

A Chromosomal-Level Genome Assembly for the insect vector for Chagas disease, *Triatoma rubrofasciata* --Manuscript Draft--

Manuscript Number:	GIGA-D-19-00028R2	
Full Title:	A Chromosomal-Level Genome Assembly for the insect vector for Chagas disease, <i>Triatoma rubrofasciata</i>	
Article Type:	Data Note	
Funding Information:	Foundation for the Development of Science and Technology Museums in China (Grant No. 2016YFC1202000)	Prof. Xiao-Nong Zhou
Abstract:	<p>Background: <i>Triatoma rubrofasciata</i> is a widespread pathogen vector for Chagas disease, an illness that affects approximately seven million people worldwide. Despite of its importance to human health, its evolutionary origin has not been conclusively determined. A reference genome for <i>T. rubrofasciata</i> is not yet available.</p> <p>Finding: We have sequenced the genome of a female <i>T. rubrofasciata</i> individual using a single molecular DNA sequencing technology (i.e., PacBio Sequel platform) and have successfully reconstructed a whole-genome (680 Mb) assembly that covers 90% of the nuclear genome (757 Mb). Through Hi-C analysis, we have reconstructed full-length chromosomes of this female individual that has 13 unique chromosomes ($2n = 24 = 22 + X1 + X2$) with a contig N50 of 2.72Mb and a scaffold N50 of 50.7 Mb. This genome has achieved a high base-level accuracy of 99.99%. This platinum-grade genome assembly has 12,691 annotated protein-coding genes. More than 95.1% BUSCO genes were single-copy completed, indicating a high level of completeness of the genome.</p> <p>Conclusion: The platinum-grade genome assembly and its annotation provide valuable information for future in-depth comparative genomics studies including sexual determination analysis in <i>T. rubrofasciata</i> and the pathogenesis of Chagas disease.</p>	
Corresponding Author:	Xiao-Nong Zhou National Institute of Parasitic Diseases Shanghai, CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	National Institute of Parasitic Diseases	
Corresponding Author's Secondary Institution:		
First Author:	Qin Liu, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Qin Liu, Ph.D. Yunhai Guo Yi Zhang Wei Hu Yuanyuan Li Dan Zhu Zhengbin Zhou Jiatong Wu	

	Lansheng Chen
	Nansheng Chen
	Xiao-Nong Zhou
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Point-to-point responses to Editors of GigaScience</p> <p>Dear Editors:</p> <p>Thank you very much.</p> <p>We have read all comments word by word, along with those corrections in the edited manuscript. The suggestions had been accepted and amended carefully in this new version. All questions had been answered in this point-to-point response, and our response to each comment was written as follows, following each comment in BLUE.</p> <p>We should be appreciated if you could take the revised version consideration to be published in GigaScience.</p> <p>Sincerely yours,</p> <p>Qin Liu (First Author) Xiao-Nong Zhou (Corresponding author)</p> <p>Replies to comments: In addition, please register any new software application in the SciCrunch.org database to receive a RRID (Research Resource Identification Initiative ID) number, and include this in your manuscript. This will facilitate tracking, reproducibility and re-use of your tool. Re : Thank you. No new software application was needed to register in this manuscript. The reference of the GCE software was inserted in Page 4 line 13.</p> <p>Reviewer reports: Reviewer #1: Most reviewer issues were addressed satisfactorily. Some minor issues is suggested below.</p> <p>Page 4, line 2: ... distribution/ratio is showed in Figure 2. Re : Thank you. "showing" was changed to "showed" in Page4, line2.</p> <p>Figure 3: Unit of X-axes? How does this relate to estimated genome size? Re : Thank you. The name and unit of X-axes in Figure 3 was changed in the new version. The genome size of <i>T. rubrofasciata</i> was estimated by using the Kmer-based method in GCE software (Liu B , Shi Y , Yuan J , et al. Estimation of genomic characteristics by analyzing k-mer frequency in de novo genome projects. <i>Quantitative Biology</i>, 2013, 35(s 1–3):62-67). We calculated and plotted the 17-mer depth distribution in Figure 3. The X-axes was the kmer count, means the peak frequency of 17-mers. The peak frequency was estimated around 41 and the genome size of <i>T. rubrofasciata</i> was estimated to be 757 Mb on the basis of the formula "G = N17-mer/D17-mer", where the N17-mer was the number of 17-mers, D17-mer denoted the peak depth of 17-mers estimated, and G represented the estimated genome size.</p> <p>Page 5, line 28: Please note that the last sentence of this paragraph still refer to the difficulty of the mollusk genome assembly. This paragraph also contained numerous small errors. Please see suggested modified paragraph below:</p>

	<p>Re : Thank you. This paragraph was changed as the reviewer's suggestion.</p> <p>Reviewer #2: The authors have sufficiently addressed all of my concerns. Re : Thank you.</p> <p>Additional: The pictures of female <i>T. rubrofasciata</i> were changed in Figure 1.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>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.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<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</p>	Yes

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](#)?

[Click here to view linked References](#)

1 A Chromosomal-Level Genome Assembly for the insect vector
2 for Chagas disease, *Triatoma rubrofasciata*

3 Qin Liu^{1#}, Yunhai Guo^{1#}, Yi Zhang^{1#}, Wei Hu^{1,2}, Yuanyuan Li¹, Dan Zhu¹, Zhengbin
4 Zhou¹, Jiatong Wu¹, Nansheng Chen^{3,4,5*}, Xiao-Nong Zhou^{1*}

5 ¹National Institute of Parasitic Diseases, Chinese Center for Disease Control and Prevention; Key
6 Laboratory of Parasite and Vector Biology, Ministry of Health; WHO Collaborating Center for
7 Tropical Diseases; Chinese Center for Tropical Diseases Research, Shanghai 200025, P. R. China

8 ²Department of Microbiology and Microbial Engineering, School of Life Sciences, Fudan,
9 Shanghai 200025, P. R. China

10 ³CAS Key laboratory of Marine Ecology and Environmental Sciences, Institute of Oceanology,
11 Chinese Academy of Sciences; Qingdao, Shandong 266071, China;

12 ⁴Laboratory for Marine Ecology and Environmental Science, Qingdao National Laboratory for
13 Marine Science and Technology, Qingdao, Shandong 266237, China;

14 ⁵Department of Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, Canada

15

16 **Abstract**

17 **Background:**

18 *Triatoma rubrofasciata* is a widespread pathogen vector for Chagas disease, an illness that affects
19 approximately seven million people worldwide. Despite of its importance to human health, its
20 evolutionary origin has not been conclusively determined. A reference genome for *T. rubrofasciata*
21 is not yet available.

