

Manuscript Number:	GIGA-D-17-00315R2	
Full Title:	Draft genome of the Peruvian scallop <i>Argopecten purpuratus</i>	
Article Type:	Data Note	
Funding Information:	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
Abstract:	<p>Background: The Peruvian scallop, <i>Argopecten purpuratus</i>, is mainly cultured in southern Chile and was 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 of 885.29 Mb), with a contig N50 size of 80.11 kb and a 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 high quality 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>	
Corresponding Author:	Chunde Wang, Ph.D Qingdao Agricultural University CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Qingdao Agricultural University	
Corresponding Author's Secondary Institution:		
First Author:	Chao Li, Ph.D	
First Author Secondary Information:		
Order of Authors:	Chao Li, Ph.D	
	Xiao Liu, D.Sc	
	Bo Liu, D.Sc	
	Bin Ma, MSc	
	Fengqiao Liu	
	Guilong Liu, MSc	
	Qiong Shi, Ph.D	
	Chunde Wang, Ph.D	
Order of Authors Secondary Information:		
Response to Reviewers:	Dear Editor,	

Thank you for the comments on the revised manuscript. We have corrected the manuscript thoroughly and have responded to the Reviewer's comments as attached below.

If you have more questions about the re-revised manuscript, please feel free to let us know. Thanks again for your kindly consideration of our paper.

Best regards,

Chunde and Chao

Reviewer reports:

Reviewer #2: I found all major questions by the reviewers were answered and sufficient information and data were added in the revised manuscript. Now I would like to ask the authors to carefully correct some typos.

There are still uncorrected typos which the reviewers pointed out. Please check the WHOLE sentences again and correct them.

Response: The manuscript has been revised carefully.

Lines 161-164 and Lines 207-212

Description of animal common names are inconsistent.

For example, "Molluscs", "Mosquito"... should be decapitalized.

Response: The names have been corrected.

Some are singular form and others are plural, which should be consistent. i.e. "molluscs", "mosquitoes"... or "mollusca", "mosquito"...

Response: The names have been corrected!

In addition, in line 161 "Molluscs" while in line 210 "mollusk". They should be consistent.

Response: Corrected.

Line 210

"Lottia gigantean" should be "Lottia gigantea"

Response: Corrected.

Line 248

Mya stands for "million years ago." So "Mya ago" should be "Mya."

Response: Corrected.

Additional Information:

Question

Response

Are you submitting this manuscript to a special series or article collection?

No

Experimental design and statistics

Yes

Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our [Minimum Standards Reporting Checklist](#). Information essential to interpreting the data presented should be made available in the figure legends.

Have you included all the information requested in your manuscript?

<p>Resources</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 Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (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> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	Yes

1 **Draft genome of the Peruvian scallop *Argopecten purpuratus***

2 Chao Li¹, Xiao Liu², Bo Liu¹, Bin Ma³, Fengqiao Liu¹, Guilong Liu¹, Qiong Shi⁴,

3 Chunde Wang^{1,*}

4
5 ¹Marine Science and Engineering College, Qingdao Agricultural University, Qingdao
6 266109, China

7 ²Key Laboratory of Experimental Marine Biology, Institute of Oceanology, Chinese
8 Academy of Sciences, Qingdao 266071, China

9 ³Qingdao Oceanwide BioTech Co., Ltd., Qingdao 266101, China

10 ⁴Shenzhen Key Lab of Marine Genomics, Guangdong Provincial Key Lab of
11 Molecular Breeding in Marine Economic Animals, BGI Academy of Marine Sciences,
12 BGI Marine, BGI, Shenzhen 518083, China

13
14 *Correspondence address. Professor Chunde Wang, Marine Science and Engineering
15 College, Qingdao Agricultural University, Qingdao 266109, China (tel:
16 +8613589227997; email: chundewang2007@163.com)

17
18 Email addresses of all authors:

19 Chao Li, leochao@163.com

20 Xiao Liu, liuxiao@qdio.ac.cn

21 Bo Liu, liubomusic@126.com

22 Bin Ma, hereiammabin@163.com

23 Fengqiao Liu, liufengqiao2014@163.com

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 Guilong Liu, 969613442@qq.com

2
3
4 25 Qiong Shi, shiqiong@genomics.cn

5
6 26 **Abstract**

7
8
9 27 **Background:** The Peruvian scallop, *Argopecten purpuratus*, is mainly cultured in
10
11
12 28 southern Chile and was introduced into China in the last century. Unlike other
13
14
15 29 *Argopecten* scallops, the Peruvian scallop normally has a long life span of up to seven
16
17
18 30 to ten years. Therefore, researchers have been employing it to develop hybrid vigor.
19
20
21 31 Here, we performed whole genome sequencing, assembly, and gene annotation of the
22
23
24 32 Peruvian scallop, with an important aim to develop genomic resources for genetic
25
26
27 33 breeding in scallops. **Findings:** A total of 463.19-Gb (Gigabase) raw DNA reads were
28
29
30 34 sequenced. A draft genome assembly of 724.78 Mb was generated (accounting for
31
32
33 35 81.87% of the estimated genome size of 885.29 Mb), with a contig N50 size of 80.11
34
35
36 36 kb and a scaffold N50 size of 1.02 Mb. Repeat sequences were calculated to reach
37
38
39 37 33.74% of the whole genome, and a total of 26,256 protein-coding genes and 3,057
40
41
42 38 non-coding RNAs were predicted from the assembly. **Conclusion:** We generated a
43
44
45 39 high quality draft genome assembly of the Peruvian scallop, which will provide a
46
47
48 40 solid resource for further genetic breeding and for the analysis of the evolutionary
49
50
51 41 history of this economically important scallop.

