

## A hybrid-hierarchical genome assembly strategy to sequence the invasive golden mussel *Limnoperna fortunei* --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-17-00124R1									
<b>Full Title:</b>	A hybrid-hierarchical genome assembly strategy to sequence the invasive golden mussel <i>Limnoperna fortunei</i>									
<b>Article Type:</b>	Data Note									
<b>Funding Information:</b>	<table border="1"> <tr> <td>CAPES (PVE (71/2013))</td> <td>Dr Mauro F Rebelo</td> </tr> <tr> <td>FAPERJ (APQ1 (2014))</td> <td>Dr Mauro F Rebelo</td> </tr> <tr> <td>FAPERJ/DFG (FAPERJ/DFG (39/2014))</td> <td>Dr Mauro F Rebelo</td> </tr> <tr> <td>Crowdfunding (www.catarse.me/genoma)</td> <td>Dr Mauro F Rebelo</td> </tr> </table>	CAPES (PVE (71/2013))	Dr Mauro F Rebelo	FAPERJ (APQ1 (2014))	Dr Mauro F Rebelo	FAPERJ/DFG (FAPERJ/DFG (39/2014))	Dr Mauro F Rebelo	Crowdfunding (www.catarse.me/genoma)	Dr Mauro F Rebelo	
CAPES (PVE (71/2013))	Dr Mauro F Rebelo									
FAPERJ (APQ1 (2014))	Dr Mauro F Rebelo									
FAPERJ/DFG (FAPERJ/DFG (39/2014))	Dr Mauro F Rebelo									
Crowdfunding (www.catarse.me/genoma)	Dr Mauro F Rebelo									
<b>Abstract:</b>	<p>Background: For more than 25 years, the golden mussel <i>Limnoperna fortunei</i> has aggressively invaded South American freshwaters, having travelled more than 5,000 km upstream across five countries. Along the way, the golden mussel has outcompeted native species and economically harmed aquaculture, hydroelectric powers, and ship transit. We have sequenced the complete genome of the golden mussel to understand the molecular basis of its invasiveness and search for ways to control it. Findings: We assembled the 1.6 Gb genome into 20548 scaffolds with an N50 length of 312 Kb using a hybrid and hierarchical assembly strategy from short and long DNA reads and transcriptomes. A total of 60717 coding genes were inferred from a customized transcriptome-trained AUGUSTUS run. We also compared predicted protein sets with those of complete molluscan genomes, revealing an exacerbation of protein-binding domains in <i>L. fortunei</i>. Conclusions: We built one of the best bivalve genome assemblies available using a cost-effective approach using Illumina pair-end, mate pair, and PacBio long reads. We expect that the continuous and careful annotation of <i>L. fortunei</i>'s genome will contribute to the investigation of bivalve genetics, evolution, and invasiveness, as well as to the development of biotechnological tools for aquatic pest control.</p>									
<b>Corresponding Author:</b>	Marcela Uliano da Silva, Ph.D Universidade Federal do Rio de Janeiro Rio de Janeiro, RJ BRAZIL									
<b>Corresponding Author Secondary Information:</b>										
<b>Corresponding Author's Institution:</b>	Universidade Federal do Rio de Janeiro									
<b>Corresponding Author's Secondary Institution:</b>										
<b>First Author:</b>	Marcela Uliano da Silva, Ph.D									
<b>First Author Secondary Information:</b>										
<b>Order of Authors:</b>	<table border="1"> <tr> <td>Marcela Uliano da Silva, Ph.D</td> </tr> <tr> <td>Francesco Dondero, Ph.D</td> </tr> <tr> <td>Thomas D Otto, Ph.D</td> </tr> <tr> <td>Igor R Costa, Msc</td> </tr> <tr> <td>Nicholas CB Lima, Ph.D</td> </tr> <tr> <td>Juliana A Americo, Ph.D</td> </tr> <tr> <td>Camila J Mazzone, Ph.D</td> </tr> <tr> <td></td> </tr> </table>		Marcela Uliano da Silva, Ph.D	Francesco Dondero, Ph.D	Thomas D Otto, Ph.D	Igor R Costa, Msc	Nicholas CB Lima, Ph.D	Juliana A Americo, Ph.D	Camila J Mazzone, Ph.D	
Marcela Uliano da Silva, Ph.D										
Francesco Dondero, Ph.D										
Thomas D Otto, Ph.D										
Igor R Costa, Msc										
Nicholas CB Lima, Ph.D										
Juliana A Americo, Ph.D										
Camila J Mazzone, Ph.D										

	Francisco Prosdocimi, Ph.D
	Mauro F Rebelo, Ph.D
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>We thank the reviewers for their attentive read of the manuscript and for suggesting revisions that have increased the overall quality of the data presentation and of the manuscript. Please find bellow each reviewers comment and the answers to them:</p> <p>Reviewers 1:</p> <p>Reviewer 1 - Line 49: could the authors provide an extended background to the readers about the arrival of this invasive species in South America?</p> <p>Response: Yes, the extended background was provided and it's situated in lines 53-55 in the new submission. It is as follows: "... Research suggests that <i>L. fortunei</i> was introduced in South America through ballast water of ships coming from Hong Kong or Korea [2]. It was found for the first time in the estuary of the La Plata River in 1991 [1]."</p> <p>Reviewer 1 - Line 66: it is maybe better to specify here "freshwater bivalves". Indeed, many other species could be considered as "invasive" in the marine environment, including <i>Mytilus</i> spp.</p> <p>Response: "Freshwater" was added at line 72.</p> <p>Reviewer 1 - Line 76: Also, <i>L. fortunei</i> is a mytiloid and other mussel species are known to display an exceptional tolerance to biotic and abiotic contamination, with remarkable capabilities of accumulation and metabolization of toxicants. It is possible that golden mussels share some of these features with marine mussels.</p> <p>Response: It's true. But we kept the introduction as it was in order to keep it concise and cohesive.</p> <p>Reviewer 1 - *Lines 96-97: The choice to use three mussels for DNA extraction and sequencing is unclear (unless this is a typo related to the use of 3 mussels for RNA extraction). Why did the authors choose to use this non-standard procedure? Was the genomic DNA extracted from three different specimens pooled in equimolar quantities and used for sequencing? Usually, as heterozygosity might represent a considerable issue, it is desirable to use a single specimen as a reference for genome assembly.</p> <p>Response: The idea was to sequence only one specimen. But it was not possible due to (i) Illumina DNA library preparation unanticipated problems and (ii) the amount of DNA necessary for PacBio sequencing. The sequencing facility responsible for producing Illumina pair-end and mate pair reads (UNESP) failed to produce the mate pairs in their first attempt, and they asked for more DNA to repeat the library preparation. As we did not have more tissue from the first specimen, we needed to extract more from a second specimen. After that, as we notice the use of only Illumina would not allow us to produce a contiguous high-quality genome, we decided to sequence PacBio. PacBio libraries need a substantial amount of high-molecular-weight-DNA, and to meet this requirements we needed to extract DNA from a third specimen.</p> <p>To clarify the use of 3 specimens for the construction of the 3 sequencing libraries, a small complement was added to the sentence in line 103-105. It's as follows "... For the genome assembly, a total of 3 individuals were sampled for DNA extraction from gills and to produce the three types of DNA libraries used in this study."</p> <p>Reviewer 1 - Lines 137-138: Please indicate what the two colors in figure 1 correspond to (I guess to two different k-mer length, but this is not specified neither in the figure itself, nor in its caption. Also, the relative size of the heterozygous peak compared to the homozygous one is particularly remarkable and indicates an extremely high heterozygosity rate, which the authors could estimate and report. This could be linked easily with the subsequent paragraph and the difficulties in assembling such a highly</p>

heterozygous genome using short reads only. Please note that these issues have been also encountered by Murgarella and colleagues in the draft assembly of the *M. galloprovincialis* genome.

Response: We have added the legend on the figures representing the colors. Red represented a the distribution of kmers size 31 and black represented the kmers of size 25. Also, we have estimated the heterozygosity rate of *L. fortunei* genome to be 2.07%, and we have included this information and some comments between the lines 150-152 It is as follows: "The rate of heterozygosity was estimated to be 2.07% and it was calculated as described by Vij et al. (2016) [18], using as input data the 25-kmer distribution plot for reads from one unique specimen".

And also we did some editings in lines 185-190 . It is as follows "...One main challenge of assembling bivalve genomes lies in the high heterozygosity and amount of repetitive elements these organisms present: (i) the mussels *L. fortunei* and *Modiolus philippinarum* and the oyster *Crassostrea gigas* genomes were estimated to have heterozygosity rates of 2.07%, 2.02 % 1.95% respectively, which is substantially higher than other animal genomes [29], and (ii) repetitive elements correspond to at least 30% of the genomes of all studied bivalves so far (Table 3) [28, 29, 30, 31, 33, 34, 35 ]. "

Reviewer 1 - \*Table 5 and Figure 3 would benefit from the inclusion of a few recently released genomes of other bivalves. Specifically, a much improved version of the *Pinctada fucata* genome has just been released on Gigascience (the authors could not have access to this resource at the time of writing their manuscript): <https://academic.oup.com/gigascience/article/4034775/The-pearl-oyster-Pinctada-fucata-martensii-genome?searchresult=1>.

At the same time, the genome of the pectinoid *Mizuhopecten yessoensis* has also been released (data is available at <http://mgb.ouc.edu.cn/pydatabase/download.php>).

The genome of the veneroid clam *Ruditapes philippinarum* is also now available:

<https://academic.oup.com/gbe/article-lookup/doi/10.1093/gbe/evx096>

In this case, while sequence data is not publicly available yet, the authors are willing to share their data upon request.

Response: The 3 new bivalve genomes (*P. fucata*, *M. yessoensis* and *R. philippinaum*) were included in all the comparative analysis of this paper: in Table 3 and Figures 3 and 4. The previous *P. fucata* data was replaced, and now comparisons were done with the new assembly presented by Du et al (<https://doi.org/10.1093/gigascience/gix059>). Table S3 was updated accordingly. And also line 272.

Reviewer 1 - Line 235: "these genomes" should be "these transcriptomes"

Response: It was corrected. Line 234.

Reviewer 1 - Line 251: the authors could add a brief comment about the 58% rate of gene whose expression could be confirmed, stating that this is a reasonable and even expected result, based on the absence of libraries gathered from developmental stages, some adult tissues (i.e. hemocytes) and mussels subjected to different stress (so that inducible gene products might be absent).

Response: The comment was introduced in line 250-255: It is as follows "...Of those, 58% had transcriptional evidence based on RNA Illumina reads (Table S2) re-mapping, rate that was expected since our RNA-Seq libraries were constructed only for 4 tissues of adult golden mussel specimens without any environmental stresses induction (Table 2). Therefore, these libraries lack transcripts for developmental stages, for some other cell types (i.e. hemocytes) and stress-inducible genes. Finally, 67% of the gene models were annotated by homology searches against Uniprot or NCBI NR (Table 6)."

Reviewer 1 - Lines 27-273: "five mussels" should be "five bivalves". Also, this data could be updated using the newly released bivalve genomes I have listed above.

Response: This was corrected and the information, Supl Table S3, and Figure 3 were updated with the new species included in the analysis. Lines 275.

Reviewer 1 - \*Line 276: "reconstruct phylogeny" needs to be detailed. What strategy was used (Bayesian, ML, NJ?), what model of molecular evolution, what software? Are the support values displayed in the tree posterior probabilities or bootstrap values?

