and constraining  
407 molecular clocks. in H, S B. & S Kumar (eds), *The timetree of Life*. Oxford  
408 University Press, 2009, pp. 35 - 86.

409 50. Mergl M, Massa D and Plauchut B. Devonian and carboniferous brachiopods  
410 and bivalves of the Djado sub-basin (North Niger, SW Libya). *Journal of the*  
411 *Czech Geological Society*. 2001; 46 3:169-88.

412 51. Erwin DH, Laflamme M, Tweedt SM, Sperling EA, Pisani D and Peterson KJ.  
413 The Cambrian conundrum: early divergence and later ecological success in the  
414 early history of animals. *Science*. 2011; 334 6059:1091-7.

415

416 51. Erwin DH, Laflamme M, Tweedt SM, Sperling EA, Pisani D and Peterson KJ.  
417 The Cambrian conundrum: early divergence and later ecological success in the  
418 early history of animals. *Science*. 2011;334 6059:1091-7.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432

Table 1. Summary of the Peruvian scallop genome assembly and annotation

<b>Genome assembly</b>	<b>Parameter</b>
Contig N50 size (kb)	80.11
Scaffold N50 size (Mb)	1.02
Estimated genome size (Mb)	885.29
Assembled genome size (Mb)	724.78
Genome coverage (×)	303.83
The longest scaffold (bp)	11,125,544
<b>Genome annotation</b>	<b>Parameter</b>
Protein-coding gene number	26,256
Average transcript length (kb)	10.53
Average CDS length (bp)	1,418.29
Average intron length (bp)	1,505.92
Average exon length (bp)	201.09
Average exons per gene	7.05

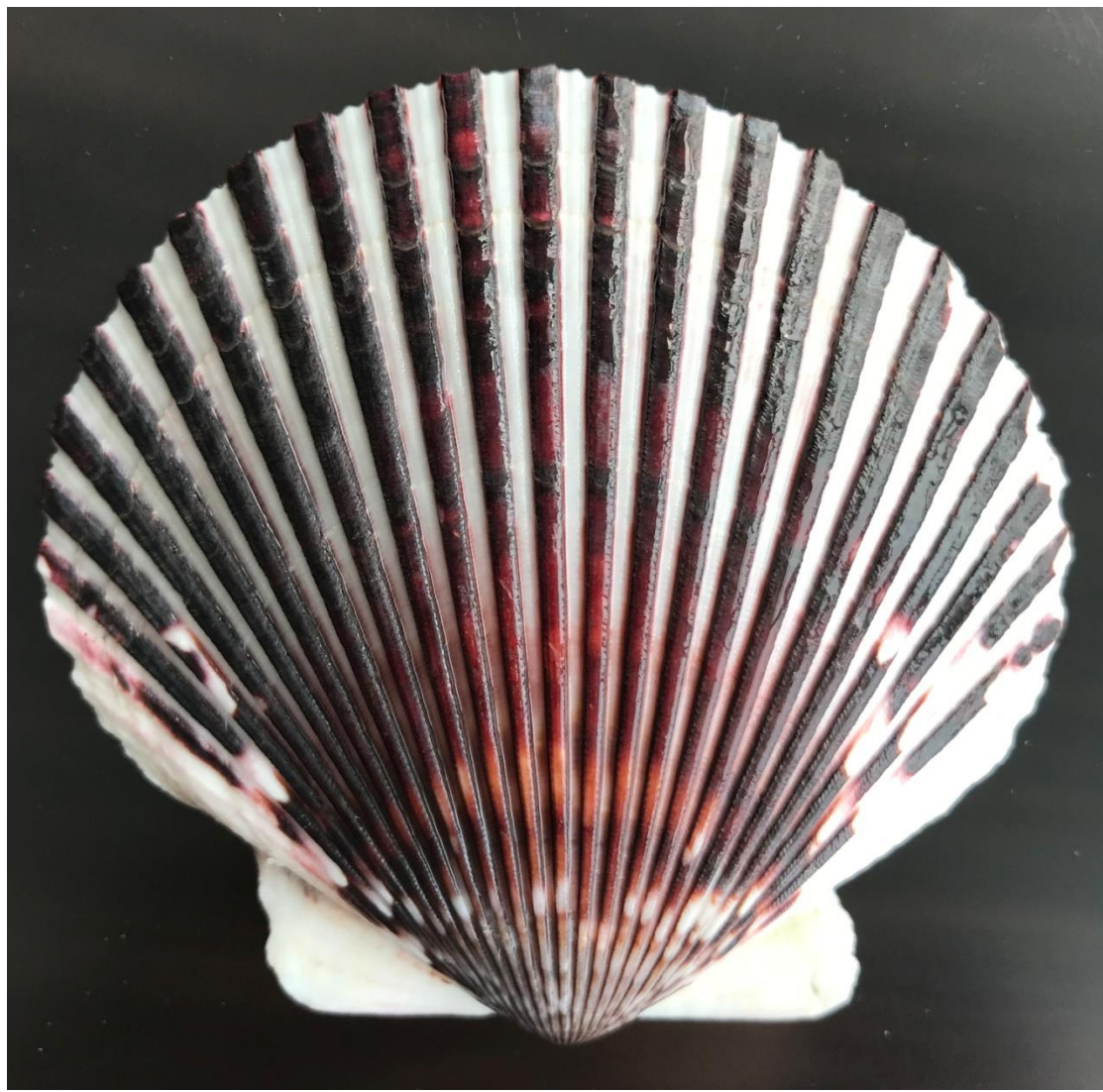
1 433  
 2  
 3  
 4 434  
 5  
 6 435  
 7  
 8  
 9 436  
 10  
 11  
 12 437  
 13  
 14  
 15  
 16  
 17  
 18  
 19  
 20  
 21  
 22  
 23 438  
 24 439  
 25  
 26 440  
 27  
 28  
 29 441  
 30  
 31  
 32 442  
 33  
 34  
 35 443  
 36  
 37 444  
 38  
 39  
 40 445  
 41  
 42  
 43 446  
 44  
 45  
 46 447  
 47  
 48  
 49 448  
 50  
 51  
 52 449  
 53  
 54  
 55 450  
 56  
 57  
 58 451  
 59  
 60 452  
 61  
 62  
 63  
 64  
 65