52
53
54
55 43 **Keywords:** *Argopecten purpuratus*; Peruvian scallop; genome assembly; annotation;
56
57
58 44 gene prediction; phylogenetic analysis

1 45

2
3 46

4
5 47 **Data description**

6
7
8 48

9
10 49 Introduction

11
12
13 50 The Peruvian scallop (*Argopecten purpuratus*), also known as the Chilean
14
15 51 scallop, is a medium-sized bivalve with a wide distribution in Peru and Chile [1]. In
16
17 52 Chile, the cultured scallops reach a commercial size of around 9 cm in shell height
18
19 53 within 14-16 months [2]. It is a relatively stenothermic species as its natural habitat is
20
21 54 largely under the influence of upwelling currents from Antarctica [3]. Unlike other
22
23 55 *Argopecten* scallops, the Peruvian scallop normally has a long life span of up to 7-10
24
25 56 years [4, 5]. This scallop was introduced into China in the late 2000's and had played
26
27 57 an important role in stock improvement of *Argopecten* scallops via inter-specific
28
29 58 hybridization with bay scallops [6, 7].
30
31
32
33
34

35 59

36
37 60 Whole genome sequencing

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

1 70 450 bp) were all sequenced on Illumina HiSeq4000 platform with paired-end 150 bp.
2
3 71 In addition, SMRTbell libraries were prepared using either 10-kb or 20-kb preparation
4
5 72 protocols. Briefly, the DNA sample was sheared by Diagenode Megaruptor2 (the
6
7 73 Kingdom of Belgium), the SMRTbell library was produced by ligating universal
8
9 74 hairpin adapters onto double-stranded DNA fragments. Adapter dimers were efficiently
10
11 75 removed using PacBio's MagBead kit. The final step of the protocol was to remove
12
13 76 failed ligation products through the use of exonucleases. After the exonuclease and
14
15 77 AMPure PB purification steps, sequencing primer was annealed to the SMRTbell
16
17 78 templates, followed by binding of the sequence polymerase to the annealed templates.
18
19 79 Subsequent sequencing was performed on PacBio Sequel instrument with Sequel™
20
21 80 Sequencing Kit 1.2.1 (Pacific Biosciences of California, USA). Finally, the 10X
22
23 81 Genomics library was constructed and sequenced with paired-end 150 bp on the
24
25 82 Illumina Hiseq platform. The Chromium™ Genome Solution (10X Genomics, USA)
26
27 83 massively partitions and molecularly barcodes DNA using microfluidics, producing
28
29 84 sequencing-ready libraries with >1,000,000 unique barcodes. In total, 463.19 gigabases
30
31 85 (Gb) raw reads were generated, including 75.72, 70.22, 19.21, 45.71, 28.34, 11.78,
32
33 86 18.01 and 194.20 Gb from the 250-bp, 450-bp, 2-kb, 5-kb, 10-kb, 20-kb libraries,
34
35 87 Pacbio sequencing library, and 10X Genomics library, respectively. The raw reads
36
37 88 were trimmed before being used for subsequent genome assembly. For Illumina HiSeq
38
39 89 sequencing, the adaptor sequences, the reads containing more than 10% ambiguous
40
41 90 nucleotides, as well as the reads containing more than 20% low quality nucleotides
42
43 91 (quality score less than 5) were all removed. For PacBio sequencing, the generated
44
45 92 polymerase reads were firstly broken at the adaptor positions and the subreads were
46
47 93 generated following removing the adaptor sequences. The subreads were then filtered
48
49 94 by minimum length = 50.
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

95

96 Estimation of the genome size and sequencing coverage

97 The 17-mer frequency distribution analysis [8] was performed on the remaining
98 clean reads to estimate the genome size of the Peruvian scallop using the following
99 formula: genome size = k-mer number / peak depth. Based on a total number of 6.22
100 10^{10} k-mers and a peak k-mer depth of 69, the estimated genome size was calculated
101 to be 885.29 Mb (Table 1) and the estimated repeat sequencing ratio was 33.74%.