Response: The methods used were more detailed in lines 277-282. Also, the updated phylogeny was performed including the new data for the *P. fucata* genome, replacing the old one used, and also including the new data recommended by the review for *R. phillapirum* and *P. yeoensis*. It is as follows: " These sequences were used to reconstruct a phylogeny: the single-copy orthologs sequences were concatenated and aligned with CLUSTALW [45] with a resulting alignment of 30755 sites in length (Figure 3B). ProtTest 3.4.2 [46] was used to estimate the best fitting substitution model, which was VT [47]. With this alignment and model we reconstructed the phylogeny using PhyML [48] and 100 bootstrap repetition, the resulting tree is shown on Figure 3B."

Reviewer 1 - \*Line 301: TIR domains do not necessarily belong to TLRs. More than half of bivalve TIR-DC proteins are indeed intracellular receptors of unknown function (but which are still likely involved in intracellular immune signaling (see Gerdol et al, DCI 2017). The interpretation of Figure S2 and the discussion contained in lines 303-309 is therefore quite difficult to be evaluated without knowing whether only proteins containing LRRs+TIR or all those containing TIR domains (with and without LRRs) were taken into account. Furthermore, BLAST is not overly useful, by itself, to classify these proteins, as it has been previously demonstrated. Considering the complexity of this topic and the fact that this goes probably beyond the scopes of this manuscript, the authors could simplify tis section by reporting and expanded complement of TIR-DC proteins and DEATH-domain containing proteins of different nature which, accordingly to the know functions of these domain and existing literature data, are likely to be involved in immune signaling. Overall the expansion of these gene families might suggest an improved resistance to infections. It is however equally curious that other immune-related gene families (e.g. FREPs and C1qDC) seem to be somewhat contracted in figure 4.

Response: Having found LRRs and TIR in the list of over-represented PFAM we looked for TLRs in Blast results, since it was logical to find many of them. However, we were completely aware that not all those Blast hits could represent a genuine TLR, since Blast is heuristically biased towards short High Scoring Pairs (HSP) that could be tagged only to a TIR domain. We, therefore, used SMART (Simple Modular Architecture Research Tool, see [http://smart.embl-heidelberg.de/help/smart\\_about.shtml](http://smart.embl-heidelberg.de/help/smart_about.shtml)) to analyze all Blast TLR hits for their modular domain architectures. Only those sequences showing a prototypical TLR architecture were further considered, i.e. N-terminal extracellular leucine-rich repeat (LRR) motifs including either a single or multiple cysteine cluster domain, a C-terminal TIR domain spaced by a single transmembrane-spanning domain (Leulier & Lemaitre, 2008). We know this analysis is not conclusive but TLR expansions in lophotrochozoa were not known until a few years ago when it has been demonstrated in anellida. This finding can contribute to stimulate TLR evolutionary studies. We added some details of the analysis in the body text to explain that those TLR we considered are representative of genuine TLRs.

We have changed a few sentences in the manuscript accordingly. Lines 319-325: It is as follows: "Overall, the expansion of these gene families might suggest an improved resistance to infections. It is, however, equally curious that other immune-related gene families such as Fribinogen\_C and C1q seem to be contracted (Supplementary Table S5). This feature may depend on the evolutionary-driven, yet random, fate of the *L. fortunei* genome and consequence of different specific duplicate genes in other species. Also, other protein families involved in toxin metabolism, especially glutathione based processes and sulfotransferases are clearly contracted (Table S5)."

Reviewer 1 - Line 555: bellow -> below  
Response: Thank you, it was corrected. Line 611.

Reviewer 1 - \*In Figure 4 legend, it is specified that transposable elements were taken into account. I guess that, depending on the annotation pipeline followed by the different genome sequencing projects these might have been either masked or not, thereby being often excluded from the final protein set. While the heat map seems to show that TEs are, in general, extremely expanded in Limnoperna, I would be very careful about this claim. This also applies to Table S4. Considering the very high number of gene predictions corresponding to TEs in Limnoperna a particular attention should be also posed into the calculations of under-representation of domains, as these were made based on relative abundance, which would be de facto lowered in Limnoperna if TEs have been masked in the other molluscan genomes.

Response: We agree with this comment, and it was, in fact, a relevant debate among us if we should include or not such retro-domains in the analysis. However, as it seems that such sequences can have a central biological role in shaping some *L. fortunei* genomic features (and maybe physiological ones), we decided to show them even knowing that in other genome studies they might have been kept out or not considered with attention. Indeed, some genomes we used for the new comparison presented in this revised ms, did include TEs in their annotation analysis, e.g. *Ruditapes philippinarum*, *Haliotis discus*, *Modiolus philippinarum* (See Table 5 of the revised ms). The golden mussel genome always outperformed these numbers. However, we tested how considering TE elements in our PFAM analysis might have biased the down-represented features. The reviewer comment has been very appropriate since it can happen and we were not aware of that. Nevertheless, we are confident of the genuinity of our analysis and results. In fact, we made some trials considering a lower total PFAM count value for frequency normalization in other mollusc genomes. When we re-normalized PFAM frequencies at 5% or 10% less counts than before, about 25% and 50% PFAMs are excluded from the original list. Considering that (i) we have estimated about 2500 PFAM countss (nearly 6%); (ii) some other annotations included in the analysis are actually using PFAM associated to TEs; (iii) we used the most conservative false discovery rate procedure, i.e. Bonferroni's; we can conclude that excluding TE from this analysis can be more detrimental than beneficial to the correct functional annotation of the golden mussel genome.

Reviewer 1 - Table S3: "4 other mollusk" -> please correct 4  
Response: Table S3 was updated.

Reviewer #2 (Kevin Kocot): Specific comments:  
There are too many very short paragraphs. A paragraph should always have at least two sentences. The paragraph spanning lines 58-65 covers two disparate topics and the introduction of the text may need to be reorganized.

Response: we tried to avoid the short paragraph as much as possible. For example, adding a short paragraph to the last line of Table 3, and then deleting it from the manuscript.

Reviewer 2: Why were multiple individuals used?

Response: The idea was to sequence only one specimen. But it was not possible due to (i) Illumina DNA library preparation unanticipated problems and (ii) the amount of DNA necessary for PacBio sequencing. The sequencing facility responsible for producing Illumina pair-end and mate pair reads (UNESP) failed to produce the mate pairs in their first attempt, and they asked for more DNA to repeat the library preparation. As we did not have more tissue from the first specimen, we needed to extract more from a second specimen. After that, as we notice the use of only Illumina would not allow us to produce a contiguous high-quality genome, we decided to sequence PacBio. PacBio libraries need a substantial amount of high-molecular-weight-DNA, and to meet this requirements we needed to extract DNA from a third specimen.

To clarify the use of 3 specimens for the construction of the 3 sequencing libraries, a small complement was added to the sentence in line 103-105. It's as follows "... For the

	<p>genome assembly, a total of 3 individuals were sampled for DNA extraction from gills and to produce the three types of DNA libraries used in this study.”</p> <p>Reviewer 2 : The recent Crown of Thorns sea star genome paper (<a href="http://www.nature.com/nature/journal/v544/n7649/full/nature22033.html?foxtrotcallback=true">http://www.nature.com/nature/journal/v544/n7649/full/nature22033.html?foxtrotcallback=true</a>) would be an appropriate citation on line 82. Response: The citation was added. It’s now present in line 88.</p> <p>Reviewer 2: Line 85: Change "U\$ " to "USD \$" Response: It was changed in line 91.</p> <p>Reviewer 2: Lines 166-167: I suggest the authors move this text to the table. Response: The small paragraph was removed and now it is presented as the last line of Table 3.</p> <p>Lines 266-273: Despite the name, OrthoMCL does not identify orthologs, it identifies gene families. These are gene family comparisons and not strict orthologs. Response: Manuscript was edited. Line 268.</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
<b>Availability of data and materials</b>	Yes

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.

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

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 **DATA NOTE**

2 **A hybrid-hierarchical genome assembly strategy to sequence the invasive golden mussel**

3 ***Limnoperna fortunei***

4 **Authors:** Marcela Uliano-Silva<sup>1,2,3\*</sup>, Francesco Dondero<sup>4</sup>, Thomas Dan Otto<sup>5,6</sup>, Igor Costa<sup>7</sup>,  
5 Nicholas Costa Barroso Lima<sup>7,8</sup>, Juliana Alves Americo<sup>1</sup>, Camila Junqueira Mazzoni<sup>2,3</sup>,  
6 Francisco Prosdocimi<sup>7</sup>, Mauro de Freitas Rebelo<sup>1\*</sup>

7 Marcela Uliano-Silva: [marcela.uliano@gmail.com](mailto:marcela.uliano@gmail.com)

8 Francesco Dondero: [francesco.dondero@uniupo.it](mailto:francesco.dondero@uniupo.it)

9 Thomas D. Otto: [tdo@sanger.ac.uk](mailto:tdo@sanger.ac.uk)

10 Igor Costa: [igor.bioinfo@gmail.com](mailto:igor.bioinfo@gmail.com)

11 Nicholas Costa Barroso Lima: [ncblima@gmail.com](mailto:ncblima@gmail.com)

12 Juliana Alves Americo: [juliana.americo@gmail.com](mailto:juliana.americo@gmail.com)

13 Camila Mazzoni: [mazzoni@izw-berlin.de](mailto:mazzoni@izw-berlin.de)

14 Francisco Prosdocimi: [prosdocimi@bioqmed.ufrj.br](mailto:prosdocimi@bioqmed.ufrj.br)

15 Mauro de Freitas Rebelo: [mrebelo@biof.ufrj.br](mailto:mrebelo@biof.ufrj.br)

16 Affiliations:

17 1 Carlos Chagas Filho Biophysics Institute (IBCCF), Universidade Federal do Rio de Janeiro,

18 Rio de Janeiro, Brazil

19 2 Department of Evolutionary Genetics, Leibniz Institute for Zoo and Wildlife Research, Berlin,

20 Germany

21 3 Berlin Center for Genomics in Biodiversity Research, Berlin, Germany

22 4 Department of Science and Technological Innovation (DiSIT), Università del Piemonte

23 Orientale Amedeo Avogadro, Vercelli-Novara-Alessandria, Italy

24 5 Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton CB10 1SA, UK

25 6 Centre of Immunobiology, Institute of Infection, Immunity & Inflammation, College of

26 Medical, Veterinary and Life Sciences, University of Glasgow, Glasgow, UK.

27 7 Leopoldo de Meis Biomedical Biochemistry Institute (IBqM), Universidade Federal do Rio de

28 Janeiro, Rio de Janeiro, Brazil



1  
2  
3  
4 29 8 Bioinformatics Laboratory (LabInfo) of the National Laboratory for Scientific Computing,  
5  
6 30 Petrópolis, Rio de Janeiro, Brazil

7  
8  
9 31 \*Correspondence: [marcela.uliano@gmail.com](mailto:marcela.uliano@gmail.com); [mrebelo@biof.ufrj.br](mailto:mrebelo@biof.ufrj.br)

10  
11 32 **ABSTRACT**

12  
13  
14 33 **Background:** For more than 25 years, the golden mussel *Limnoperna fortunei* has aggressively  
15  
16 34 invaded South American freshwaters, having travelled more than 5,000 km upstream across five  
17  
18 35 countries. Along the way, the golden mussel has outcompeted native species and economically  
19  
20 36 harmed aquaculture, hydroelectric powers, and ship transit. We have sequenced the complete  
21  
22 37 genome of the golden mussel to understand the molecular basis of its invasiveness and search for  
23  
24 38 ways to control it. **Findings:** We assembled the 1.6 Gb genome into 20548 scaffolds with an  
25  
26 39 N50 length of 312 Kb using a hybrid and hierarchical assembly strategy from short and long  
27  
28 40 DNA reads and transcriptomes. A total of 60717 coding genes were inferred from a customized  
29  
30 41 transcriptome-trained AUGUSTUS run. We also compared predicted protein sets with those of  
31  
32 42 complete molluscan genomes, revealing an exacerbation of protein-binding domains in *L.*  
33  
34 43 *fortunei*. **Conclusions:** We built one of the best bivalve genome assemblies available using a  
35  
36 44 cost-effective approach using Illumina pair-end, mate pair, and PacBio long reads. We expect  
37  
38 45 that the continuous and careful annotation of *L. fortunei*'s genome will contribute to the  
39  
40 46 investigation of bivalve genetics, evolution, and invasiveness, as well as to the development of  
41  
42 47 biotechnological tools for aquatic pest control.