22 **Finding:**

23 We have sequenced the genome of a female *T. rubrofasciata* individual using a single molecular
24 DNA sequencing technology (i.e., PacBio Sequel platform) and have successfully reconstructed a
25 whole-genome (680 Mb) assembly that covers 90% of the nuclear genome (757 Mb). Through
26 Hi-C analysis, we have reconstructed full-length chromosomes of this female individual that has
27 13 unique chromosomes ($2n = 24 = 22 + X1 + X2$) with a contig N50 of 2.72Mb and a scaffold
28 N50 of 50.7 Mb. This genome has achieved a high base-level accuracy of 99.99%. This
29 platinum-grade genome assembly has 12,691 annotated protein-coding genes. More than 95.1%

1 BUSCO genes were single-copy completed, indicating a high level of completeness of the
2 genome.

3 **Conclusion:**

4 The platinum-grade genome assembly and its annotation provide valuable information for future
5 in-depth comparative genomics studies including sexual determination analysis in *T. rubrofasciata*
6 and the pathogenesis of Chagas disease.

7

8 **Key Words:** *Triatoma rubrofasciata*, PacBio Sequel platform, Hi-C, chromosomal-level
9 assembly, comparative genomics, RNA-Seq, Iso-Seq

1 **Data description**

2 **Introduction**

3 The insect *T. rubrofasciata* (De Geer) (Hemiptera, Triatominae) is the first Triatominae species
4 formally described, initially with the name *Cimex rubrofasciatus* De Geer, 1773 [1]. This insect
5 presents anthropogenic habits with its dispersion favored by the interaction between residential
6 settlement and human activities [2]. It is considered of global epidemiological importance since it
7 has a pantropical widespread distribution which is found in approximately 45 countries from the
8 Old World to the New World [3]. It is one of the 151 species of Triatominae that has 18 genera
9 currently described worldwide that can transmit American trypanosomiasis known as Chagas
10 disease [4]. This condition has great impact on public health, with 7-8 million people estimated to
11 be infected worldwide, mostly in Latin America. It has become a global health issue in this
12 century with the spread to the non-endemic countries due to growing population movements [5].

13 Due to growing population movements, important epidemiological changes have occurred in
14 recent decades, and the disease has now spread to many non-endemic countries [6]. The
15 widespread of *T. rubrofasciata* emerges as a potential risk of outbreaks in these regions, which
16 demands urgent studies through comprehensive sampling and comparative studies. The lack of a
17 high-quality reference genome represents a major hurdle for such efforts. Here, we present a
18 platinum-grade reference genome for *T. rubrofasciata*, which will be valuable for developing
19 vector control programs.

20 **Sample description and DNA sequencing**

21 An adult female insect *T. rubrofasciata* (Figure 1) was used for reference genome construction in
22 this study. This insect was the second generation offspring of a population that was established
23 from the eggs of single female adult collected in Shunde County, Foshan City, Guangdong
24 Province (22°42'44.63"N, 113°08'45.34"E), China, in 2016 [7]. DNA was extracted from this
25 individual using the traditional phenol/chloroform extraction method and was quality checked
26 using agarose gel electrophoresis. A single band was observed, indicating high integrity of DNA
27 molecules for library construction for the Illumina X Ten (Illumina Inc., San Diego, CA, USA)
28 and the PacBio Sequel (Pacific Biosciences of California, Menlo Park, CA, USA) sequencing
29 platforms.

30 Using the DNA preparation, a library with the insertion length of 350 bp was constructed for
31 Illumina sequencing platform according to the manufacturer's protocol. 46.75 Gb short reads were
32 obtained from the Illumina X Ten DNA sequencing platform (Table 1). 39.32 Gb filtered reads
33 were used for the following genome survey analysis, and for final-stage base-level genome
34 sequence polishing. Meanwhile, 20 Kb-libraries were constructed for PacBio Sequel sequencing.
35 Using fourteen SMRT cells, 8.23 million reads were generated, with the total length of 69.38 Gb

1 (Table 1). The mean length of these subreads was 8.43 Kb and the plot of the read length
2 distribution/ratio was showed in Figure 2.

3 **Genome features estimation through Kmer analysis**

4 With sequencing data from the Illumina HiSeq X Ten DNA sequencing platform, several genome
5 features were evaluated for the genome of *T. rubrofasciata*. To ensure the quality of the analysis,
6 ambiguous bases and low-quality reads were first trimmed and filtered using the HTQC package
7 [8]. First, the quality of bases at two read ends was checked. Bases in sliding 5 bp windows were
8 deleted if the average quality of the window was below 20. Second, reads were filtered if the
9 average quality were smaller than 20 or the read length was shorter than 75 bp. Third, the mate
10 reads were also removed if the corresponding reads were filtered.

11 The processed reads were used for genome assessment. We calculated the number of each
12 17-mer from the sequencing data using the jellyfish software (v2.1.3) [9], and the distribution was
13 analyzed with GCE software [10]. We estimated the genome size of 757 Mb with the
14 heterozygosity of 1.01% and repeat content of 55.49% in the genome. Kmer analysis was using to
15 estimate the genome size which showed the PacBio assembly was of good quality (Figure 3). The
16 genome size of *T. rubrofasciata* is similar to that of *Rhodnius prolixus*, another insect vector of
17 Chagas disease, which has a predicted 733 Mb genome size [11].

18

19 **Genome assembly using PacBio long reads**

20 FALCON [12] was employed using the length_cut_off and length_cutoff_pr parameters of 3 Kb
21 and 3 Kb, respectively. We first obtained 677.72 Mb genome with a contig N50 of 2.71 Mb. The
22 genome sequences were subsequently polished using PacBio long reads using arrow [13] and
23 Illumina short reads by pilon [14] to correct base errors.

24

25 ***In situ* Hi-C library construction and chromosome assembly using Hi-C data**

26 A separate female individual *T. rubrofasciata* was used for library construction for Hi-C analysis
27 as described previously [15, 16]. The library was sequenced with 150 bp paired-end mode on the
28 Illumina HiSeq X Ten platform (San Diego, CA, United States).