102

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

104 All the pair-end Illumina reads were first assembled into scaffolds using
105 Platanus_v1.2.4 (Platanus, RRID:SCR_015531) [9], and the gaps were then filled by
106 GapCloser_v1.12-r6 (GapCloser, RRID:SCR_015026) [10]. Subsequently, the Pacbio
107 data were used for additional gap filling by PBJelly_v14.1 (PBJelly,
108 RRID:SCR_012091) with default parameters [11], and then all the Illumina reads
109 were employed to correct the genome assembly by Pilon_v1.18 (Pilon,
110 RRID:SCR_014731) for two rounds [12]. After that, the 10X linked-reads were used
111 to link scaffolds by fragScaff_140324.1 [13]. Particularly, in order to solve the issue
112 of heterozygosity, in our assembly process, we chose 19-kmer to draw k-mer
113 distribution histogram, and classified all the kmers into homozygous kmer and
114 heterozygous kmer according to the coverage depth. Secondly, we utilized 45-kmer to
115 construct the de bruijn figure and combine the bubbles for heterozygous sites,
116 according to the sequences with longer length and deeper coverage depth. Then the
117 pair-end information was used to determine the connection between the heterozygous
118 parts, and filter the contigs lacking support. Finally, the heterozygous contigs and
119 homozygous contigs were distinguished based on contig coverage depth. After

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

120 assembly, the reads from short insert length libraries were mapped onto the assembled
121 genome. And only one peak was observed in the sequencing depth distribution
122 analysis with the average sequencing depth of $148.2\times$, which is consistent with the
123 sequencing depth, indicating high quality of the assembled scallop genome. Finally, a
124 draft genome of 724.78 Mb was assembled (accounting for 81.87% of the estimated
125 genome size of 885.29 Mb), with a contig N50 size of 80.11 kb and scaffold N50 size
126 of 1.02 Mb (Table 1).

127 With this initial assembly, we mapped the short insert library reads onto the
128 assembled genome using BWA_0.6.2 (BWA, RRID:SCR_010910) software [14] to
129 calculate the mapping ratio and assess the assembly integrity. In summary, 91.05% of
130 the short reads were mapped onto the assembled genome with a coverage of 89.40%,
131 indicating high reliability of genome assembly. CEGMA_v2.5 (Core Eukaryotic
132 Genes Mapping Approach; CEGMA, RRID:SCR_015055) defines a set of conserved
133 protein families that occur in a wide range of eukaryotes, and presents a mapping
134 procedure to accurately identify their exon-intron structures in a novel genomic
135 sequence [15]. A protein is classified as complete if the alignment of the predicted
136 protein to the HMM profile represents at least 70% of the original KOG domain, or
137 otherwise classified as partial. Through mapping to the 248 core eukaryotic genes, a
138 total of 222 genes (89.52%) were identified. BUSCO_v3 (Benchmarking Universal
139 Single-Copy Orthologs; BUSCO, RRID:SCR_015008) provides quantitative
140 measures for the assessment of genome assembly completeness, based on
141 evolutionarily-informed expectations of gene content from near-universal single-copy
142 orthologs [16]. We confirmed that 89% of the 843 single-copy genes were identified,
143 indicating a good integrity of the genome assembly.

144

145 Repeat sequence analysis of the genome assembly

146 We searched transposable elements (TEs) in the assembled genome through
147 *ab-initio* and homology-based methods. For the first method, we applied
148 RepeatModeler_1.0.4 (RepeatModeler, RRID:SCR_015027) [17] (the parameter set
149 as '--engine_db wublast') to build a specific repeat database. For the second method,
150 we employed known repeat library (Repbase) [18] to identify repeats with
151 RepeatMasker_open-4.0 [19] (the parameter set as '-a -nolow -no_is -norna -parallel 3
152 -e wublast --pvalue 0.0001') and RepeatProteinMask (the parameter set as
153 '-noLowSimple -pvalue 0.0001 -engine wublast') [19]. Tandem repeats finder_4.04
154 (TRF) was used to find tandem repeats with the parameters setting as 'Match = 2,
155 Mismatching penalty = 7, Delta = 7, PM = 80, PI = 10, Minscore = 50, MaxPeriod =
156 2,000' [20]. Finally, we summarized that the total repeat sequences are 294,496,811
157 bp, accounting for 40.63% of the assembled genome, and including 11.46% of tandem
158 repeats, which is consistent with our above-mentioned estimation (Table 2).

160 Gene annotation

161 (1) Annotation of protein coding genes

162 The annotation strategy for protein-coding genes integrated *de novo* prediction
163 with homology and transcriptome data based evidence. Homology sequences from
164 African malaria mosquito (*Anopheles gambiae*), Ascidian (*Ciona intestinalis*), Florida
165 lancelet (*Branchiostoma floridae*), Fruit fly (*Drosophila melanogaster*), Human
166 (*Homo sapiens*), Leech (*Helobdella robusta*), Nematode (*Caenorhabditis elegans*),
167 Octopuse (*Octopus bimaculoides*), Owl limpet (*Lottia gigantea*), Pacific oyster
168 (*Crassostrea gigas*), Sea urchin (*Strongylocentrotus purpuratus*) were downloaded
169 from Ensemble [21]. The protein sequences of homology species were aligned to the