Table 2. The prediction of repeats elements in the Peruvian scallop genome.

Type	Repeat Size (bp)	% of genome
TRF	83,037,380	11.46
RepeatMasker	237,471,691	32.76
RepeatProteinMask	21,719,425	3.00
Total	294,496,811	40.63

1 453  
2  
3  
4 454  
5  
6 455  
7  
8  
9 456  
10

11  
12 457 **Figure 1** Picture of a representative Peruvian scallop in China.  
13



14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50 458  
51  
52 459  
53  
54  
55 460  
56  
57  
58 461  
59  
60  
61  
62  
63  
64  
65

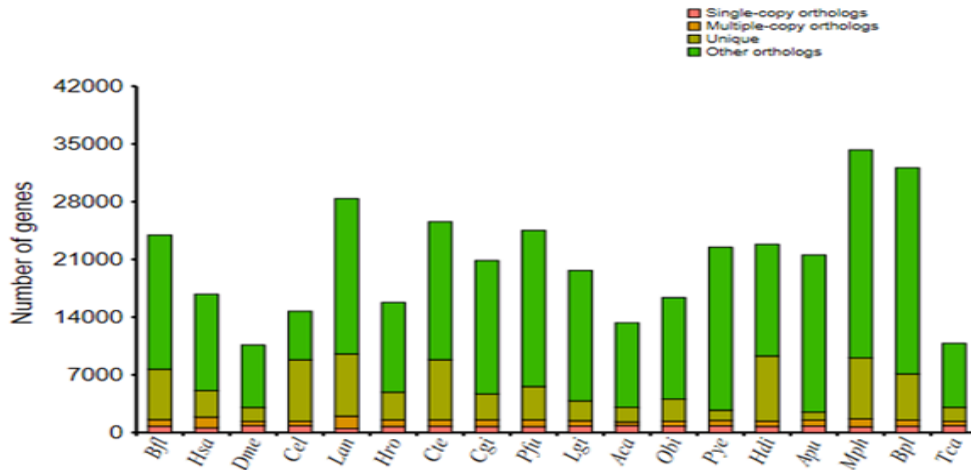
1 462

4 463

7 464

10 465

11 **Figure 2. Distribution of genes in different species.** Abbreviations: Aca, *Aplysia*  
 12 *californica*; Apu, *Argopecten purpuratus*; Bfl, *Branchiostoma floridae*; Bpl,  
 13 *Bathymodiolus platifrons*; Cel, *Caenorhabditis elegans*; Cgi, *Crassostrea gigas*; Cte,  
 14 *Capitella teleta*; Dme, *Drosophila melanogaster*; Has, *Homo sapiens*; Hdi, *Haliotis*  
 15 *discus*; Hro, *Helobdella robusta*; Lan, *Lingula anatine*; Lgi, *Lottia gigantean*; Mph,  
 16 *Modiolus philippinarum*; Obi, *Octopus bimaculoides*; Pfu, *Pinctada fucata*; Tca,  
 17 *Tribolium castaneum*.



38 473

40 474

42 475

45 476

48 477

51 478

54 479

57 480

60 481

1 482

2

3 483

4

5 484

6

7 485

8 486

9 487

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41 488

42

43 489

44

45

46 490

47

48

49 491

50

51 492

52

53 493

54

55

56

57

58

59

60

61

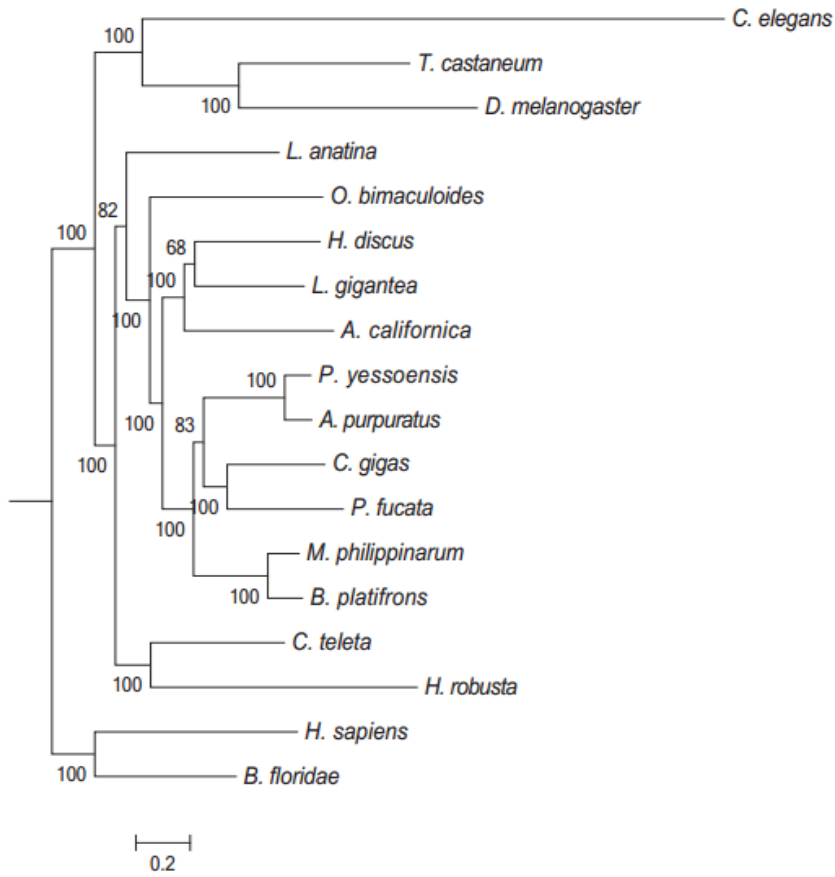
62

63

64

65

**Figure 3. Bootstrap support of phylogenetic tree.** A ML tree was constructed by RAxML based on 108 single-copy protein-coding genes of the related species. The total number of bootstrap was 100.



41 488

42

43 489

44

45

46 490

47

48

49 491

50

51 492

52

53 493

54

55

56

57

58

59

60

61

62

63

64

65



Click here to access/download  
**Supplementary Material**  
RESPONSE TO REVIEWERS.docx