43  
44  
45 48 **KEYWORDS:** Amazon; binding domain; bivalves; genomics; TLR; transposon.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 51 **DATA DESCRIPTION**

52           The golden mussel *Limnoperna fortunei* is an Asian bivalve that arrived in the southern  
53 part of South America about 25 years ago [1]. Research suggests that *L. fortunei* was introduced  
54 in South America through ballast water of ships coming from Hong Kong or Korea [2]. It was  
55 found for the first time in the estuary of the La Plata River in 1991 [1]. Since then, it has moved  
56 ~5,000 km, invading upstream continental waters and reaching northern parts of the continent [3]  
57 leaving behind a track of great economic impact and environmental degradation [4]. The latest  
58 infestation was reported in 2016 in the São Francisco River, one of the main rivers in the  
59 Northeast of Brazil, with a 2,700 km riverbed that provides water to more than 14 million  
60 people. At Paulo Afonso, one of the main hydroelectric power plants in the São Francisco River,  
61 maintenance due to clogging of pipelines and corrosion caused by the golden mussel is estimated  
62 to cost US\$ 700,000 per year (*personal communication, Mizael Gusmã, Chief Maintenance*  
63 *Engineer for Centrais Hidrelétricas do São Francisco – CHESF*).

64           A recent review has shown that, before arriving in South America, *L. fortunei* was  
65 already an invader in China. Originally from the Pearl River Basin, the golden mussel has  
66 traveled 1,500 km into the Yang Tse and the Yellow River basins, being limited further north  
67 only by the extreme natural barriers of Northern China [5]. Today, *L. fortunei* is found in the  
68 Paraguaizinho River, located only 150 km from the Teles-Pires River that belongs to the Alto  
69 Tapajós River Basin and is the first to directly connect with the Amazon River Basin [6]. Due to  
70 its fast dispersion rates, it is very likely that *L. fortunei* will reach the Amazon River Basin in the  
71 near future.

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

72           The reason why some freshwater bivalves, such as *L. fortunei*, *Dreissena polymorpha*,  
73 and *Corbicula fluminea*, are aggressive invaders is not fully understood. These bivalves present  
74 characteristics such as (i) tolerance to a wide range of environmental variables, (ii) short life  
75 span, (iii) early sexual maturation, and (iv) high reproductive rates that allow them to reach  
76 densities as high as 150,000 ind.m<sup>-2</sup> over a year [7, 8] that may explain the aggressive behavior.  
77 On the other hand, these traits are not exclusive to invasive freshwater bivalves and do not  
78 explain how they outcompete native species and disperse so widely.

79           To the best of our knowledge, there are no reports of successful strategies to control the  
80 expansion of mussel invasion in industrial facilities. Bivalves can sense chemicals in the water  
81 and close their valves as a defensive response [9], making them tolerant to a wide range of  
82 chemical substances, including strong oxidants like chlorine [10]. Microencapsulated chemicals  
83 have shown better results in controlling mussel populations in closed environments [10, 11] but  
84 it is unlikely they would work in the wild. Currently, there is no effective and efficient approach  
85 to control the invasion by *L. fortunei*.

86           The genome sequence is one of the most relevant and informative descriptions of species  
87 biology. The genetic substrate of invasive populations, upon which natural selection operates,  
88 can be of primary importance to understand and control a biological invader [12, 13].

89           We have partially funded the golden mussel genome sequencing through a pioneer  
90 crowdfunding initiative in Brazil ([www.catarse.me/genoma](http://www.catarse.me/genoma)). In this campaign, we could raise  
91 around USD\$ 20,000.00 at the same time we promoted scientific education and awareness in  
92 Brazil.

1  
2  
3  
4 93 Here we present the first complete genome dataset for the invasive bivalve *Limnoperna*  
5  
6 94 *fortunei*, assembled from short and long DNA reads and using a hybrid and hierarchical  
7  
8  
9 95 assembly strategy. This high-quality reference genome represents a substantial resource for  
10  
11 96 further studies of genetics and evolution of mussels, as well as for the development of new tools  
12  
13  
14 97 for plague control.

15  
16 98  
17  
18  
19 99 **Genome sequencing in short Illumina and long PacBio reads**

20  
21 100 *Limnoperna fortunei* mussels were collected from the Jacui River, Porto Alegre, Rio  
22  
23 101 Grande do Sul, Brazil (29°59'29.3"S 51°16'24.0"W). Voucher specimens were housed at the  
24  
25  
26 102 zoological collection (specimen number: 19643) of the Biology Institute at the Universidade  
27  
28  
29 103 Federal do Rio de Janeiro, Brazil. For the genome assembly, a total of 3 individuals were  
30  
31 104 sampled for DNA extraction from gills and to produce the three types of DNA libraries used in  
32  
33 105 this study. DNA was extracted using DNeasy Blood & Tissue Kit (Qiagen, Hilden, Germany) to  
34  
35  
36 106 prepare libraries for Illumina Nextera paired-end reads, with ~180bp and ~500bp of insert size,  
37  
38 107 (ii) Illumina Nextera mate-pair reads with insert sizes from 3 to 15 Kb, and (iii) Pacific  
39  
40  
41 108 Biosciences long reads (**Table 1**). Illumina libraries were sequenced respectively in a HiScanSQ  
42  
43 109 or HiSeq 1500 machine, and Pacific Biosciences reads were produced with the P4C6 chemistry  
44  
45  
46 110 and sequenced in 10 SMRT Cells. All Illumina reads were submitted to quality analysis with  
47  
48 111 FastQC (FastQC, RRID:SCR\_014583) followed by trimming with Trimmomatic (Trimmomatic,  
49  
50  
51 112 RRID:SCR\_011848) [14]. Pacific Biosciences adaptor-free subreads sequences were used as  
52  
53 113 input data for the genome assembly.

54  
55 114  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

115

116

117

118 **Table 1 - DNA reads produced for *L. fortunei* genome assembly**

Library technology			Raw data		Trimmed Data*	
	Reads insert size	Pairs	Number of reads	Number of bases	Number of reads	Number of bases
<b>Illumina</b> <b>Nextera</b>	Paired end – 180 bp	R1	209542721	21060365702	209036571	21001101404
		R2	209542721	21049308698	209036571	20991650008
	Paired end – 500 bp	R1	153948902	15472966961	153482290	15423123500
		R2	153948902	15462883157	153482290	15414813589
	Mate pair 3-12 Kb	R1	178392944	18017687344	58157933	5822572152
		R2	178392944	18017687344	58157933	5811310412
<b>Pacific</b> <b>Biosciences</b>	P4C - 10/SMTRC	Subreads	1663730	11171487485		

119

120 \*trimmomatic parameters for Illumina reads - ILLUMINACLIP:NexteraPE-PE.fa:2:30:10  
 121 SLIDINGWINDOW:4:2 LEADING:10 TRAILING:10 CROP:101 HEADCROP:0 MINLEN:80

122

123 For transcriptome sequencing, RNA was sampled from four tissues (gills, adductor  
 124 muscle, digestive gland, and foot) of three different golden mussel specimens. RNA was  
 125 extracted using NEXTflex Rapid Directional RNA-Seq Kit (Bioo Scientifics, TX, USA) and 12  
 126 barcodes from NEXTflex Barcodes compatible with Illumina NexSeq Machine. Resulting reads  
 127 (**Supplementary Table S1**) were submitted to FastQC quality analysis (FastQC,  
 128 RRID:SCR\_014583) and trimmed with Trimmomatic (Trimmomatic, RRID:SCR\_011848) [14]  
 129 for all NEXTflex adaptors and barcodes. A total of 3 sets of *de novo* assembled transcriptomes

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  
2  
3  
4 130 were generated using Trinity (Trinity, RRID:SCR\_013048) (**Table 2**); one set for each specimen  
5  
6 131 was a pool of the 4 tissue samples to avoid assembly bias due to intraspecific polymorphism  
7  
8  
9 132 [15]. All generated sequences are deposited in the SRA Archive under the following accession  
10  
11  
12 133 numbers: SRR5188384, SRR5195098, SRR5188200, SRR5195097, SRR5188315, and  
13  
14 134 SRR5181514. Also this Whole Genome Shotgun project has been deposited in the  
15  
16 135 DDBJ/ENA/GenBank under accession number NFUK00000000. The version described in this  
17  
18  
19 136 paper is version NFUK01000000. Genome files are available in the Gigascience database.  
20  
21  
22 137

23  
24 138 **Table 2 - Trinity assembled transcripts used in the assembly and annotation of *L. fortunei***  
25  
26 139 **genome**

Sample	Pooled tissues	Number of reads prior assembly	Number of Trinity Transcripts	Number of Trinity Genes	Average Contig Length	GC%
Mussel 1	Gills, mantle, digestive gland, foot	406589144	433197	303172	854	34
Mussel 2	Gills, mantle, digestive gland, foot	376577660	435054	298117	824	34
Mussel 3	Gills, mantle, digestive gland, foot	334316116	499392	351649	844	34

27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45 140  
46  
47  
48 141  
49  
50 142  
51  
52 143 **Genome assembly using a hybrid and hierarchical strategy**  
53  
54  
55  
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

144 The Jellyfish software [16] was used to count and determine the distribution frequency of  
145 lengths 25 and 31 k-mers (**Figure 1**) for the Illumina DNA paired-end and mate-pair reads  
146 (**Table 1**). Genome size was estimated to be 1,6 Gb by using the 25 k-mer distribution plot as  
147 total k-mer number and then subtracting erroneous reads (starting k-mer counts from 12 times  
148 coverage), to further divide by the homozygous coverage-peak depth (45 times coverage), as  
149 performed by Li *et al.* (2010) [17]. A double-peak k-mer distribution was used as evidence of  
150 genome diploidy (**Figure 1**) and high heterozygosity. The rate of heterozygosity was estimated  
151 to be 2.07% and it was calculated as described by Vij *et al.* (2016) [18], using as input data the  
152 25-kmer distribution plot for reads from one unique specimen.

153 Initially, we attempted to assemble the golden mussel genome using only short Illumina  
154 reads of different insert sizes (paired-end and mate-pairs, **Table 1**) using traditional *de novo*  
155 assembly software such as ALLPATHS [19], SOAPdenovo [20], and Masurca [21]. All these  
156 attempts resulted in very fragmented genome drafts, with an N50 no higher than 5 Kb and a total  
157 of 4 million scaffolds. To reduce fragmentation, we further sequenced additional long reads (10  
158 PacBio SMTR Cells, **Table 1**) and performed a hybrid and hierarchical *de novo* assembly  
159 described below and depicted in **Figure 2**.