170 assembled genome with TBLASTn (e-value $\leq 10^{-5}$) [22], and gene structures were
171 predicted with GeneWise_2.4.1 (GeneWise, RRID:SCR_015054) (the parameter set
172 as '-genesf') [23]. The transcriptome data were generated from adductor muscle,
173 hepatopancreas and mantle on Illumina HiSeq4000 platform. And Tophat_2.1.1 (the
174 parameter set as '--max-intron-length 500000 -m 2 --library-type fr-unstranded') [24]
175 was utilized to map the transcriptome data onto genome assembly and then
176 Cufflinks_2.1.0 (Cufflinks, RRID:SCR_014597), the parameter set as
177 '--multi-read-correct', [25] was employed to generate gene model according to the
178 pair-end relationships and the overlap between aligned reads. The *de novo* prediction
179 of genes was carried out with four programs: Augustus_3.0.3 (Augustus: Gene
180 Prediction, RRID:SCR_008417), the parameter set as '-uniqueGeneId true
181 --noInFrameStop=true --gff3 on --genemodel complete --strand both' [26]; GENSCAN
182 (GENSCAN, RRID:SCR_012902), with default parameter [27]; GlimmerHMM_3.0.2
183 (GlimmerHMM, RRID:SCR_002654), the parameter set as '-f -g' [28]; and SNAP
184 (the default parameter) [29]. All evidences of gene model were integrated using
185 EvidenceModeler_r2012-06-25 (EVM) [29]. Finally, we identified 26,256
186 protein-coding genes in Peruvian scallop genome. In detail, a total of 26,513 genes
187 were predicted through the *de novo* method, 19,394 genes were annotated by RNA
188 transcripts or raw RNA reads, and 15,608 genes were supported by homolog
189 evidences. The average transcript length, CDS length and intron length were 10,534
190 bp, 1,418 bp and 1,505 bp, respectively (Table 1).

191

192 (2) Gene functional annotation

193 Gene functions were predicted from the best BLASTP (e-value $\leq 10^{-5}$) hits in
194 SwissProt databases [30]. Gene domain annotation was performed by searching the

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

195 InterPro (InterPro, RRID:SCR_006695) database [31]. All genes were aligned
196 against Kyoto Encyclopedia of Genes and Genomes (KEGG, RRID:SCR_012773)
197 [32] to identify the best hits for pathways. Gene Ontology (GO) terms for genes were
198 obtained from the corresponding InterPro entry [33]. Finally, among these annotated
199 genes, 70.3% of encoded proteins showed homology to proteins in the SwissProt
200 database, 91.1% were identified in the non-redundant (Nr) database, 70.4% were
201 identified in the KEGG database, 72.1% were identified in the InterPro, and a total of
202 92.1% could be mapped onto the functional databases.

203

204 (3) Non-coding RNA annotation

205 The non-coding RNA genes, including miRNAs, rRNAs, snRNAs and tRNAs,
206 were identified. The tRNAscan-SE_2.0 (tRNAscan-SE, RRID:SCR_010835) software
207 with eukaryote parameters [34] was employed to predict tRNA genes. The miRNA
208 and snRNA genes in the assembled genome were extracted by INFERNAL_1.1.2
209 software [35] against the Rfam (Rfam, RRID:SCR_007891) database [36] with
210 default parameters. Finally, 1,132 miRNAs, 1,664 tRNAs, 41 rRNAs and 220
211 snRNAs were discovered from the Peruvian scallop genome.

212

213 Global gene family classification

214 Protein-coding genes from the Peruvian scallop and other sequenced species,
215 including Brachiopod (*Lingula anatina*), Brown mussel (*Modiolus philippinarum*),
216 California sea hare (*Aplysia californica*), cold seep mussel (*Bathymodiolus platifrons*),
217 Florida lancelet (*B. floridae*), Fruit fly (*D. melanogaster*), Human (*H. sapiens*), Leech
218 (*H. robusta*, *Capitella teleta*), Nematode (*C. elegans*), Octopus (*O. bimaculoides*),
219 Owl limpet (*L. gigantea*), Pacific abalone (*Haliotis discus*), Pacific oyster (*C. gigas*),

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

220 Pearl oyster (*Pinctada fucata*), Red flour beetle (*Tribolium castaneum*), and Yesso
221 scallop (*Patinopecten yessoensis*) were analyzed. All data were downloaded from
222 Ensemble [21] or NCBI [37]. For each protein-coding gene with alternative splicing
223 isoforms, only the longest protein sequence was kept as the representative.

224 Gene family analysis based on the homolog of gene sequences in related species
225 was initially implemented by the alignment of an "all against all" BLASTP (with a
226 cutoff of 1e-7) and subsequently followed by alignments with high-scoring segment
227 pairs conjoined for each gene pair by TreeFam_3.0 [38]. To identify homologous
228 gene pairs, we required more than 30% coverage of the aligned regions in both
229 homologous genes. Finally, homologous genes were clustered into gene families by
230 OrthoMCL-5 (OrthoMCL DB: Ortholog Groups of Protein Sequences,
231 RRID:SCR_007839) [39] with the optimized parameter of '-inflation 1.5'. All
232 protein-coding genes from the examined 18 genomes were employed to assign gene
233 families. In total, the protein-coding genes were classified into 45,268 families and
234 108 strict single-copy orthologs (Figure 2).