29 From the Illumina HiSeq X Ten platform, 103.61 Gb reads were obtained for the Hi-C library
30 and 99.28 Gb filtered reads were used for the following Hi-C analysis. The reads were mapped to
31 the above *T. rubrofasciata* genome with Bowtie [17], with both ends of paired reads being mapped
32 to the genome separately. To increase the interactive Hi-C reads ratio, an iterative mapping
33 strategy was performed as previous studies, and only read pairs that both ends uniquely mapped
34 were used for the following analysis. From the alignment of the paired ends, self-ligation,
35 non-ligation and other sorts of invalid reads, including StartNearRsite, PCR amplification, random

1 break, LargeSmallFragments and ExtremeFragments, were filtered out by Hi-C lib and the method
2 was described in a previous study [15]. Through the recognition of restriction sites in sequences,
3 contact counts among contigs were calculated and normalized.

4 By clustering the contigs using the contig contact frequency matrix, we were able to correct
5 some minor errors in the FALCON assembly results. Contigs with errors were corrected by
6 breaking into shorter contigs, we obtained a chromosome-level genome assembly of 680.73 Mb
7 with 2,126 contigs, and a contig N50 of 2.72 Mb. The longest contig was 10.27 Mb in size (Table
8 2). Among these 2,126 contigs, 626 contigs were mounted to 13 chromosomes with Lachesis [18]
9 using the agglomerative hierarchical clustering method. Lachesis was further applied to order and
10 orient the clustered contigs according to the contact matrix. Contigs anchored to chromosomes
11 accounted for 92.51% of the total genome bases (Figure 4). The number of chromosomes matched
12 nicely to previously published karyotype of a female *T. rubrofasciata* individual ($2n = 24 = 11 * 2$
13 $+ X1 + X2$) [1]. Taken together, we have successfully reconstructed the first chromosomal-level
14 assembly of *T. rubrofasciata* of 680.73 Mb, with 2,126 contigs, a contig N50 of 2.72 Mb, a
15 scaffold N50 of 50.70 Mb (Table 2).

16 **Genome quality evaluation**

17 We assessed the quality of genome of *T. rubrofasciata* in three aspects: sequence continuity,
18 genome completeness and base level accuracy.

19 First of all, we compared the contig/scaffold number and N50 length of contig of *T.*
20 *rubrofasciata* with insect species with sequenced genomes and found that our assembly has much
21 improved quality over other insects (Figure 5). We attributed the improvement to the application
22 of the PacBio long reads for genome assembly. With Hi-C data analysis, we successfully
23 assembled *T. rubrofasciata* genome to chromosome-level with just one individual. Like previous
24 studies, insect genome heterozygosity was one of the biggest challenges for genome assembly,
25 both in terms of contig and scaffold assembly. Traditional chromosomal genome assembly
26 requires physical maps and genetic maps, which is enormously time and labor-consuming. Our
27 work illustrated that the genome assembly using PacBio long sequencing data was not only
28 affordable but also effective for overcoming the difficulties presented by insect genome assembly.

29 Second, the assembled genome was subjected to the BUSCO v.3.0.2 (Benchmarking
30 Universal Single-Copy Orthologs, RRID:SCR_015008) [19] to assess the completeness of the
31 genome assembly. We used “insect_obd9” gene set. 98.2% of the BUSCO genes were identified in
32 *T. rubrofasciata* genome. More than 95.1% BUSCO gene were single-copy completed in our
33 genome, illuminating a high level of completeness of the genome.

34 Third, NGS short reads were aligned to the genome using BWA [20]. About 98.1% of reads
35 were aligned to the genome, of which 98.0% were reads paired aligned. The insertion length
36 distribution of read pairs exhibited a single peak around 300 bp, which was consistent with the

1 design for the Illumina sequencing library construction. Note that the NGS data, which was used
2 for error correction, was not used in contig assembly. Therefore, the insertion length distribution
3 of NGS data illustrated the high quality of our assembly at the contig level. From the NGS reads
4 alignment, we detected 8,478 homologous SNP loci using GATK [21], demonstrating the high
5 base-level accuracy of 99.99%.

6 **Repeat element and gene annotation**

7 Tandem Repeat Finder (TRF) [22] was used for repetitive element identification in *T.*
8 *rubrofasciata* genome. A *de novo* method applying RepeatModuler
9 (<http://www.repeatmasker.org/RepeatModeler.html>) was used to detect transposable elements
10 (TEs). The resulting *de novo* data, combined with known repeat library from Repbase [23], were
11 used to identify TEs in the *T. rubrofasciata* genome by RepeatMasker [24].

12 Protein-coding genes in the *T. rubrofasciata* genome were annotated using the *de novo*
13 program Augustus (RRID:SCR_008417) [25]. Protein sequences of the closely related species
14 including *Rhodnius prolixus* (from VectorBase), *Halyomorpha halys* (from NCBI), *Oncopeltus*
15 *fasciatus* (from USDA), *Cimex lectularius* (from NCBI), and *Drosophila melanogaster* (from
16 NCBI), were aligned to the *T. rubrofasciata* genome with tblastn. Full-length transcripts obtained
17 using Iso-Seq were mapped to the genome using Gmap [26]. Finally, gene models predicted from
18 all above methods were combined by MAKER [27], resulting in 12,691 protein-coding genes. The
19 gene number, gene length, CDS length, exon length and intron length distribution were all
20 comparable with the related insects (Figure 6).

21 To functionally annotate protein-coding genes in the *T. rubrofasciata* genome, we searched
22 all predicted gene sequences to NCBI non-redundant protein (NR), InterPro (InterProScan,
23 RRID:SCR_005829) [28], GO (Gene Ontology), KEGG (RRID:SCR_012773) [29], Swissprot
24 [30], TrEMBL databases [30] by BLASTN [31] and BLASTX [32]. A threshold of e-value of 1e-5
25 was used for all BLAST applications. Finally, 12,063 genes were functionally annotated (Table 3).