160 First, (i) trimmed paired-end and mate-pair DNA Illumina reads (**Table 1**) were  
161 assembled into contigs using the software Sparse Assembler [22] with parameters *LD 0*  
162 *NodeCovTh 1 EdgeCovTh 0 k 31 g 15 PathCovTh 100 GS 1800000000*. Next, (ii) the resulting  
163 contigs were assembled into scaffolds using Pacific Biosciences long subreads data and the  
164 PacBio-correction-free assembly algorithm DBG2OLC [23] with parameters *LD1 0 k 17*  
165 *KmerCovTh 10 MinOverlap 20 AdaptiveTh 0.01*. Finally, (iii) resulting scaffolds were submitted

1  
2  
3  
4 166 to 6 iterative runs of the program L\_RNA\_Scaffolder [24] that uses exon-distance information  
5  
6 167 from *de novo* assembled transcripts (**Table 2**) to fill gaps and connect scaffolds whenever  
7  
8  
9 168 appropriate. At the end, (iv) the final genome scaffolds were corrected for Illumina and Pacific  
10  
11 169 Biosciences sequencing errors with the software PILON [25]: all DNA and RNA short Illumina  
12  
13  
14 170 reads were re-aligned back to the genome with BWA aligner (BWA , RRID:SCR\_010910) [26]  
15  
16 171 and resulting sam files were BAM-converted, sorted, and indexed with samtools package  
17  
18  
19 172 (SAMTOOLS, RRID:SCR\_002105) [27]. Pilon [25] identifies INDELS and mismatches by  
20  
21 173 coverage of reads and yields a final corrected genome draft. Pilon was run with parameters --  
22  
23  
24 174 *diploid –duplicates*.

25  
26 175 The final genome was assembled in 20,548 scaffolds, with an N50 of 312 Kb and a total  
27  
28  
29 176 assembly length of 1.6 Gb (**Table 3**).

30  
31 177  
32  
33 178

**Table 3: Assembly statistics for *Limnoperna fortunei*'s genome**

<b>Parameter</b>	<b>Value</b>
Estimated genome size by k-mer analysis	1.6 Gb
Total size of assembled genome	1.673 Gb
Number of scaffolds	20548
Number of contigs	61093
Scaffold N50	312 Kb
Maximum scaffold length	2.72 Mb
Percentage of genome in scaffolds > 50 Kb	82,55%
Masked percentage of total genome	33 %
Mapping percentage of Illumina reads back to scaffolds	91 %

34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54 179  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65



1  
2  
3  
4 180 The golden mussel genome presents 81% of all Benchmarking Universal Single Copy  
5  
6 181 Orthologs (BUSCO version 3.3 analysis with Metazoa database) (BUSCO, RRID:SCR\_015008)  
7  
8  
9 182 (**Table 4**) and, compared to the mollusk genomes currently available [28, 29, 30, 31, 32, 33, 34  
10  
11 183 35] it represents one of the best assemblies of molluscan genomes so far also in terms of scaffold  
12  
13  
14 184 N50 and contiguity (**Table 5**).

15  
16 185 One main challenges of assembling bivalve genomes lies in the high heterozygosity and  
17  
18  
19 186 amount of repetitive elements these organisms present: (i) the mussels *L. fortunei* and *Modiolus*  
20  
21 187 *philippinarum* and the oyster *Crassostrea gigas* genomes were estimated to have heterozygosity  
22  
23  
24 188 rates of 2.07%, 2.02 % 1.95% respectively, which is substantially higher than other animal  
25  
26 189 genomes [29], and (ii) repetitive elements correspond to at least 30% of the genomes of all  
27  
28  
29 190 studied bivalves so far (**Table 3**) [28, 29, 30, 31, 33, 34, 35 ]. Also, retroelements might be active  
30  
31 191 in some species such as *L. fortunei* (refer to the retroelements-related section of this paper) and  
32  
33  
34 192 *C. gigas* [29], allowing genome rearrangements that may hinder for genome assembly. One  
35  
36 193 exception seems to be the deep-sea mussel *B. platifrons* which has lower heterozygosity rates  
37  
38  
39 194 compared to other bivalves [31]. Sun *et al.*, (2017) [31] suggested it might be due to recurrent  
40  
41 195 population bottlenecks happened after events of population extinction and recolonization in the  
42  
43 196 extreme environment [31]. Nevertheless, most of the bivalve genome projects relying only on  
44  
45  
46 197 short Illumina reads are likely to present fragmented initial drafts [28, 30]. PacBio long reads  
47  
48 198 allowed us to increase the N50 to 32 Kb and to reduce the number of scaffolds from millions to  
49  
50  
51 199 61102, using the DBG2OLC [23] assembler. Finally, interactive runs of L\_RNA\_scaffolder [24]  
52  
53 200 using the transcriptomes (**Table 2**) rendered the final result of N50 312 Kb in 20548 scaffolds.  
54  
55 201 Thus, our assembly strategy of Illumina contigs, low coverage of PacBio reads, transcriptome  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 202 and Illumina re-mapping for final correction (**Figure 2**) represents an option for cost-efficient  
5  
6  
7 203 assembly of highly heterozygous genomes of nonmodel species such as bivalves.  
8  
9 204

10  
11 **Table 4: Summary statistics of Benchmarking Universal Single-Copy Orthologs**  
12 **(BUSCO) analysis for *L. fortunei* genome run for Metazoans**  
13 206  
14 207

Categories	Number of Genes	Percentage (%)
Total BUSCO groups searched	978	--
Complete BUSCOs	801	81.9%
Complete and single-copy BUSCOs	769	78.62%
Complete and duplicated BUSCOs	32	3.27%
Fragmented BUSCOs	72	7.36%
Missing BUSCOs	105	10.73%

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33 208  
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

209 **Table 5: Comparison of genome assembly statistics for molluscan genomes.**

	<i>Haliotis discus hannai</i>	<i>Lottia gigantea</i>	<i>Aplysia californica</i>	<i>Ruditapes philippinarum</i>	<i>Patinopecten yessoensis</i>	<i>Crassostrea gigas</i>	<i>Pinctada fucata</i>	<i>Mytillus galloprovincialis</i>	<i>Bathymodiolus platifrons</i>	<i>Modiolus philippinarum</i>	<i>Limnoperna fortunei</i>
<b>Estimated genome size</b>	1.65Gb	359.5 Mb	1.8Gb	1.37 Gb	1.43 Gb	545 Mb	1.15 Gb	1.6 Gb	1.64Gb	2.38 Gb	1.6 Gb
<b>Number of scaffolds</b>	80,032	4,475	8,766	223,851	82,731	11,969	7997	1,746,447	65,664	74,575	<b>20,548</b>
<b>Total size of scaffolds</b>	1,865,475,499	359,512,207	715,791,924	2,561,070,351	987,685,017	558,601,156	915,721,316	1,599,211,957	1,659,280,971	2,629,649,654	<b>1,673,125,894</b>
<b>Longest scaffold</b>	2,207,537	9,386,848	1,784,514	572,939	7,498,238	1,964,558	5,897,787	67,529	2,790,175	715382	<b>2,720,304</b>
<b>Shortest scaffold</b>	854	1000	5001	500	200	100	1807	100	292	205	<b>558</b>
<b>Number of scaffolds &gt; 1 K nt</b>	79,923 (99.9%)	4,471 (99.9%)	8,766 (100.0%)	138,771	16,004	5,788 (48.4%)	7997 (100%)	393,685 (22.5%)	38,704 (58.9%)	44,921 (60.2%)	<b>20,547 (100%)</b>
<b>Number of scaffolds &gt; 1 M nt</b>	67 (0.1%)	98 (2.2%)	27 (0.3%)	0 (0.0%)	248 (0.3%)	60 (0.5%)	27 (0.3%)	0 (0.0%)	164 (0.2%)	0 (0%)	<b>95 (0.5%)</b>
<b>Mean scaffold size</b>	23,309	80,338	81,655	11,441	11,939	46,671	114,508	916	25,269	35,262	<b>81,425</b>
<b>Median scaffold size</b>	1,697	3,622	13,763	1,327	362	824	14,683	258	1,284	13,722	<b>22,134</b>
<b>N50 scaffold length</b>	200,099	1,870,055	264,327	48,447	803,631	401,319	345,846	2,651	343,373	100,161	<b>312,020</b>
<b>Sequencing coverage</b>	322 X	8.87 X	11 X	39.7 X	297 X	155 X	234 X	32 X	319 X	209.5 X	<b>60 X</b>
<b>Sequencing Technology</b>	Illumina + PacBio	Sanger	Sanger	Illumina	Illumina	Illumina	Illumina + BACs	Illumina	Illumina	Illumina	<b>Illumina + PacBio</b>

1  
2  
3  
4 212  
5 213  
6  
7 214 **Around 10% of repetitive elements are transposons**  
8

9 215 Initial masking of *L. fortunei* genome was done using RepeatMasker program  
10  
11 216 (RepeatMasker, RRID:SCR\_012954) [36] with parameter *-species bivalves* and masked 3.4% of  
12  
13  
14 217 the total genome. This content was much lower than the masked portion of other molluscan  
15  
16 218 genomes: 34% in *C. gigas* [29] and 36% in *M. galloprovincialis* [28], suggesting that the fast  
17  
18  
19 219 evolution of interspersed elements limits the use of repeat libraries from divergent taxa [37].  
20  
21 220 Thus, we generated a *de novo* repeat library for *L. fortunei* using the program RepeatModeler  
22  
23 221 (RepeatModeler, RRID:SCR\_015027) [38] and its integrated tools (RECON [39], TRF [40], and  
24  
25  
26 222 RepeatScout [41]). This *de novo* repeat library was the input to RepeatMasker together with the  
27  
28  
29 223 first masked genome draft of *L. fortunei*, and resulted in a final masking of 33.4% of the genome.  
30  
31 224 Even though more than 90% of the repeats were not classified by RepeatMasker  
32  
33 225 (**Supplementary Table S2**), 8.85% of the repeats were classified as LINEs, Class I transposable  
34  
35  
36 226 elements. In addition, large numbers of reverse-transcriptases (824 counts, Pfam RVT\_1  
37  
38 227 PF00078), transposases (177 counts, Pfam HTH\_Tnp\_Tc3\_2 PF01498), and integrases (501  
39  
40  
41 228 counts, Pfam Retroviral integrase core domain PF00665) and other related elements were  
42  
43 229 detected; over 98% of these had detectable transcripts.  
44  
45

46 230  
47  
48 231 **More than 30,000 sequences identified by gene prediction and automated**  
49  
50 232 **annotation.**

51  
52  
53 233 To annotate the golden mussel genome, we sequenced a number of transcriptomes (**Table S1**),  
54  
55 234 *de novo* assembled (**Table 2**) and aligned these transcriptomes to the genome scaffolds, and  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 235 created gene models with the PASA pipeline [36]. These models were used to train and run the  
5  
6 236 *ab initio* gene predictor AUGUSTUS (Augustus: Gene Prediction, RRID:SCR\_008417) [37]  
7  
8  
9 237 (**Supplementary Figure S1**). The complete gene models yielded by PASA [42] were BLASTed  
10  
11 238 (e-value 1e-20) against the Uniprot database (UniProt, RRID:SCR\_002380) and those with 90%  
12  
13  
14 239 or more of their sequences showing in the BLAST hit alignment were considered for further  
15  
16 240 analysis. Next, all the necessary filters to run an AUGUSTUS [43] personalized training were  
17  
18  
19 241 performed: (i) only gene models with more than 3 exons were maintained, (ii) sequences with  
20  
21 242 90% or more overlap were withdrawn and only the longest sequences were retained, and (iii)  
22  
23  
24 243 only gene models free of repeat regions, as indicated by BLASTN similarity searches with *de*  
25  
26 244 *novo* library of repeats, were maintained. These curated data yielded a final set of 1,721 gene  
27  
28  
29 245 models on which AUGUSTUS [35] was trained in order to predict genes in the genome using the  
30  
31 246 default AUGUSTUS [43] parameters. Once the gene models were predicted, a final step was  
32  
33  
34 247 performed by using the PASA pipeline [42] once again in the *update* mode (parameters -c -A -g -  
35  
36 248 t). This final step compared the 55,638 gene models predicted by AUGUSTUS [43] with the  
37  
38 249 40,780 initial transcript-based gene-structure models from PASA [42] to generate the final set of  
39  
40  
41 250 60,717 gene models for *L. fortunei*. Of those, 58% had transcriptional evidence based on RNA  
42  
43 251 Illumina reads (**Table S2**) re-mapping, rate that was expected since our RNA-Seq libraries were  
44  
45  
46 252 constructed only for 4 tissues of adult golden mussel specimens without any environmental  
47  
48 253 stresses induction (**Table 2**). Therefore, these libraries lack transcripts for developmental stages,  
49  
50  
51 254 for some other cell types (i.e. hemocytes) and stress-inducible genes. Finally, 67% of the gene  
52  
53 255 models were annotated by homology searches against Uniprot or NCBI NR (**Table 6**).