235 236 Phylogenetic analysis

237 Evolutionary analysis was performed using these single-copy protein-coding
238 genes from the examined 18 species. Amino acid and nucleotide sequences of the
239 ortholog genes were aligned by the multiple alignment software MUSCLE (MUSCLE,
240 RRID:SCR_011812) with default parameters [40]. A total number of 108 single-copy
241 ortholog alignments were concatenated into a super alignment matrix of 242,085
242 nucleotides. A Maximum Likelihood method (ML) deduced tree was inferred based
243 on the matrix of nucleotide sequences using RAxML_v8.0.19 (RAxML,
244 RRID:SCR_006086) with default nucleotide substitution

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

245 model-PROTGAMMAAUTO [41]. Clade support was assessed using bootstrapping
246 algorithm in the RAxML package with 100 alignment replicates (Figure 3) [42]. The
247 constructed phylogenetic tree (Figure 3) indicated that the Peruvian scallop and Yesso
248 scallop were clustered closely first and then clustered with oysters and mussels, which
249 is in consistent with their putative evolution relationships [43-46].

250

251 The estimation of divergence time

252 The species divergence times were inferred with MCMCTree included in PAML
253 v4.7a (PAML, RRID:SCR_014932) [47] with the parameter set as 'burn-in=1,000,
254 sample-number=1000,000, sample-frequency=2', and evolutionary analysis was
255 performed using single-copy protein-coding genes from the 18 examined species.
256 Based on the phylogenetic tree (Figure 3), the molecular clock was calibrated based
257 on the fossil records according to previous studies [48-51]. Finally, we estimated that
258 the divergence between the Peruvian scallop and Yesso scallop happened at 113.6
259 Mya.

260

261 **Conclusion**

262 In the present study, we report the first whole genome sequencing, assembly and
263 annotation of the Peruvian scallop (*A. purpuratus*), an economically important bivalve
264 in Chile, Peru and China. The assembled draft genome of 724.78 Mb accounts for
265 81.87% of the estimated genome size (885.29 Mb). A total of 26,256 protein-coding
266 genes and 3,057 non-coding RNAs were predicted from the genome assembly. This
267 genome assembly will provide solid support for in-depth biological studies. With the
268 availability of these genomic data, subsequent development of genetic markers for
269 further genetic selection and molecular breeding of scallops could be realized. The

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

270 current genome data will also facilitate genetic analyses of the evolutionary history of
271 the abundant scallops in the world.

272

273 **Availability of Data**

274 Supporting data are available in the *GigaScience* database [52]. Raw data have been
275 deposited in NCBI with the project accession PRJNA418203. BioSample accessions:
276 SAMN08022140 (genome); SAMN08731415 (transcriptome; muscle)
277 SAMN08731411 (transcriptome; mantle); SAMN08731410 (transcriptome;
278 hepatopancreas).

279

280 **Acknowledgements**

281 This work was supported by National Natural Science Foundation of China
282 (31572618 granted to C. Wang and 41676152 granted to X. Liu) and the earmarked
283 fund for Shandong Modern Agro-Industry Technology Research System (SDAIT-14)
284 granted to C. Wang.

285

286 **Conflicts of interest**

287 The authors declare that they have no competing interests.

288

289 **Author's contributions**

290 C.W., X.L. and C.L. designed the project. B.M., F.L. and G.L. collected the samples
291 and prepared the quality control. C.L., C.W. and X.L. were involved in the data
292 analysis. C.W., X.L., C.L. and Q.S. wrote the manuscript. All authors read and
293 approved the final manuscript.

294

295 **Reference**

- 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
- 296 1. Dall WH. The mollusca and branchiopoda. Report of dredging operation, Albatros' 1891. Bulletin Mollusca Comparative Zoology, 1909; 37: 147-294.
 - 297
298 2. Gonzalez ML, Lopez DA, Perez MC, Riquelme VA, Uribe JM and Le PM. Growth and the scallop, *Argopecten purpuratus* (Lamarck, 1819), in southern Chile. Aquaculture. 1999; 175 3–4:307-16.
 - 299
300
301 3. Genética D, Morfológica Y, Dos E, Del P, Argopecten P, De P, et al. Genetic and morphological differentiation between two pectinid populations of *Argopecten purpuratus* from the northern Chilean coast. Estudios Oceanologicos. 2001; 1:51-60.
 - 302
303
304
305 4. Disalvo LH, Alarcon E, Martinez E and Uribe E. Progress in mass culture of *Chlamys (Argopecten) purpurata* Lamarck (1819) with notes on its natural history. Revista Chilena de Historia Natural. 1984, 57:35-45.
 - 306
307
308 5. Estabrooks SL. The possible role of telomeres in the short life span of the bay scallop, *Argopecten irradians irradians* (Lamarck 1819). Journal of Shellfish Research. 2007; 26 2:307-13.
 - 309
310
311 6. Wang C, Liu B, Li J and Liu S. Inter-specific hybridization between *Argopecten purpuratus* and *Argopecten irradians irradians*. Marine Sciences. 2009; 33 10:84-75.
 - 312
313
314 7. Wang C, Liu B, Li J, Liu S, Hu L, Fan X, et al. Introduction of the Peruvian scallop and its hybridization with the bay scallop in China. Aquaculture. 2011; 310 3–4:380-7.
 - 315
316
317 8. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. Bioinformatics. 2011; 27 6:764.
 - 318
319 9. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. Genome Research. 2014; 24 8:1384-95.
 - 320
321
322 10. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-efficient short-read *de novo* assembler. Gigascience. 2012; 1 1:18.
 - 323
324
325 11. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the gap: upgrading genomes with Pacific Biosciences RS long-read sequencing technology. Plos One. 2012; 7 11:e47768.
 - 326
327
328 12. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. Plos One. 2014; 9 11:e112963.
 - 329
330
331 13. Adey A, Kitzman JO, Burton JN, Daza R, Kumar A, Christiansen L, et al. In vitro, long-range sequence information for *de novo* genome assembly via transposase contiguity. Genome Research. 2014; 24 12:2041.
 - 332
333
334 14. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. 2009; 25 14:1754-1760..
 - 335
336 15. Parra G, Bradnam K and Korf I. CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes. Bioinformatics. 2007; 23 9:1061.
 - 337

