# GigaScience
# Whole-genome sequence of the oriental lung fluke Paragonimus westermani
## --Manuscript Draft--

| Manuscript Number: | GIGA-D-18-00193R1 | |
|---|---|---|
| Full Title: | Whole-genome sequence of the oriental lung fluke Paragonimus westermani | |
| Article Type: | Data Note | |
| Funding Information: | the QIMR Berghofer Medical Research Institute (Chenhall Estate) | A/Prof Lutz Krause |
| | Australian Infectious Diseases Research Centre | A/Prof Lutz Krause |
| | National Health and Medical Research Council | Prof Donald P. McManus |

| Abstract: | Background |
|---|---|
| | Foodborne infections caused by lung flukes of the genus Paragonimus are a significant and widespread public health problem in tropical areas. Around 50 Paragonimus species have been reported to infect animals and humans, but Paragonimus westermani is responsible for the bulk of human disease. Despite their medical and economic importance, no genome sequence for any Paragonimus species is available. |
| | Results |
| | We sequenced and assembled the genome of P. westermani, which is among the largest of the known pathogen genomes with an estimated size of 1.1 Gb. A 922.8 Mb genome assembly was generated from Illumina and PacBio sequence data, covering 84% of the estimated genome size. The genome has a high proportion (45%) of repeat-derived DNA, particularly of the LINE and LTR subtypes, and the expansion of these elements may explain some of the large size. We predicted 12,852 protein coding genes, showing a high level of conservation with related trematode species. The majority of proteins (80%) had homologs in the human liver fluke Opisthorchis viverrini with an average sequence identity of 64.1%. Assembly of the P. westermani mitochondrial genome from long PacBio reads resulted in a single high-quality circularized 20.6 kb contig. The contig harboured a 6.9 kb region of non-coding repetitive DNA comprised of three distinct repeat units. Our results suggest that the region is highly polymorphic in P. westermani, possibly even within single worm isolates. |
| | Conclusions |
| | The generated assembly represents the first Paragonimus genome sequence and will facilitate future molecular studies of this important, but neglected, parasite group. |

| Corresponding Author: | Lutz Krause |
|---|---|
| | AUSTRALIA |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Harald Oey |
| First Author Secondary Information: | |
| Order of Authors: | Harald Oey |
| | Martha Zakrzewski |
| | Kanwar Narain |
| | K. Rekha Devi |
| | Takeshi Agatsuma |
| | Sujeevi Nawaratna |

| | Geoffrey Gobart |
| --- | --- |
| | Malcolm Jones |
| | Mark A. Ragan |
| | Donald P. McManus |
| | Lutz Krause |
| **Order of Authors Secondary Information:** | |
| **Response to Reviewers:** | Reviewer: 1<br>1 – I suggest a small change in the manuscript title: "Draft Whole genome sequence of the oriental lung fluke…" or just "Whole genome sequence of the …". The term complete for nuclear genome sequence means that it is the final version (in chromosome level with no gaps), not the case here where the genome is still in 30,977 pieces, so complete should be not used here. The mitochondrial indeed looks complete.<br>Response: We have changed the manuscript title to "Whole-genome sequence of the oriental lung fluke Paragonimus westermani" as suggested by the reviewer.<br><br>2 – The authors did not mention how they removed potential contamination or how they maintained the pathogen for the DNA extraction (Please add this information)<br>Response: Comparison of assembled scaffolds with public genome sequence data identified contamination by rat (the experimental host) and the bacterium Delftia sp. All sequences mapped to these genomes were removed. We have now added this information to the methods section.<br><br>3 – Table 1 could be used as supplemental material<br>Response: We intend to submit the manuscript as Data Note. We believe that Table 1 is important for a Data Note and suggest we keep it in the main manuscript.<br><br>4 – The assembly was performed by well-known genome assemblers, but there was any particular reason to not use any of the two most used PacBio assemblers (HGAP and CANU?)<br>Response: We have used CANU for several parasite genomes. The program worked well for other parasites (manuscript in review), but did not perform well on this particular genome. Mira worked better for Paragonimus and generated a single complete mitochondrial contig, whereas CANU resulted in multiple shorter contigs.<br><br>5 – The authors choose to use for the Illumina assembly the ABYSS assembler. From my personal experience and from some colleagues there are several other assemblers that give a better job than ABYSS (Spades, MIRA, Velvet and SoapDenovo2). I know that it varies depending of the nature of the organism and sample used for the assay, but since the group used for the gapfilling step the soapDenovo gapcloser, I would like to see in the manuscript some information about why these pipelines were choosen beside others<br>Response: We have evaluated several assembly programs and ABYSS performed best for this particular genome. ABYSS is also one of the few assemblers that allow inclusion of long-read data to guide scaffolding. The program is still widely used and well maintained. We have an established pipeline using SoapDenovo2, which has been used for the assembly of other parasite genomes (manuscript in review). However, SoapDenovo2 did not perform well for this particular genome, with a large size and many repetitive regions. Additionally, the Paragonimus genome was sequenced from 50 individual worms, resulting in a low-level sequence heterogeneity and assembly of this data proofed to be challenging. ABYSS performed particularly well for the assembly of contigs for this genome. However, the ABYSS gap filler is not well suited for closing gaps larger than 1kb (according to the ABYSS manual and our own experience), whereas the soapDeNovo gapfiller is well suited for this task and performed particularly well on this genome. Additional information has been added to the methods section.<br><br>6 – Line 179 - REAPR typo. I would also suggest the authors to perform for this final polishing genome correction step Pilon or ICORN2 using the Illumina reads generated |