54  
55 256  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 257  
5  
6  
7 258  
8  
9 259  
10  
11 260  
12  
13  
14 261  
15 262  
16  
17 263  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40 264  
41 265  
42 266  
43  
44  
45 267  
46  
47 268  
48  
49  
50 269  
51  
52 270  
53  
54  
55 271  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 6: Summary of gene annotation against various databases for *L. fortunei* whole genome-predicted genes**

<b>Total number of genes</b>	60,717
<b>Total number of exons</b>	220,058
<b>Total number of proteins</b>	60,717
<b>Average protein size</b>	304 aa
<b>Number of protein BLAST hits* with Uniprot</b>	26,198
<b>Number of protein BLAST hits* with NR NCBI (no hits with Uniprot)</b>	14,810
<b>Number of protein HMMER hits* with Pfam.A</b>	24,513
<b>Number with proteins with KO assigned by KEGG</b>	8,387
<b>Number of proteins with BLAST hits* with EggNOG</b>	36,868

\*all considered hits had a minimum e-value of 1e-05

**Protein clustering indicates evolutionary proximity among mollusks species.**

Gene family relationships were assigned using reciprocal best BLAST and OrthoMCL software (version 1.4) [44] between *L. fortunei* proteins and the total protein set predicted for nine other mollusks: the mussels *M. galloprovincialis*, *M. philippinarum* and *B. platifrons*, the clam *Ruditapes philippinarum*, the scallop *Patinopecten yessoensis*, the 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

272 the pearl oyster *Pinctada fucata* (genome version from Du *et al* [35]), and the gastropods *Lottia*  
273 *gigantea* and *Haliotis discus hannai* (see **Supplementary Table S3** for detailed information on  
274 the comparative data). **Figure 3A** presents orthologs relationships for five of the bivalves  
275 analyzed. A total of 6,337 orthologs groups are shared among the five bivalve species.

276 Of all the orthologous found for the total 10 species, 44 groups are composed of single-  
277 copy orthologs containing one representative protein sequence of each species. These sequences  
278 were used to reconstruct a phylogeny: the single-copy orthologs sequences were concatenated  
279 and aligned with CLUSTALW [45] with a resulting alignment of 30755 sites in length (**Figure**  
280 **3B**). ProtTest 3.4.2 [46] was used to estimate the best fitting substitution model, which was VT  
281 [47]. With this alignment and model we reconstructed the phylogeny using PhyML [48] and 100  
282 bootstrap repetition, the resulting tree is shown on **Figure 3B**.

283  
284 **Protein domain analysis shows expansion of binding domain in *L. fortunei*.**

285 We performed a quantitative comparison of protein domains predicted from whole  
286 genome projects of 10 molluscan species. The complete protein sets of *M. galloprovincialis*, *M.*  
287 *philippinarum* and *B. platifrons*, *Ruditapes philippinarum*, *Patinopecten yessoensis*, *C. gigas*,  
288 *Pinctada fucata*, *Lottia gigantea* and *Haliotis discus hannai* (**Supplementary Table S3**) were  
289 submitted to domain annotation using HMMER against Pfam-A database (e-value 1e-05).  
290 Protein expansions in *L. fortunei* were rendered using the normalized Pfam count value  
291 (average) obtained from the other nine mollusks, according to a model based on the Poisson  
292 cumulative distribution. Bonferroni correction ( $p \leq 0.05$ ) was applied for false discovery and  
293 absolute frequencies of Pfam-assigned-domains were initially normalized by the total count

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

294 number of Pfam-assigned-domains found in *L. fortunei* to compensate for discrepancies in  
295 genome size and annotation bias.

296 For *L. fortunei*, the annotation against Pfam.A classified 40127 domains in 24513 gene  
297 models of which 83 and 67 were respectively expanded or contracted in comparison with the  
298 other mollusks (**Supplementary Table S4 and S5; Figure 4A**). The 83 overrepresented domains  
299 were further analyzed for functional enrichment using domain-centric Gene Ontology (**Figure**  
300 **4B**). The analysis shows a prominent expansion of binding domains in *L. fortunei*, such as  
301 Thrombospondin (TSP\_1), Collagen, Immunoglobulins (Ig, I-set,Izumo-Ig Ig\_3), and Ankyrins  
302 (Ank\_2, Ank\_3, and Ank\_4). These repeats have a variety of binding properties and are involved  
303 in cell-cell, protein-protein and receptor-ligand interactions driving evolutionary improvement of  
304 complex tissues and immune defense system in metazoans [49, 50, 51, 52, 53]. An evolutionary  
305 pressure towards the development of a diversificated innate immune system is also suggested by  
306 the high amount of Leucine Rich Repeats (LRR) and Toll/interleukin-1 receptor homology  
307 domains (TIR). Death, another over-represented PFAM, is also part of TLR signaling, being  
308 present in several docking proteins such as Myd88, Irak4 and Pelle [54]. Interestingly, BLAST  
309 analysis of *L. fortunei* gene models against Uniprot identified two types of Toll Like Receptors  
310 (TLRs) whose prototypical architecture of N-terminal extracellular leucine-rich repeat (LRR)  
311 motifs including either a single or multiple cysteine cluster domain, a C-terminal TIR domain  
312 spaced by a single transmembrane-spanning domain [55] could be correctly identified using the  
313 Simple Modular Architecture Research Tool (SMART) [56]. Indeed, we confirmed 141  
314 sequences with similarity to single cysteine clusters TLRs (scc) typical of vertebrates, and 29  
315 sequence hits with the multiple cysteine cluster TLRs (mcc) typical of *Drosophila* [55].



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

316 Phylogenetic analysis of all sequences (using PhyML [48], model JTT) (**Supplementary Figure**  
317 **S2**) shows evidence for TLRs clade separation in *L. fortunei*; the scc TLRs exhibit a higher  
318 degree of amino acid changes, higher molecular evolution, and diversification than the mcc  
319 TLRs. Overall, the expansion of these gene families might suggest an improved resistance to  
320 infections. It is, however, equally curious that other immune-related gene families such as  
321 Fribinogen\_C and C1q seem to be contracted (**Supplementary Table S5**). This feature may  
322 depend on the evolutionary-driven, yet random, fate of the *L. fortunei* genome and consequence  
323 of different specific duplicate genes in other species. Also, other protein families involved in  
324 toxin metabolism, especially glutathione based processes and sulfotransferases are clearly  
325 contracted (**Table S5**).

**326 Final considerations**

327 Here we have described the first version of the golden mussel complete genome and its  
328 automated gene prediction that were funded through a crowdfunding initiative in Brazil. This  
329 genome contains valuable information for further evolutionary studies of bivalves and metazoa  
330 in general. Additionally, our team will further search for the presence of proteins of  
331 biotechnology interest such as the adhesive proteins produced by the foot gland that we have  
332 described elsewhere [57], or genes related to the reproductive system that have been shown to be  
333 very effective for invertebrate plague control [58]. The golden mussel genome and the predicted  
334 proteins are available for download in the Gigabase repository and the scientific community is  
335 welcome to further curate the gene predictions.

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

336 As the golden mussel advances towards the Amazon river basin, the information provided in this  
337 study may be used to help developing biotechnological strategies that may control the expansion  
338 of this organism in both industrial facilities and open environment.

339

#### 340 **Availability of supporting data**

341 *Limnoperna fortunei*'s genome and transcriptome data are available in the Sequence  
342 Read Archive (SRA) as BioProject PRJNA330677 and under the accession numbers  
343 **SRR5188384, SRR5195098, SRR518800, SRR5195097, SRR5188315, SRR5181514**. Also  
344 this Whole Genome Shotgun project has been deposited in the DDBJ/ENA/GenBank under  
345 accession number NFUK00000000. The version described in this paper is version  
346 NFUK01000000.

347

#### 348 **Additional files**

349 **Supplementary Table S1.** RNA raw reads sequenced for 3 *L. fortunei* specimens, 4 tissues each.

350 **Supplementary Table S2:** RepeatMasker classification of repeats predicted in *L. fortunei*  
351 genome.

352 **Supplementary Table S3:** Details of the online availability of the data used for ortholog  
353 assignment and protein domain expansion analysis.

354 **Supplementary Table S4:** Expanded protein families in *L. fortunei* genome.

355 **Supplementary Table S5:** Contracted protein families in *L. fortunei* genome.

356 **Supplementary Table S6:** Fantasy names given to *L. fortunei* genes and proteins from the  
357 backers that have supported us through crowdfunding ([www.catarse.me/genoma](http://www.catarse.me/genoma)).

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

**Supplementary Figure 1:** Steps performed for the prediction and annotation of *L. fortunei* genome.

**Supplementary Figure 2:** Phylogenetic tree of Toll-like (TLRs) receptors found in *L. fortunei* genome.

**List of Abbreviations**

BUSCO: Benchmarking Universal Single-Copy Orthologs; SRA: Sequence Read Archive; KEGG: Kyoto Encyclopedia of Genes and Genomes.

**Competing interests**

The authors declare that they have no competing interests.

**Authors' contribution**

Conceived and designed the experiments: MR, MU, TO, CM, FD. Performed the experiments: MU, JA. Analyzed the data: MU, TO, CM, FD, FP, NC, IC, MR. Contributed reagents/materials/analysis tools: MR, FP, CM. Wrote the paper: MU, FD, MR. All authors read and approved the final manuscript.

**Funding**

This work was supported by the Brazilian Government agencies CAPES (PVE 71/2013), FAPERJ APQ1 (2014), and FAPERJ/DFG (39/2014). Also, this work was funded through crowdfunding with the support of 346 people ([www.catarse.me/genoma](http://www.catarse.me/genoma)).

**Acknowledgements**

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

378 We thank Susan Mbedi and Kirsten Richter from BeGenDiv for RNA-Seq library preparation  
379 and sequencing. We thank Dr. Loris Bennett for IT support while performing bioinformatics  
380 analysis.