26 **Phylogenetic analysis of *T. rubrofasciata* with other insects**

27 OrthMCL was used to cluster gene families. First, proteins from *T. rubrofasciata* and the closely
28 related insects, including *Rhodnius prolixus*, *Oncopeltus fasciatus*, *Halyomorpha halys*, *Cimex*
29 *lectularius*, *Drosophila melanogaster*, *Gerris buenoi*, *Homalodisca vitripennis*, *Acyrtosiphon*
30 *pisum*, *Culex quinquefasciatus*, *Glossina palpalis*, *Apis mellifera* and *Heliconius melpomene* were
31 all-to-all blasted by BLASTP [32] utility with an e-value threshold of 1e-5. Only proteins from the
32 longest transcript were used for genes with alternative splices. We identified 21,850 gene families
33 for *T. rubrofasciata* and the related species, among them 330 single-copy orthologs families.

34 Using single-copy orthologs, we probed the phylogenetic relationships for the *T.*
35 *rubrofasciata* and other insects. To this end, protein sequences of single-copy genes were aligned

1 using MUSCLE [33]. Guided by the protein multi-sequence alignment, the alignment of the
2 coding DNA sequences (CDS) for those genes were generated and concatenated for the following
3 analysis. The phylogenetic relationships were constructed using PhyML [34] using the
4 concatenated nucleotide alignment with the JTT+G+F model. We first obtained divergent times
5 for all pair using the phylogenetic tree using r8s [35], which were used as input, together with
6 molecular clock data from the divergence time from the TimeTree database [36], to estimate
7 species divergence time for all pairs of species in the phylogenetic tree using MCMCtree program
8 (from PAML) [37]. We found that *T. rubrofasciata* was most closely related to *R. prolixus*, and the
9 two species diverged from their common ancestor around 60.00-95.00 million years ago (MYA)
10 (Figure 7).

11 **Conclusion**

12 We reconstructed the first high-quality, chromosome-level assembly of *T. rubrofasciata* using an
13 integrated strategy of PacBio, Illumina and Hi-C technologies. Using the long reads from PacBio
14 Sequel platform and short reads from the Illumina HiSeq X Ten platform, we successfully
15 constructed contig assembly for *Triatoma*. Leveraging contact information among contigs from
16 Hi-C technology, we further improved the assembly to the chromosome-level quality. We
17 annotated 12,691 protein-coding genes in the *T. rubrofasciata* genome, 12,063 of which were
18 functionally annotated. With 330 single-copy orthologs from *T. rubrofasciata* and other related
19 insects, we construct the phylogenetic relationship of these insects, and found that *T. rubrofasciata*
20 might have diverged from its common ancestor of *R. prolixus* around 60.00-95.00 MYA. Given
21 the increasing interests in insect genome evolution and the biological importance of *T.*
22 *rubrofasciata* as the vector for Chagas disease, our genomic and transcriptome data provide
23 valuable genetic resource for the following functional genomics investigations for the research
24 community.

25 **Availability of supporting data**

26 The raw data from our genome project was deposited in the NCBI Sequence database with
27 Bioproject IDs PRJNA516044. The Illumina, PacBio and Hi-C sequencing data are available from
28 NCBI via the accession number of SRR8466736, SRR8466737 and SRR8466756, respectively.
29 The Illumina transcriptome sequencing data were deposited to NCBI via the accession number of
30 SRR8468315 and SRR8468316. Other data further supporting this work are available in the
31 GigaScience repository, GigaDB [38].

32 **Ethics Statement**

33 This study was approved by the Animal Care and Use committee of National Institute of Parasitic
34 Diseases, Chinese Center for Disease Control and Prevention. All participants consent the study
35 under the 'Ethics, consent and permissions' heading. All participants consent to publish the work

1 under the 'Consent to publish' heading.

2 Competing interests

3 The authors declare that they have no competing interests.

4 Funding

5 This work was supported by the National Key Research and Development Program of China
6 (Grant No. 2016YFC1202000), the National Science and Technology Project (No. 201810101002)
7 and the CAS Pioneer Hundred Talents Program (to N.S.C.) and Taishan Scholar Project Special
8 Fund (to N.S.C.).

9 Author Contributions

10 Z.X.N., L.Q., Z.Y. and H.W. conceived the project. L.Q., G.Y.H., Z.Y., Z.D., L.Y.Y., W.J.T. and
11 Z.Z.B. collected the samples and extracted the DNA and RNA. L.Q., G.Y.H., Z.Y. performed the
12 genome assembly and data analysis. C.N.S. performed the data analysis. L.Q and C.N.S. wrote the
13 paper. Z.X.N. revised the manuscript. All authors read, edited and approved the final version of
14 the manuscript.

15 Acknowledgements

16 We thank for Frasergen Bioinformatics for providing technical support for this work.

17

18 References

- 19 1. Alevi, K.C.C., et al., *Cytogenetic Characterisation of Triatoma rubrofasciata (De Geer)*
20 *(Hemiptera, Triatominae) Spermatocytes and Its Cytotaxonomic Application*. African
21 Entomology, 2016. **24**(1): p. 4.
- 22 2. Hypsa, V., et al., *Phylogeny and biogeography of Triatominae (Hemiptera: Reduviidae):*
23 *molecular evidence of a New World origin of the Asiatic clade*. Mol Phylogenet Evol, 2002.
24 **23**(3): p. 447-57.
- 25 3. Galvão, C., et al., *A checklist of the current valid species of the subfamily Triatominae Jeannel,*
26 *1919 (Hemiptera, Reduviidae) and their geographical distribution, with nomenclatural and*
27 *taxonomic notes*. Zootaxa, 2003. **202**: p. 36.
- 28 4. Justi, S.A. and C. GALVÃO, *The Evolutionary Origin of Diversity in Chagas Disease Vectors*.
29 Trends in Parasitology, 2017. **33**(1): p. 11.
- 30 5. Carod-Artal, F.J., *Chapter 7 - American trypanosomiasis*. Handbook of Clinical Neurology, 2013.
31 **114**: p. 120.
- 32 6. Coura, J.R. and P.A. Vinas, *Chagas disease: a new worldwide challenge*. Nature, 2010. **465**: p.
33 2.
- 34 7. Liu, Q., et al., *First records of Triatoma rubrofasciata (De Geer, 1773) (Hemiptera, Reduviidae)*