- 1 338 16. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM.
2 339 BUSCO: assessing genome assembly and annotation completeness with
3 340 single-copy orthologs. *Bioinformatics*. 2015; 31 19:3210.
- 4 341 17. Grundmann N, Demester L and Makalowski W. TEclass-a tool for automated
5 342 classification of unknown eukaryotic transposable elements. *Bioinformatics*.
6 343 2009; 25 10:1329.
- 7 344 18. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, et al. Repbase Update, a
8 345 database of eukaryotic repetitive elements. *Cytogenetic & Genome Research*.
9 346 2005 110:462-7.
- 10 347 19. Tarailogrovac M and Chen N. Using RepeatMasker to identify repetitive
11 348 elements in genomic sequences. *Current protocols in bioinformatics*. 2009;
12 349 Chapter 4 Unit 4:Unit 4.10.
- 13 350 20. Benson G. Tandem repeats finder: a program to analyze DNA sequences.
14 351 *Nucleic Acids Research*. 1999; 27 2:573-80.
- 15 352 21. Kersey PJ, Allen JE, Allot A, Barba M, Boddu S, Bolt BJ, et al. Ensembl
16 353 Genomes 2018: an integrated omics infrastructure for non-vertebrate species.
17 354 *Nucleic Acids Res*. 2017; doi:10.1093/nar/gkx1011.
- 18 355 22. Kent WJ. BLAT--the BLAST-like alignment tool. *Genome Research*. 2002;
19 356 12 4:656-64.
- 20 357 23. Birney E, Clamp M and Durbin R. GeneWise and Genomewise. *Genome*
21 358 *Research*. 2004; 14 5:988.
- 22 359 24. Trapnell C, Pachter L and Salzberg SL. TopHat: discovering splice junctions
23 360 with RNA-Seq. *Bioinformatics*. 2009; 25 9:1105-11.
- 24 361 25. Pertea G. Transcript assembly and quantification by RNA-Seq reveals
25 362 unannotated transcripts and isoform switching during cell differentiation.
26 363 2010.
- 27 364 26. Stanke M and Waack S. Gene prediction with a hidden Markov model and a
28 365 new intron submodel. *Bioinformatics*. 2003; 19 suppl_2:215-25.
- 29 366 27. Salamov AA and Solovyev VV. Ab initio gene finding in *Drosophila* genomic
30 367 DNA. *Genome Research*. 2000; 10 4:516.
- 31 368 28. Majoros WH, Pertea M and Salzberg SL. TigrScan and GlimmerHMM: two
32 369 open source ab initio eukaryotic gene-finders. *Bioinformatics*. 2004; 20
33 370 16:2878-9.
- 34 371 29. Korf I. Gene finding in novel genomes. *Bmc Bioinformatics*. 2004; 5 1:59.
- 35 372 30. Bairoch A and Apweiler R. The SWISS-PROT protein sequence database and
36 373 its supplement TrEMBL in 2000. *Nucleic Acids Research*. 2000; 28 1:45-8.
- 37 374 31. Mulder N and Apweiler R. InterPro and InterProScan: tools for protein
38 375 sequence classification and comparison. *Methods in Molecular Biology*. 2007;
39 376 396:59.
- 40 377 32. Kanehisa M and Goto S. KEGG: Kyoto Encyclopedia of genes and genomes.
41 378 *Nucleic Acids Research*. 2000; 27 1:29-34.
- 42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1 379 33. Sherlock G. Gene Ontology: tool for the unification of biology. Canadian
2 380 Institute of Food Science & Technology Journal. 2009; 22 4:415.
- 3 381 34. Lowe TM and Eddy SR. tRNAscan-SE: a program for improved detection of
4 382 transfer RNA genes in genomic sequence. Nucleic Acids Research. 1997; 25
5 383 5:955-64.
- 6 384 35. Nawrocki EP, Kolbe DL and Eddy SR. Infernal 1.0: inference of RNA
7 385 alignments. Bioinformatics. 2009; 25 10:1335.
- 8 386 36. Griffiths-Jones S, Moxon S, Marshall M, Khanna A, Eddy SR and Bateman A.
9 387 Rfam: annotating non-coding RNAs in complete genomes. Nucleic Acids
10 388 Research. 2005; 33 Database issue:121-4.
- 11 389 37. Sayers EW, Barrett T, Benson DA, Bryant SH, Canese K, Chetvernin V, et al.
12 390 Database resources of the national center for biotechnology information. In:
13 391 *Haptics Symposium 2010*, pp.199-205.
- 14 392 38. Ruan J, Li H, Chen Z, Coghlan A, Coin LJ, Guo Y, et al. TreeFam: 2008
15 393 Update. Nucleic Acids Research. 2008; 36 Database issue:D735-D40.
- 16 394 39. Li L, Stoeckert CJ and Roos DS. OrthoMCL: Identification of ortholog groups
17 395 for eukaryotic genomes. Genome Research. 2003; 13 9:2178.
- 18 396 40. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and
19 397 high throughput. Nucleic Acids Research. 2004; 32 5:1792-7.
- 20 398 41. Stamatakis A. RAxML Version 8: A tool for Phylogenetic Analysis and
21 399 post-analysis of large phylogenies. Bioinformatics. 2014; 30 9:1312.
- 22 400 42. Stamatakis A, Ott M and Ludwig T. RAxML-OMP: An Efficient Program for
23 401 Phylogenetic Inference on SMPs. In: *International Conference on Parallel
24 402 Computing Technologies 2005*, pp.288-302.
- 25 403 43. Sun J, Zhang Y, Xu T, Mu H, Lan Y, Fields CJ, et al. Adaptation to deep-sea
26 404 chemosynthetic environments as revealed by mussel genomes. Nature Ecology
27 405 & Evolution. 2017;1 5:121.
- 28 406 44. Shi W, Zhang J, Jiao W, Ji L, Xun X, Yan S, et al. Scallop genome provides
29 407 insights into evolution of bilaterian karyotype and development. Nature
30 408 Ecology & Evolution. 2017; 1 5:120.
- 31 409 45. Kocot KM, Cannon JT, Todt C, Citarella MR, Kohn AB, Meyer A, et al.
32 410 Phylogenomics reveals deep molluscan relationships. Nature. 2011; 477
33 411 7365:452.
- 34 412 46. Smith SA, Wilson NG, Goetz FE, Feehery C, Andrade SC, Rouse GW, et al.
35 413 Resolving the evolutionary relationships of molluscs with phylogenomic tools.
36 414 Nature. 2011; 480 7377:364-7.
- 37 415 47. Yang Z. PAML: a program package for phylogenetic analysis by maximum
38 416 likelihood. Computer Applications in the Biosciences Cabios. 1997; 13 5:555.
- 39 417 48. Simakov O, Kawashima T, Marlétaz F, Jenkins J, Koyanagi R, Mitros T, et al.
40 418 Hemichordate genomes and deuterostome origins. Nature. 2015; 527
41 419 7579:459-65.
- 42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 420 49. Benton, M, Donoghue, P & Asher, R. Calibrating and constraining
2 421 molecular clocks. in H, S B. & S Kumar (eds), *The timetree of Life*. Oxford
3 422 University Press, 2009, pp. 35 - 86.