381 We especially want to thank the 346 backers that supported the sequencing of the golden mussel  
382 through crowdfunding, in a 2013 campaign that raised U\$ 20,000.00 ([www.catarse.me/genoma](http://www.catarse.me/genoma)).  
383 We decided to give fantasy names to the genes and proteins that we found in the genome, to  
384 thank the backers for their support. The name list is available in **Supplementary Table S6**.

**Consent for publication**

Does not apply.

**Ethics approval**

*Limnoperna fortunei* specimens used for DNA extraction and sequencing were collected in the  
Jacuí River (29°59'29.3"S 51°16'24.0"W), southern Brazil. This bivalve is an exotic species in  
Brazil and is not characterized as an endangered or protected species.

**References**

1. Pastorino G, Darrigran G, et al., *Limnoperna fortunei* (Dunker, 1857) (Mytilidae), nuevo bivalvo invasor em águas Del Rio de la Plata. *Neotropica*. 1993;39:101–2.
2. Darrigran G. Potential impact of filter-feeding invaders on temperate inland freshwater environments. *Biol Invasions* 2002; 4:145–156.

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

399 3. Uliano-Silva M, Fernandes F da C, Holanda IBB, Rebelo MF. Invasive species as a threat  
400 to biodiversity: The golden mussel *Limnoperna fortunei* approaching the Amazon River  
401 basin. In: Exploring Themes Aquat Toxicol Alodi, S, editor. Research Signpost; 2013.  
402  
403 4. Boltovskoy D, Correa N. Ecosystem impacts of the invasive bivalve *Limnoperna fortunei*  
404 (golden mussel) in South America. *Hydrobiologia*. 2015;746(1):81–95.  
405  
406 5. Xu M. Distribution and Spread of *Limnoperna fortunei* in China. In: *Limnoperna fortunei*  
407 Boltovskoy D, editor. Cham: Springer International Publishing; 2015 p. 313–20.  
408  
409 6. Oliveira M, Hamilton S, Jacobi C. Forecasting the expansion of the invasive golden  
410 mussel *Limnoperna fortunei* in Brazilian and North American rivers based on its  
411 occurrence in the Paraguay River and Pantanal wetland of Brazil. *Aquat Invasions*.  
412 2010;5(1):59–73.  
413  
414 7. Karatayev AY, Boltovskoy D, Padilla DK, Burlakova LE. The invasive bivalves  
415 *dreissena polymorpha* and *limnoperna fortunei*: parallels, contrasts, potential spread and  
416 invasion impacts. *J Shellfish Res*. 2007 1;26(1):205–13.  
417  
418 8. Orensanz JM (Lobo), Schwindt E, Pastorino G, Bortolus A, Casas G, Darrigran G, et al.  
419 No Longer The Pristine Confines of the World Ocean: A Survey of Exotic Marine  
420 Species in the Southwestern Atlantic. *Biol Invasions*. 2002 1;4(1–2):115–43.  
421

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 9. Claudi R and Mackie GL. Practical manual for zebra mussel monitoring and control.  
423 Lewis Publishers, Boca. Raton, Florida, 1994. p 227

424

425 10. Calazans SHC, Americo JA, Fernandes F da C, Aldridge DC, Rebelo M de F.  
426 Assessment of toxicity of dissolved and microencapsulated biocides for control of the  
427 Golden Mussel *Limnoperna fortunei*. *Mar Environ Res.* 2013 91:104–8.

428

429 11. Aldridge DC, Elliott P, Moggridge G.D. Microencapsulated biobullets for the control of  
430 biofouling zebra mussels. *Environ. Sci. Technol.* 2006 40:975-979.

431

432 12. Cox GW. Alien species and evolution: the evolutionary ecology of exotic plants, animals,  
433 microbes, and interacting native species. Washington: Island Press; 2004. 377 p.

434

435 13. Hall MR, Kocot KM, Baughman KW, Fernandez-Valverde SL, Gauthier MEA,  
436 Hatleberg WL, et al. The crown-of-thorns starfish genome as a guide for biocontrol of  
437 this coral reef pest. *Nature.* 2017 13;544(7649):231–4.

438

439 14. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina  
440 Sequence Data. *Bioinformatics.* 2014 1;170.

441

442 15. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, et al. Differential gene and  
443 transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat*  
444 *Protoc.* 2012 1;7(3):562–78.

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

446 16. Marçais G, Kingsford C. A Fast, Lock-free Approach for Efficient Parallel Counting of  
447 Occurrences of K-mers. *Bioinformatics*. 2011 Mar;27(6):764–770.

448  
449 17. Li R, Fan W, Tian G, Zhu H, He L, Cai J, et al. The sequence and de novo assembly of  
450 the giant panda genome. *Nature*. 2010 Jan 21;463(7279):311–7.

451  
452 18. Vij S, Kuhl H, Kuznetsova IS, Komissarov A, Yurchenko AA, Heusden PV, et al.  
453 Chromosomal-Level Assembly of the Asian Seabass Genome Using Long Sequence  
454 Reads and Multi-layered Scaffolding. *PLOS Genet*. 2016 Apr 15;12(4):e1005954.

455  
456 19. Gnerre S, Maccallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, et al. High-  
457 quality draft assemblies of mammalian genomes from massively parallel sequence data.  
458 *Proc Natl Acad Sci U S A*. 2011 25;108(4):1513–8.

459  
460 20. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically  
461 improved memory-efficient short-read de novo assembler. *GigaScience*. 2012;1:18.

462  
463 21. Zimin AV, Marçais G, Puiu D, Roberts M, Salzberg SL, Yorke JA. The MaSuRCA  
464 genome assembler. *Bioinformatics*. 2013 1;29(21):2669–77.

465  
466 22. Ye C, Ma Z, Cannon CH, Pop M, Yu DW. Exploiting sparseness in de novo genome  
467 assembly. *BMC Bioinformatics*. 2012;13(Suppl 6):S1.

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

469 23. Ye C, Hill CM, Wu S, Ruan J, Ma Z (Sam). DBG2OLC: Efficient Assembly of Large  
470 Genomes Using Long Erroneous Reads of the Third Generation Sequencing  
471 Technologies. *Sci Rep.* 2016 30;6:31900.

472 24. Xue W, Li J-T, Zhu Y-P, Hou G-Y, Kong X-F, Kuang Y-Y, et al. L\_RNA\_scaffolder:  
473 scaffolding genomes with transcripts. *BMC Genomics.* 2013;14:604.

474 25. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: An  
475 Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly  
476 Improvement. *PLOS ONE.* 2014 19;9(11):e112963.

477 26. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler  
478 transform. *Bioinformatics.* 2009 15;25(14):1754–60.

479 27. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
480 Alignment/Map format and SAMtools. *Bioinformatics.* 2009 15;25(16):2078–9.

481 28. Murgarella M, Puiu D, Novoa B, Figueras A, Posada D, Canchaya C. A First Insight into  
482 the Genome of the Filter-Feeder Mussel *Mytilus galloprovincialis*. *PLOS ONE.* 2016  
483 15;11(3):e0151561.

484 29. Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al. The oyster genome reveals stress  
485 adaptation and complexity of shell formation. *Nature.* 4;490(7418):49–54.



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

493 30. Takeuchi T, Kawashima T, Koyanagi R, Gyoja F, Tanaka M, Ikuta T, et al. Draft genome  
494 of the pearl oyster *Pinctada fucata*: a platform for understanding bivalve biology. *DNA*  
495 *Res Int J Rapid Publ Rep Genes Genomes*. 2012 19(2):117–30.

496 31. Sun J, Zhang Y, Xu T, Zhang Y, Mu H, Zhang Y, et al. Adaptation to deep-sea  
497 chemosynthetic environments as revealed by mussel genomes. *Nat Ecol Evol*. 2017 Apr  
498 3;1(5):0121

500 32. Nam B-H, Kwak W, Kim Y-O, Kim D-G, Kong HJ, Kim W-J, et al. Genome sequence of  
501 pacific abalone (*Haliotis discus hannai*): the first draft genome in family Haliotidae.  
502 *GigaScience*. 2017 May;6(5):1–8.

503 33. Wang S, Zhang J, Jiao W, Li J, Xun X, Sun Y, et al. Scallop genome provides insights  
504 into evolution of bilaterian karyotype and development. *Nat Ecol Evol*. 2017 Apr  
505 3;1(5):0120.

506 34. Mun S, Kim Y-J, Markkandan K, Shin W, Oh S, Woo J, et al. The Whole-Genome and  
507 Transcriptome of the Manila Clam (*Ruditapes philippinarum*). *Genome Biol Evol*. 2017  
508 Jun;9(6):1487–98

509 35. Du X, Fan G, Jiao Y, Zhang H, Guo X, Huang R, et al. The pearl oyster *Pinctada fucata*  
510 *martensii* genome and multi-omic analyses provide insights into biomineralization.  
511 *GigaScience*. 2017 Aug;6(8):1–12

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

517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541

36. Smit AF., Hubley R, Green PJ. RepeatMasker Open-3.0. 1996 2010.

37. Fu H, Dooner HK. Intraspecific violation of genetic colinearity and its implications in maize. Proc Natl Acad Sci U S A. 2002 9;99(14):9573–8.

38. Smith AFA, Hubley R. RepeatModeler Open-1.0. [Internet]. 2014. Available from: <http://www.repeatmasker.org>

39. Bao Z, Eddy SR. Automated De Novo Identification of Repeat Sequence Families in Sequenced Genomes. Genome Res. 2002 12(8):1269–76.

40. Benson G. Tandem repeats finder: a program to analyze DNA sequences. Nucleic Acids Res. 1999 1;27(2):573–80.

41. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes. Bioinformatics. 2005 1;21(Suppl 1):i351–8

42. Haas BJ, Delcher AL, Mount SM, Wortman JR, Jr RKS, Hannick LI, et al. Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. Nucleic Acids Res. 2003 1;31(19):5654–66.

43. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically mapped cDNA alignments to improve de novo gene finding. Bioinforma Oxf Engl. 2008 1;24(5):637–44.

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

542 44. Li L, Stoeckert CJ, Roos DS. OrthoMCL: Identification of Ortholog Groups for  
543 Eukaryotic Genomes. *Genome Res.* 2003 1;13(9):2178–89.

544  
545 45. Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: improving the sensitivity of  
546 progressive multiple sequence alignment through sequence weighting, position-specific  
547 gap penalties and weight matrix choice. *Nucleic Acids Res.* 1994 11;22(22):4673–80.

548  
549 46. Darriba D, Taboada GL, Doallo R, Posada D. ProtTest 3: fast selection of best-fit models  
550 of protein evolution. *Bioinformatics*, 2011 27:1164-1165.

551  
552 47. Müller T, Vingron M. Modeling amino acid replacement. *J. Comput Biol.* 2000. 7:761-  
553 776.

554  
555 48. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large  
556 phylogenies by maximum likelihood. *Syst Biol.* 2003. 52: 696-704.

557  
558 49. Björklund ÅK, Ekman D, Elofsson A. Expansion of Protein Domain Repeats. *PLoS*  
559 *Comput Biol.* 2006;2(8):e114.

560  
561 50. Rennemeier C, Hammerschmidt S, Niemann S, Inamura S, Zähringer U, Kehrel BE.  
562 Thrombospondin-1 promotes cellular adherence of gram-positive pathogens via  
563 recognition of peptidoglycan. *FASEB J Off Publ Fed Am Soc Exp Biol.* 2007  
564 21(12):3118–32.