- 1 *in Foshan, Guangdong Province, Southern China*. Infect Dis Poverty, 2017. **6**(1): p. 129.
- 2 8. Neff, K.L., et al., *Mojo Hand, a TALEN design tool for genome editing applications*. BMC
3 Bioinformatics, 2013. **14**: p. 1.
- 4 9. Marçais, G. and C. Kingsford, *A fast, lock-free approach for efficient parallel counting of
5 occurrences of k-mers*. Bioinformatics, 2011. **27**(6): p. 764-70.
- 6 10. Liu, B., et al., *Estimation of genomic characteristics by analyzing k-mer frequency in de novo
7 genome projects*. Quantitative Biology, 2013. **35**: p. 3.
- 8 11. Mesquita, R.D., et al., *Genome of Rhodnius prolixus, an insect vector of Chagas disease,
9 reveals unique adaptations to hematophagy and parasite infection*. Proc Natl Acad Sci U S A,
10 2015. **112**(48): p. 14936-41.
- 11 12. Chin, C.S., et al., *Phased diploid genome assembly with single-molecule real-time sequencing*.
12 Nat Methods, 2016. **13**(12): p. 1050-1054.
- 13 13. Chin, C.S., et al., *Nonhybrid, finished microbial genome assemblies from long-read SMRT
14 sequencing data*. Nat Methods, 2013. **10**(6): p. 563-9.
- 15 14. Walker, B.J., et al., *Pilon: an integrated tool for comprehensive microbial variant detection
16 and genome assembly improvement*. PLoS One, 2014. **9**(11): p. e112963.
- 17 15. Gong, G., et al., *Chromosomal-level assembly of yellow catfish genome using third-generation
18 DNA sequencing and Hi-C analysis*. Gigascience, 2018. **7**(11).
- 19 16. Xu, S., et al., *A draft genome assembly of the Chinese sillago (Sillago sinica), the first
20 reference genome for Sillaginidae fishes*. Gigascience, 2018. **7**(9).
- 21 17. Langmead, B., et al., *Ultrafast and memory-efficient alignment of short DNA sequences to the
22 human genome*. Genome Biol, 2009. **10**(3): p. R25.
- 23 18. Near, T.J., et al., *Phylogeny and tempo of diversification in the superradiation of spiny-rayed
24 fishes*. Proceedings of the National Academy of Sciences of the United States of America,
25 2013. **110**(31): p. 12738.
- 26 19. Simao, F.A., et al., *BUSCO: assessing genome assembly and annotation completeness with
27 single-copy orthologs*. Bioinformatics, 2015. **31**(19): p. 3210-2.
- 28 20. Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows-Wheeler transform*.
29 Bioinformatics, 2009. **25**(14): p. 1754-60.
- 30 21. McKenna, A., et al., *The Genome Analysis Toolkit: a MapReduce framework for analyzing
31 next-generation DNA sequencing data*. Genome Res, 2010. **20**(9): p. 1297-303.
- 32 22. Benson, G., *Tandem repeats finder: a program to analyze DNA sequences*. Nucleic Acids Res,
33 1999. **27**(2): p. 573-80.
- 34 23. Bao, W., K.K. Kojima, and O. Kohany, *Repbase Update, a database of repetitive elements in
35 eukaryotic genomes*. Mob DNA, 2015. **6**: p. 11.
- 36 24. Chen, N., *Using RepeatMasker to identify repetitive elements in genomic sequences*, in
37 *Current Protocols in Bioinformatics* 2004. p. 14.
- 38 25. Stanke, M., et al., *AUGUSTUS: ab initio prediction of alternative transcripts*. Nucleic Acids Res,
39 2006. **34**(Web Server issue): p. W435-9.
- 40 26. Wu, T.D., et al., *GMAP and GSNAP for Genomic Sequence Alignment: Enhancements to Speed,
41 Accuracy, and Functionality*. Methods Mol Biol, 2016. **1418**: p. 283-334.
- 42 27. Campbell, M.S., et al., *Genome Annotation and Curation Using MAKER and MAKER-P*. Curr
43 Protoc Bioinformatics, 2014. **48**: p. 4 11 1-39.
- 44 28. Quevillon, E., et al., *InterProScan: protein domains identifier*. Nucleic Acids Res, 2005. **33**(Web

1 Server issue): p. W116-20.

2 29. Kanehisa, M. and S. Goto, *KEGG: kyoto encyclopedia of genes and genomes*. Nucleic Acids Res,
3 2000. **28**(1): p. 27-30.

4 30. Boeckmann, B., et al., *The SWISS-PROT protein knowledgebase and its supplement TrEMBL in*
5 *2003*. Nucleic Acids Res, 2003. **31**(1): p. 365-70.

6 31. Gertz, E.M., et al., *Composition-based statistics and translated nucleotide searches:*
7 *improving the TBLASTN module of BLAST*. BMC Biol, 2006. **4**: p. 41.

8 32. Camacho, C., et al., *BLAST+: architecture and applications*. BMC Bioinformatics, 2009. **10**: p.
9 421.

10 33. Edgar, R.C., *MUSCLE: multiple sequence alignment with high accuracy and high throughput*.
11 Nucleic Acids Res, 2004. **32**(5): p. 1792-7.

12 34. Guindon, S. and O. Gascuel, *A simple, fast, and accurate algorithm to estimate large*
13 *phylogenies by maximum likelihood*. Syst Biol, 2003. **52**(5): p. 696-704.

14 35. Sanderson, M.J., *r8s: inferring absolute rates of molecular evolution and divergence times in*
15 *the absence of a molecular clock*. Bioinformatics, 2003. **19**(2): p. 301-2.

16 36. Kumar, S., et al., *TimeTree: A Resource for Timelines, Timetrees, and Divergence Times*. Mol
17 Biol Evol, 2017. **34**(7): p. 1812-1819.

18 37. Yang, Z., *PAML: a program package for phylogenetic analysis by maximum likelihood*.
19 Computer Applications in the Biosciences, 1997. **13**(5): p. 2.