4 423 50. Mergl M, Massa D and Plauchut B. Devonian and carboniferous brachiopods
5 424 and bivalves of the Djado sub-basin (North Niger, SW Libya). Journal of the
6 425 Czech Geological Society. 2001; 46 3:169-88.

7 426 51. Erwin DH, Laflamme M, Tweedt SM, Sperling EA, Pisani D and Peterson KJ.
8 427 The Cambrian conundrum: early divergence and later ecological success in the
9 428 early history of animals. Science. 2011; 334 6059:1091-7.

10 429 51. Erwin DH, Laflamme M, Tweedt SM, Sperling EA, Pisani D and Peterson KJ.
11 430 The Cambrian conundrum: early divergence and later ecological success in the
12 431 early history of animals. Science. 2011;334 6059:1091-7.

13 432 52. Li C, Liu X, Liu B, Ma B, Liu F, Liu G, Shi Q, Wang C. Supporting data for
14 433 "Draft genome of the Peruvian scallop *Argopecten purpuratus*." GigaScience
15 434 Database 2018. <http://dx.doi.org/10.5524/100419>
16 435
17
18
19
20
21
22
23
24
25 436
26
27
28 437
29
30
31 438
32
33
34 439
35
36
37 440
38
39
40 441
41
42 442
43
44
45 443
46
47
48 444
49
50
51 445
52
53
54 446
55
56
57 447
58
59
60 448
61
62
63
64
65

1 449

2
3
4 450 Table 1. Summary of the Peruvian scallop genome assembly and annotation.

Genome assembly	Parameter
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 (x)	303.83
The longest scaffold (bp)	11,125,544

Genome annotation	Parameter
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

36
37 451

38
39
40 452

41
42
43 453

44
45 454

46
47
48 455

49
50
51 456

52
53
54 457

55
56
57 458

58
59
60 459

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

460

461

462
463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

Table 2. The prediction of repeat 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

1 480

2
3
4 481 **Figure 1** Picture of a representative Peruvian scallop in China.