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

565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589

51. Schmucker D, Chen B. Dscam and DSCAM: complex genes in simple animals, complex animals yet simple genes. *Genes Dev.* 2009 15;23(2):147–56.

52. Pancer Z, Amemiya CT, Ehrhardt GRA, Ceitlin J, Larry Gartland G, Cooper MD. Somatic diversification of variable lymphocyte receptors in the agnathan sea lamprey. *Nature.* 2004 Jul 8;430(6996):174–80.

53. Tucker RP. The thrombospondin type 1 repeat superfamily. *Int J Biochem Cell Biol.* 2004 36(6):969–74.

54. Park HH, Lo YC, Lin SC, Wang L, Yang, JK, Wu H. The death domain superfamily in intracellular signaling of apoptosis and inflammation. *Annu. Rev. Immunol.* 2007 25, 561-586.

55. Leulier F, Lemaitre B..Toll-like receptors—taking an evolutionary approach. *Nature Reviews Genetics*, 2008 9,3: 165-178.

56. Schultz J, Milpetz F, Bork P, Ponting CP. SMART, a simple modular architecture research tool: identification of signaling domains. *Proc Natl Acad Sci U S A.* 1998 26;95(11):5857–64

57. Uliano-Silva M, Americo JA, Brindeiro R, Dondero F, Prosdocimi F, Rebelo M de F. Gene discovery through transcriptome sequencing for the invasive mussel *Limnoperna fortunei*. *PloS One.* 2014;9(7):e10297.

1  
2  
3  
4 590  
5  
6 591 58. Hammond A, Galizi R, Kyrou K, Simoni A, Siniscalchi C, Katsanos D, et al. A CRISPR-  
7  
8 592 Cas9 gene drive system targeting female reproduction in the malaria mosquito vector  
9  
10 593 *Anopheles gambiae*. Nat Biotechnol. 2015 7;34(1):78–83.

11  
12  
13 594  
14 595 59.

15  
16  
17 596 **Figure 1:** K-mer distribution of *Limnoperna fortunei* Illumina DNA reads (Table 1).

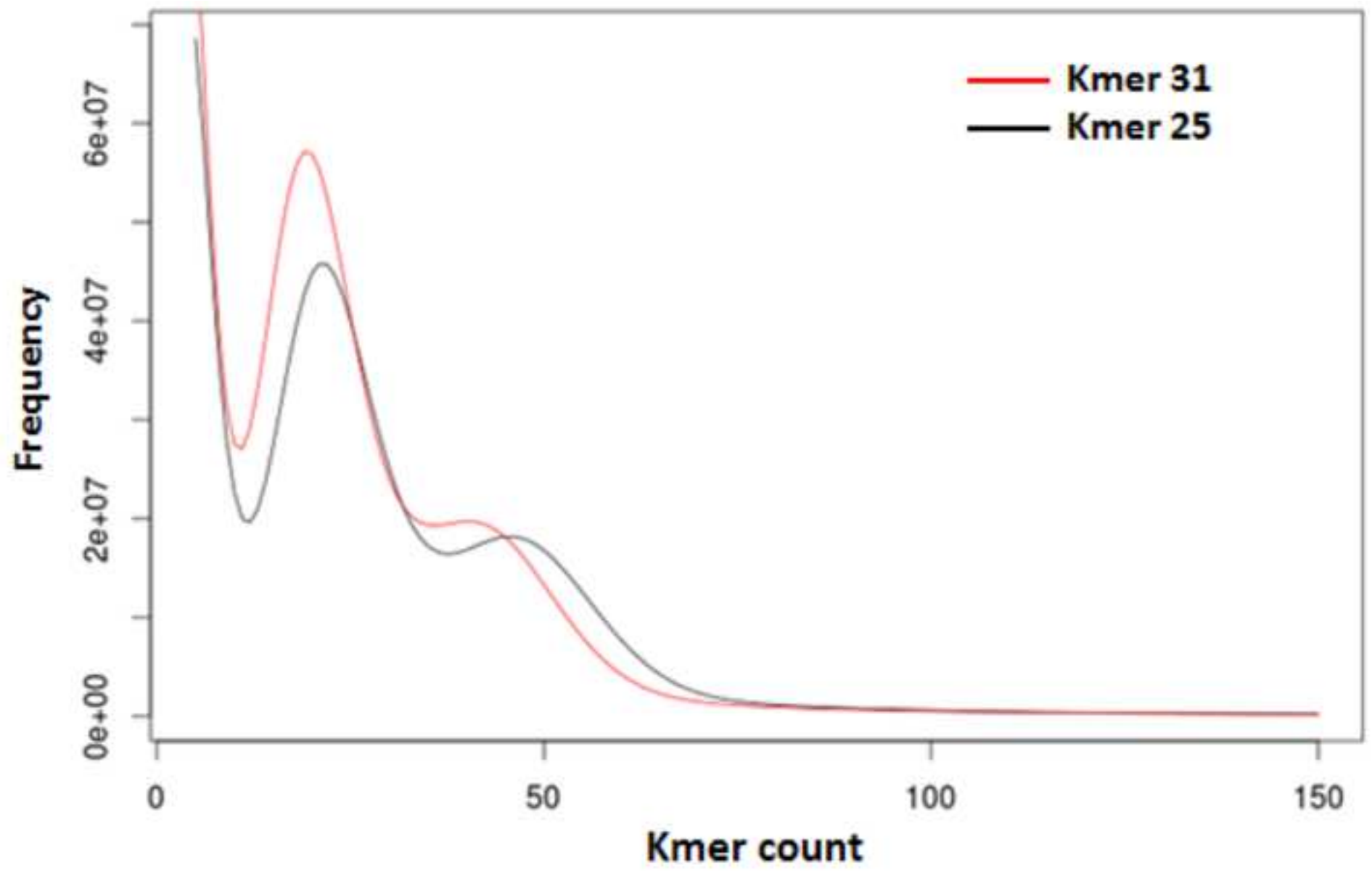
18 597  
19 598  
20 599 **Figure 2: Hierarchical assembly strategy employed for the golden mussel genome**  
21 600 **assembly.** Trimmed Illumina reads were assembled to the level of contigs with Sparse  
22 601 Assembler algorithm (**Step 1**). Then, Illumina contigs and PacBio reads were used to build  
23 602 scaffolds with DBG2OLC assembler, that anchors Illumina contigs to erroneous PacBio  
24 603 subreads, correcting them and building longer scaffolds (**Step 2**), followed by transcriptome  
25 604 joining scaffolds using L\_RNA\_scaffolder (**Step 3**). Final scaffolds were corrected by re-  
26 605 aligning all Illumina DNA and RNA-seq reads back to them and calling consensus with Pilon  
27 606 software (**Step 4**). In bold is bioinformatics software used in each step. Red blocks indicate  
28 607 PacBio errors, which are represented by insertions and/or deletions found in approximately 12%  
29 608 of PacBio subreads.

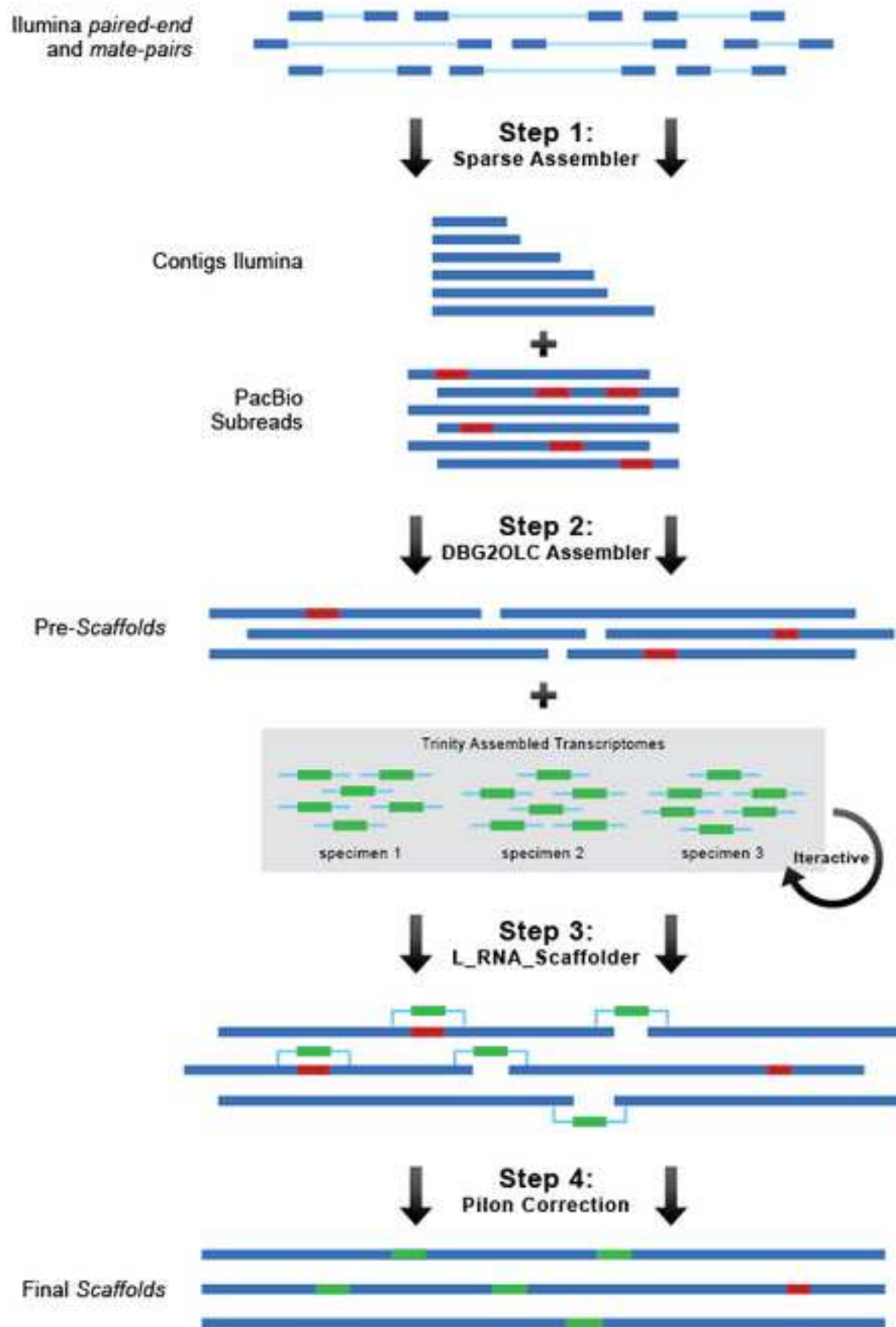
30  
31 609  
32 610 **Figure 3A:** Gene family assigned with OrthoMCL for the total set of proteins predicted  
33 611 from five mussel genome projects. Outside the Venn diagram its represented the species name  
34 612 and below it is the number of proteins / number of clustered proteins / number of clusters. **B:**  
35 613 Phylogeny of the concatenated data set using 44 single-copy orthologs extracted from ten  
36 614 molluscan genomes. The VT model was estimated to be best fitting substitution model with  
37 615 ProtTest 3.4.2. We reconstructed the phylogeny using PhyML and 100 bootstrap repetition.