20 38. Liu Q; Guo YH; Zhang Y; Hu W; Li YY; Zhu D; Zhou ZB; Wu JT; Chen NS; Zhou XN: Supporting
21 data for "A Chromosomal-Level Genome Assembly for the insect vector for Chagas disease,
22 *Triatoma rubrofasciata*" GigaScience Database. 2019. <http://dx.doi.org/10.5524/100614>
23
24

1 Tables and Figures

2 **Tables**

3 **Table 1: Sequencing data generated for *Triatoma rubrofasciata* genome assembly and**
 4 **annotation**

5

Library type	Platform	Library size (bp)	Data size (Gb)	Application
Short reads	HiSeq X Ten	350	46.75	Genome survey and genomic base correction
Long reads	PacBio Sequel	20,000	69.38	Genome assembly
Hi-C	HiSeq X Ten	300-500	103.61	Chromosome construction

6

7 **Table 2: Statistics for genome assembly of *Triatoma rubrofasciata***

8

Sample ID	Length	Scaffold (bp)	Number	Scaffold
	Contig** (bp)		Contig**	
Total	680,314,598	680,726,098	2,126	1,303
Max	10,270,547	97,329,580	-	-
N50	2,722,109	50,700,875	76	6
N60	2,121,675	50,415,845	104	7
N70	1,587,961	46,556,423	140	8
N80	1,038,484	37,928,883	193	10
N90	338,786	20,341,594	301	12

9

10 **Table 3: Statistics for genome annotation of *Triatoma rubrofasciata***

11

Database	Number	Percent
NR	11,451	90.23
InterPro	9,625	75.84
GO	7,180	56.58
KEGG ALL	10,867	85.63
KEGG KO	6,112	48.16
Swissprot	9,448	74.45
TrEMBL	11,989	94.47
Total	12,063	95.05

12

13

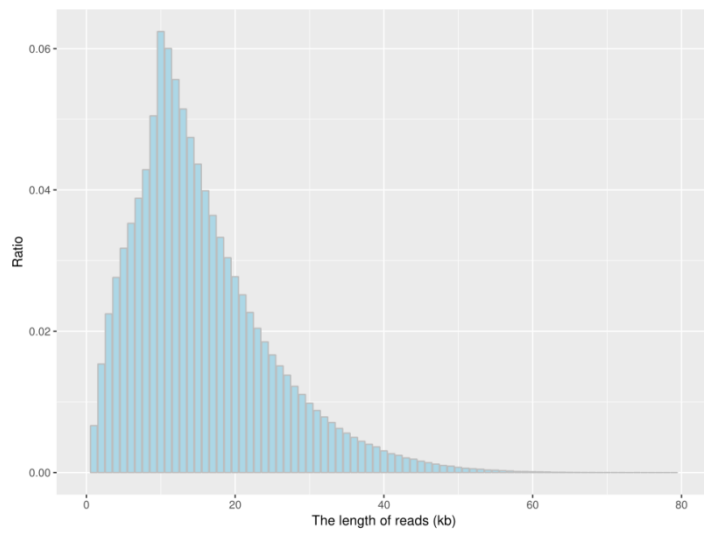
1 **Figures**

2 **Figure 1. Dorsal (left) and ventral (right) views of a female *T. rubrofasciata*.**

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



1 **Figure 2. The plot of the read length distribution/ratio of the subreads.**



2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

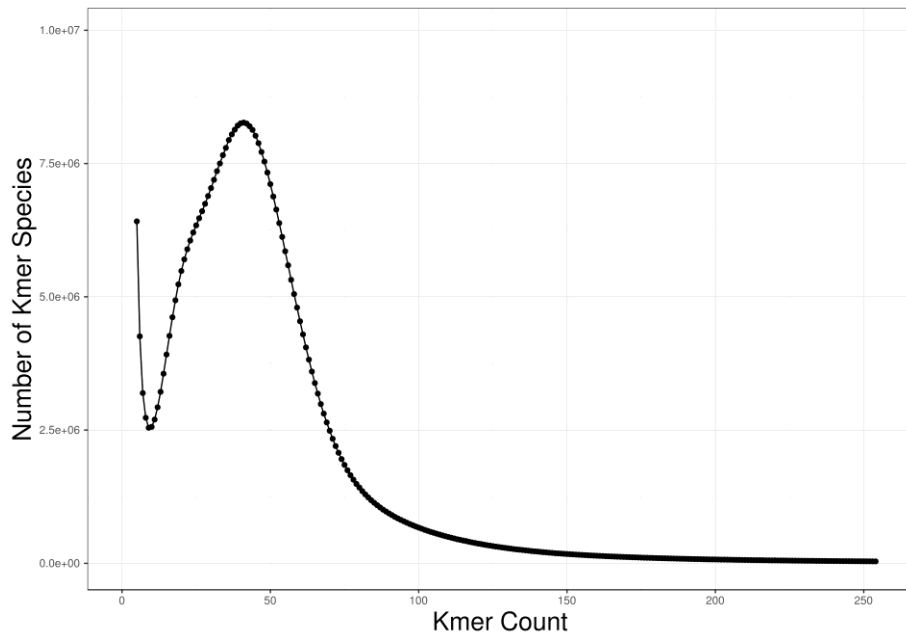
17

18

19

20

1 **Figure 3. 17-mer depth distribution for genome size estimation analysis of *T. rubrofasciata*.**