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
41 482

42
43 483

44
45
46 484

47
48
49 485

50
51
52 486

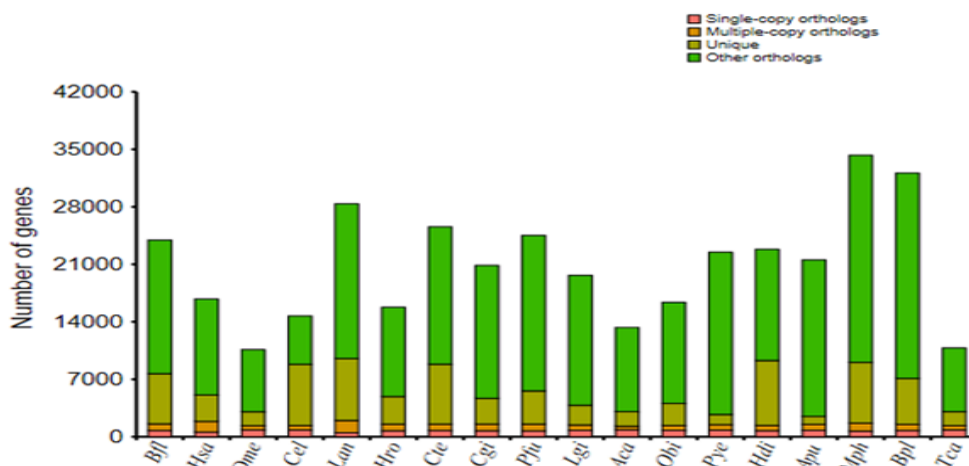
53
54
55 487

56
57
58 488

59
60
61
62
63
64
65

489

490 **Figure 2. Distribution of genes in different species.** Abbreviations: Aca, *Aplysia*
491 *californica*; Apu, *Argopecten purpuratus*; Bfl, *Branchiostoma floridae*; Bpl,
492 *Bathymodiolus platifrons*; Cel, *Caenorhabditis elegans*; Cgi, *Crassostrea gigas*; Cte,
493 *Capitella teleta*; Dme, *Drosophila melanogaster*; Hsa, *Homo sapiens*; Hdi, *Haliotis*
494 *discus*; Hro, *Helobdella robusta*; Lan, *Lingula anatina*; Lgi, *Lottia gigantea*; Mph,
495 *Modiolus philippinarum*; Obi, *Octopus bimaculoides*; Pfu, *Pinctada fucata*; Pye,
496 *Patinopecten yessoensis*; Tca, *Tribolium castaneum*.



497

498

499

500

501

502

503

504

505

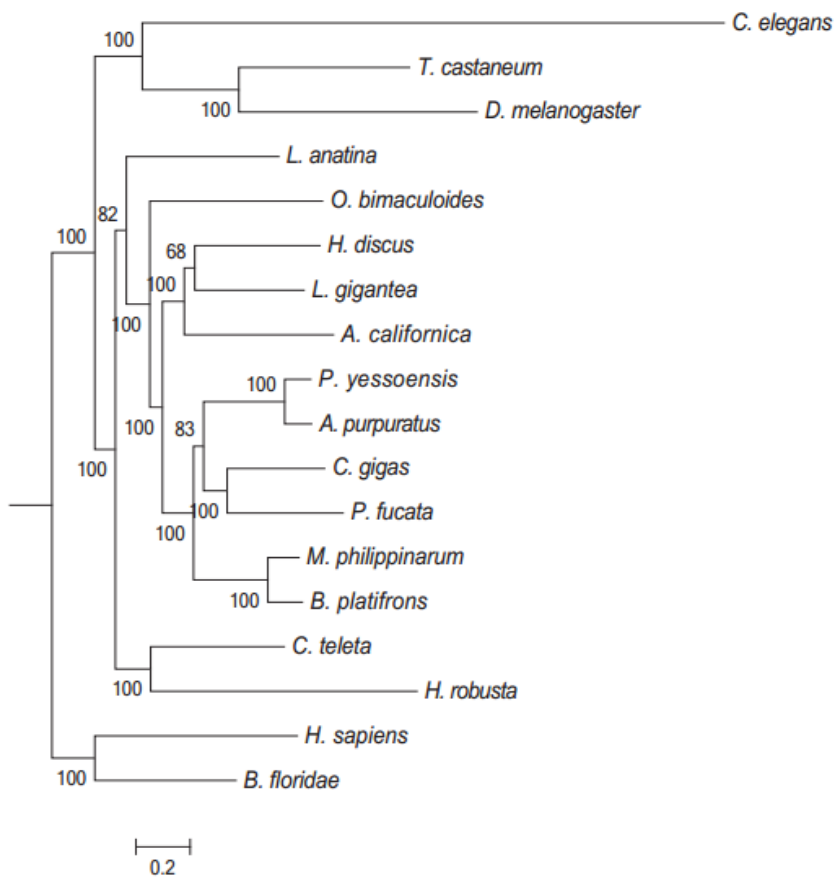
506

507

508

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

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



512

513

514

515

516

517

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