38 616  
39 617  
40 618 **Figure 4: Gene family representation analysis in the *L. fortunei* genome. Panel A.**  
41 619 **PFAM hierarchical clustering, heatmap.** Features were selected according to a model based on  
42 620 the Poisson cumulative distribution of each PFAM count in the golden mussel genome vs the  
43 621 normalized average values found in the other nine molluscan genomes (Bonferroni correction,  $P$   
44 622  $\leq 0.05$ ). Transposable elements were included in the analysis. Colors depict the log2 ratio  
45 623 between PFAM counts found in each single genome and the corresponding mean value. The  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

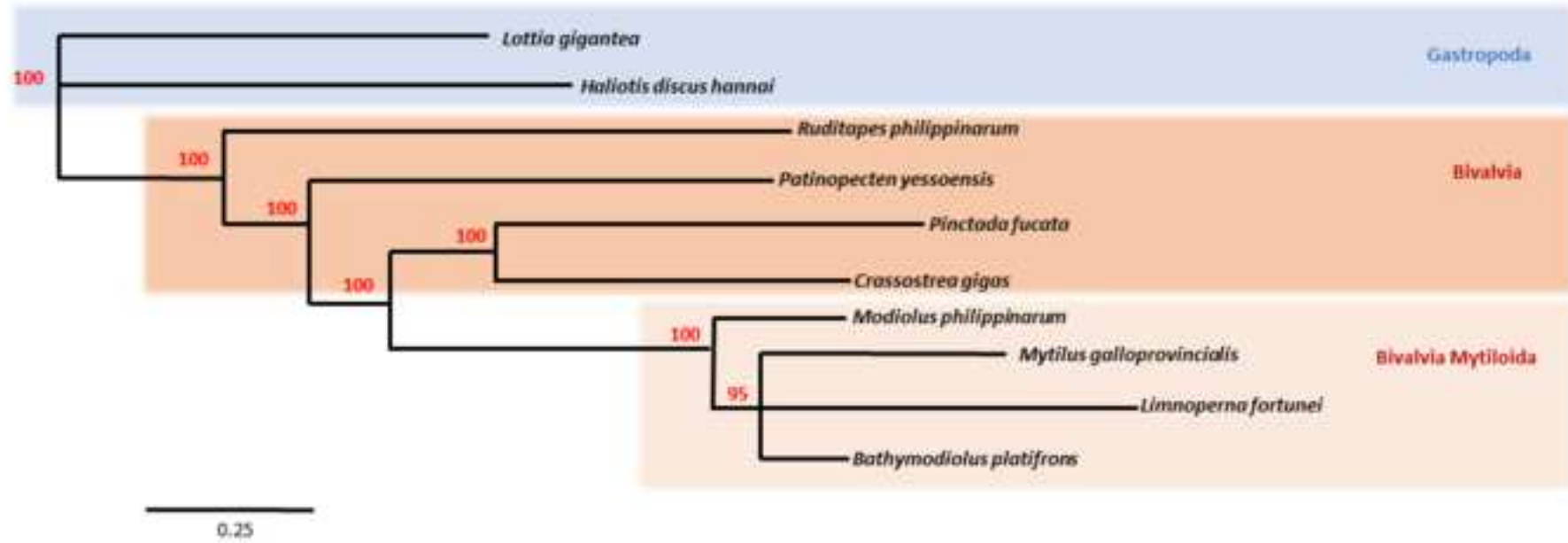
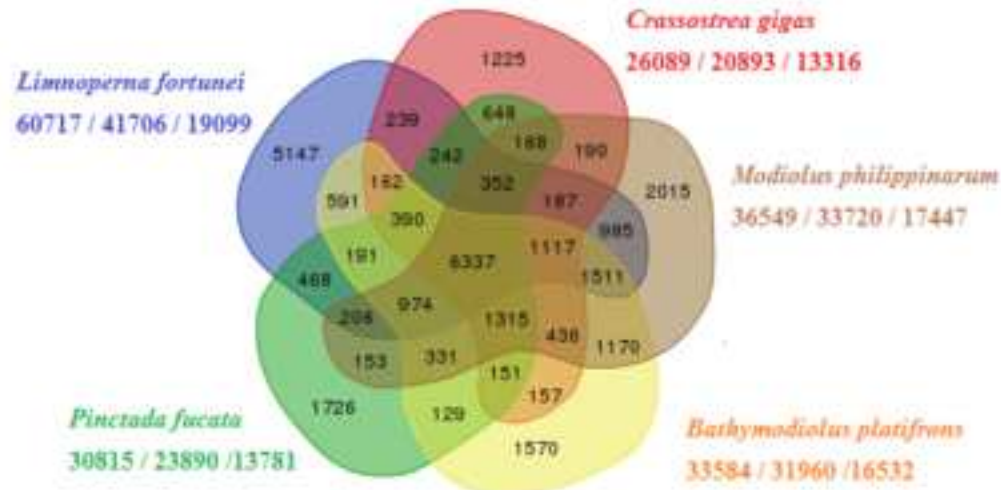
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  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

branching. Legend: Lf, *L. fortunei*; Bp, *Bathymodiolus platifrons*; Mg, *Mytilus galloprovincialis*; Mp, *Modioulus philippinarum*; Cg, *Crassostrea gigas*; Pf, *Pinctada fucata*; Py, *Patinopecten yessoensis*; Rp, *Ruditapes philippinarum*; Hd, *Haliotis discus hannai*; Lg, *Lottia gigantea* **Panel B. Gene ontology analysis of expanded gene families (PFAMs), semantic scatter plot.** Shown are cluster representatives after redundancy reduction in a two-dimensional space applying multidimensional scaling to a matrix of semantic similarities of GO term. Color indicates the GO enrichment level (legend in upper left-hand corner); size indicates the relative frequency of each term in the UNIPROT database (larger bubbles represent less specific processes).

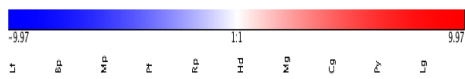




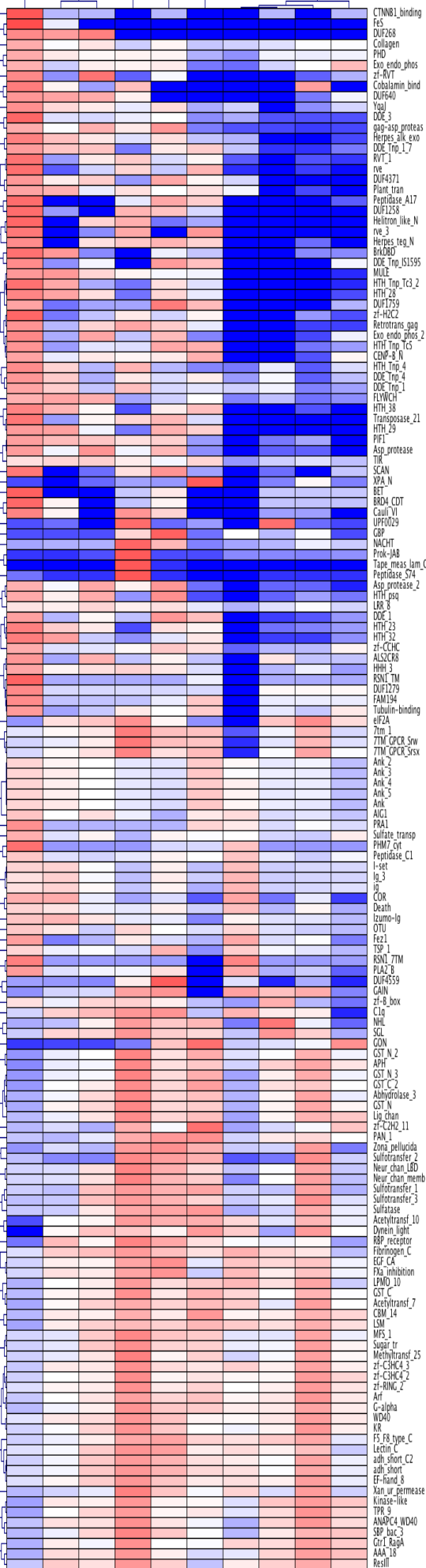




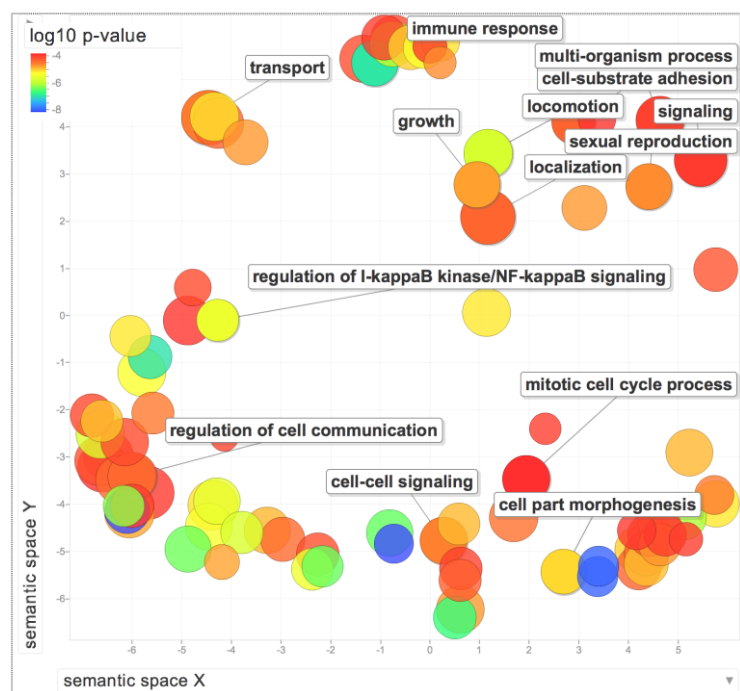
A




Li Bp Mp Pf Rp Hd Ma Cp Py Lp




B





Click here to access/download  
**Supplementary Material**  
TableS1.docx

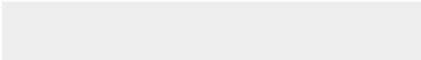





Click here to access/download  
**Supplementary Material**  
Figure-S1.png

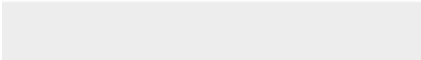



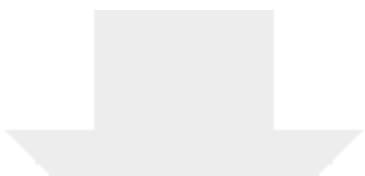
Click here to access/download  
**Supplementary Material**  
TableS2.docx






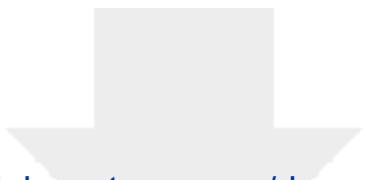
Click here to access/download  
**Supplementary Material**  
TableS3-revised.docx






Click here to access/download  
**Supplementary Material**  
Table-S4-final.xlsx

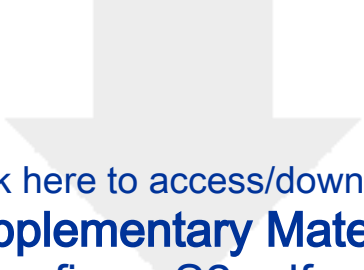




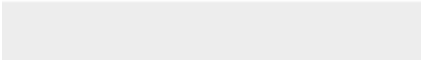
Click here to access/download  
**Supplementary Material**  
TableS5.xlsx








Click here to access/download  
**Supplementary Material**  
figureS2.pdf





Click here to access/download  
**Supplementary Material**  
tableS6-.docx