2

3

4

5

6

7

8

9

10

11

12

13

14

15

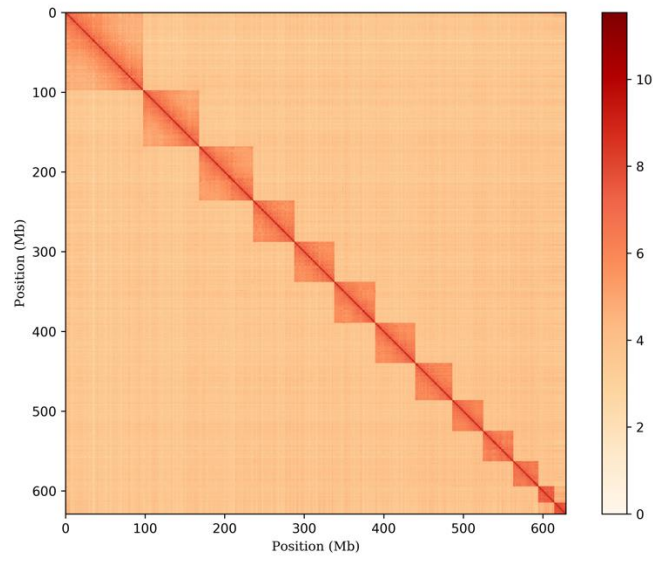
16

17

18

19

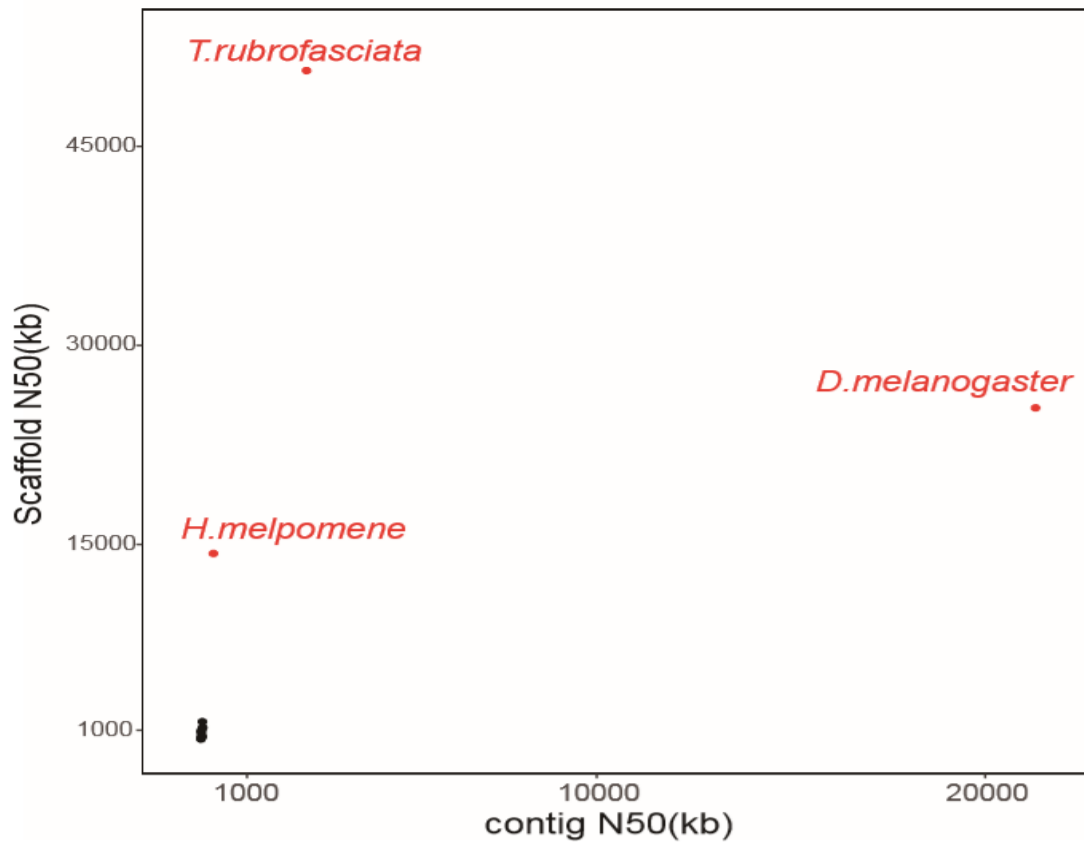
1 **Figure 4. DNA interaction heatmap generated in HiC analysis (resolution: 500 Kb)**



2

3

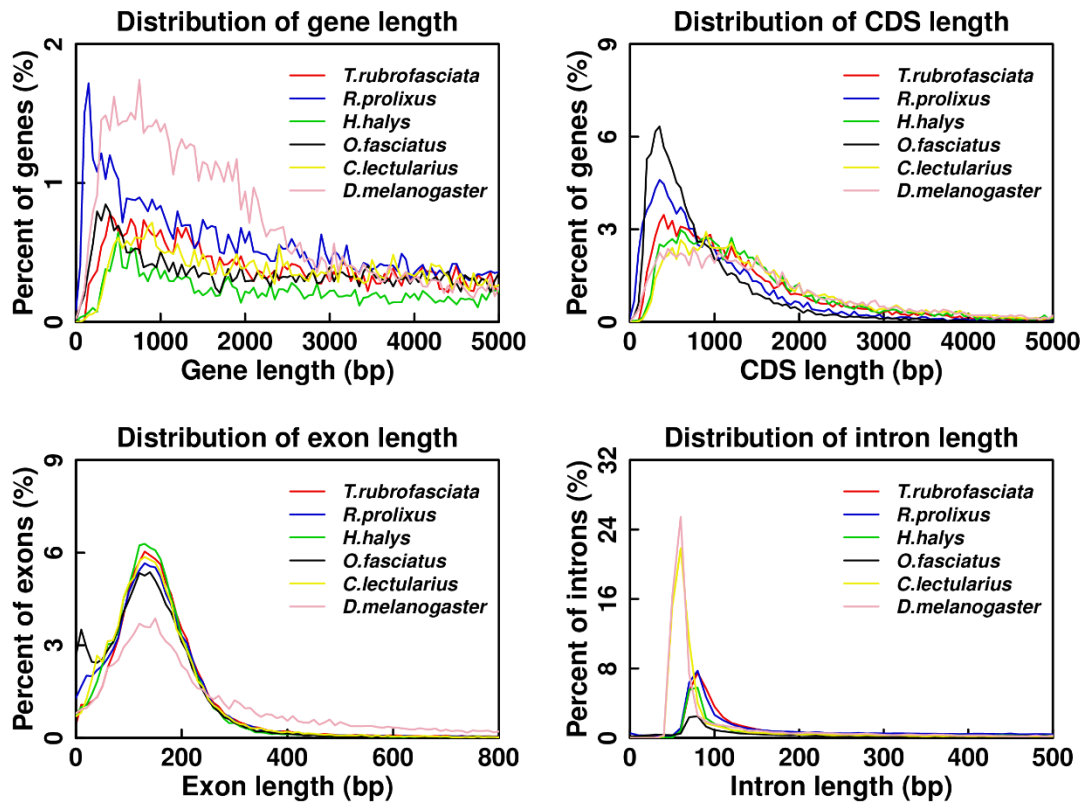
1 **Figure 5: Genome assembly comparison of *T. rubrofasciata* with other sequenced insect genomes**
2 **(*A. mellifera*, *A. pisum*, *C. lectularius*, *C. quinquefasciatus*, *D. melanogaster*, *G. buenoi*, *G. palpalis*,**
3 ***H. halys*, *H. melpomene*, *H. vitripennis*, *O. fasciatus*, *R. prolixus*). The x- and y-axis represent the**
4 **contig and scaffold N50s, respectively. The genomes both contig and scaffold N50s less than 2M**
5 **are highlighted in black.**



6
7
8

1 **Figure 6: Length distribution comparison on total gene, CDS, exon, and intron of annotated gene**
 2 **models of *T. rubrofasciata* with other closely related insect species. Length distribution of total**
 3 **gene (A), CDS (B), exon (C), and intron (D) were compared to those of *R. prolixus*, *H. halys*, *O.***
 4 ***fasciatus*, *C. lectularius* and *D. melanogaster*.**

5



6

7

8

9

10

11

12

13

14

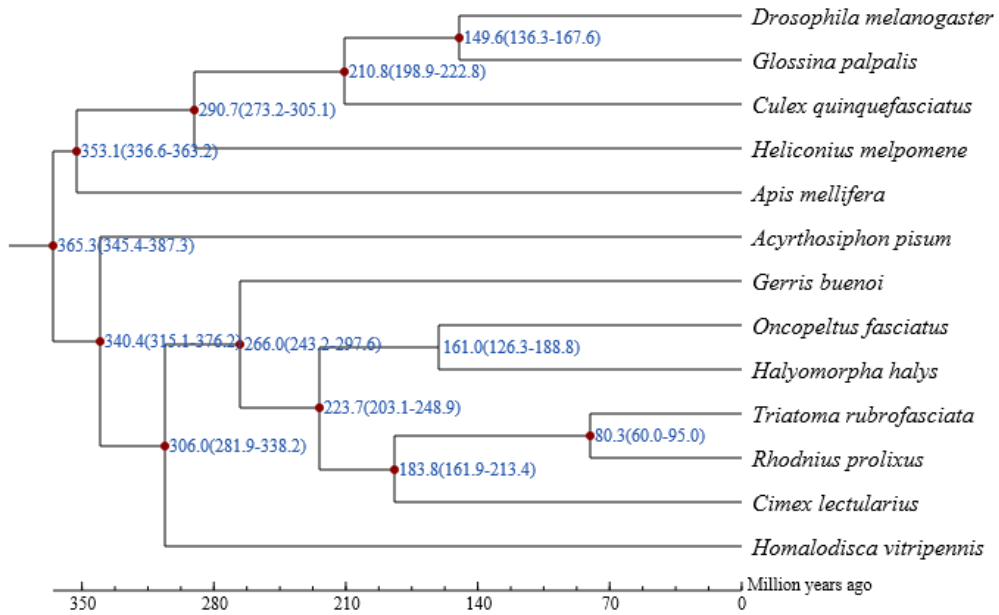
15

16

17

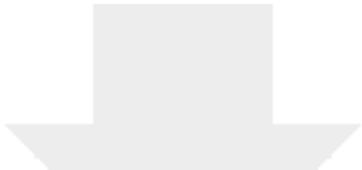
1 **Figure 7: Phylogenetic analysis of *T. rubrofasciata* with other insect species. The estimated species**
 2 **divergence time (million years ago) and the 95% confidential intervals are labeled at each branch**
 3 **site. The divergence used for time recalibration is illuminated as red dots in the tree.**

4

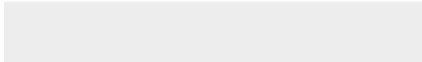


5

6



Click here to access/download
Supplementary Material
N50.xls





中国疾病
预防控制中心

寄生虫病预防控制所

National Institute of Parasitic Disease, Chinese Center For Disease Control and Prevention

207 Ruijin Er Rd., Shanghai
Shanghai 200025, P.R. China
Website: www.ipd.org.cn

Point-to-point responses to Editors of GigaScience

Dear Editors:

Thank you very much.

We have read all comments word by word, along with those corrections in the edited manuscript. The suggestions had been accepted and amended carefully in this new version. All questions had been answered in this point-to-point response, and our response to each comment was written as follows, following each comment in BLUE.

We should be appreciated if you could take the revised version consideration to be published in GigaScience.

Sincerely yours,

Qin Liu (First Author)

Xiao-Nong Zhou (Corresponding author)



中国疾病
预防控制中心

寄生虫病预防控制所

National Institute of Parasitic Disease, Chinese Center For Disease Control and Prevention

207 Ruijin Er Rd., Shanghai
Shanghai 200025, P.R. China
Website: www.ipd.org.cn

Replies to comments:

In addition, please register any new software application in the SciCrunch.org database to receive a RRID (Research Resource Identification Initiative ID) number, and include this in your manuscript. This will facilitate tracking, reproducibility and re-use of your tool.

Re: Thank you. No new software application was needed to register in this manuscript. The reference of the GCE software was inserted in Page 4 line 13.

Reviewer reports:

Reviewer #1: Most reviewer issues were addressed satisfactorily. Some minor issues is suggested below.

Page 4, line 2: ... distribution/ratio is showed in Figure 2.

Re: Thank you. “showing” was changed to “showed” in Page4, line2.

Figure 3: Unit of X-axes? How does this relate to estimated genome size?

Re: Thank you. The name and unit of X-axes in Figure 3 was changed in the new version. The genome size of *T. rubrofasciata* was estimated by using the Kmer-based method in GCE software (Liu B , Shi Y , Yuan J , et al. Estimation of genomic characteristics by analyzing k-mer frequency in de novo genome projects. Quantitative Biology, 2013, 35(s 1–3):62-67). We calculated and plotted the 17-mer depth distribution in Figure 3. The X-axes was the kmer count, means the peak frequency of 17-mers. The peak frequency was estimated around 41 and the genome size of *T. rubrofasciata* was estimated to be 757 Mb on the basis of the formula “G



中国疾病
预防控制中心

寄生虫病预防控制所

National Institute of Parasitic Disease, Chinese Center For Disease Control and Prevention

207 Ruijin Er Rd., Shanghai
Shanghai 200025, P.R. China
Website: www.ipd.org.cn

= $N_{17\text{-mer}}/D_{17\text{-mer}}$, where the $N_{17\text{-mer}}$ was the number of 17-mers, $D_{17\text{-mer}}$ denoted the peak depth of 17-mers estimated, and G represented the estimated genome size.

Page 5, line 28: Please note that the last sentence of this paragraph still refer to the difficulty of the mollusk genome assembly. This paragraph also contained numerous small errors. Please see suggested modified paragraph below:

Re: Thank you. This paragraph was changed as the reviewer's suggestion.

Reviewer #2: The authors have sufficiently addressed all of my concerns.

Re: Thank you.

Additional:

The pictures of female *T. rubrofasciata* were changed in Figure 1.