Response: We thank the reviewer for this suggestion. However, Pilon does not seem to perform well for this particular genome. Genome polishing using Pilon with a variety of different settings actually resulted in a slight reduction of BUSCO scores (original assembly: 65.3% complete proteins; after Pilon: 63.9% complete proteins), indicating that Pilon did not improve the overall quality of this particular genome assembly. We manually investigated Pilon results and postulate that Pilon was misled by low-level sequence heterogeneity caused by the pooling of 50 individual worms. As the genome has already been deposited in NCBI and passed all manual QC checks we believe that the questionable improvements by Pilon do not justify re-submission of an updated genome to NCBI.

7 – Please add more information about the genome assembly statistics in table 2 (L50 and number of Ns), a quick run on QUAST should give you this information. And please explain if these gaps are just generated during the scaffolding by the mate pair evidence or it was also generated for unknown size gaps (100Ns). This information is really important to show that some regions could be missing in this draft genome assembly, so future studies could be aware of this fact;
Response: We have run the assembly through QUAST, as requested, and added the L50 to Table 2. The number of Ns can already be inferred from Table 2 as we provide the size of the genome both with and without counting Ns ("Assembly size" and "Total base pairs"). We have re-named "Total base pairs" to "Ungapped size" to make this clearer. The Gaps are generated both during contig assembly (abyss) and scaffolding (SSPACE) and represent the estimated size of the gaps. We have added a sentence to the manuscript to make this clear.

8 – Line 250 - Since the ncRNA information was so important in the mitochondrial annotation, and the group already characterized the tRNAs, please add the method to predict these tRNAs (like tRNAscan) and also, I suggest adding an Aragorn or inferno ncRNA prediction run to improve even more the annotation
Response: The program Mitos, which was used to characterize the mitochondrial genome, identifies both non-coding RNAs and proteins. However, Aragorn was also run to identify any additional tRNAs in the mitochondrial genome (added to methods).

9 – Line 258 - no problem with the methodology, but Cufflinks has a substitute, StringTie (Petera at el., 2015). It will do a much better job to assembly the transcriptome
Response: StringTie was not available when the project started, but we thank the reviewer for this suggestion and will evaluate StringTie for future projects. Cufflinks is well established (>5,000 citations), proven to generate accurate results and is still widely used. We agree that there are many alternative tools that could have been used for transcriptome assembly, but our group has an established and well tested pipeline using cufflinks. We have extensive experience with cufflinks and have optimized the parameters to generate robust and high-quality results. We would further like to point out that we don't publish the assembled cDNA data.

10 – Genome Comparison - I understand that this was not the focus of this manuscript, but sequence identity besides important is a too general comparison method. I suggest add a orthology analysis and maybe generate a Circos synteny plot comparing the new genome with the most similar species available
Response: We will submit the manuscript as Data Note and therefore believe that additional comparative analysis are not required.

11 – Phylogeny - Add a Modeltest run to check if Jones-Taylor-Thornton (JTT) was the best substitution method to be used. For the ML analysis I suggest using PhyML instead of Phylip again, the software used is good but better and newer ones were developed
Response: As suggested, we have now repeated the phylogenetic analysis using PhyML and a model test found the LG substitution model with decorations +G+I+F as optimal. The JTT model was the second best model. PhyML using the LG+G+I+F

model resulted in exactly the same tree topology as our previous analysis using Phylip with the JTT model, demonstrating the robustness of our inferred phylogenetic relationships.

12 – Bayesian method - MCMCTREE in PAML is good, but since Bayesian methods tend to vary, I suggest the group to run another test using the most known softwares (BEAST or mrBayes), to check if these mrca inferences are matching properly
Response: As suggested, we have now estimated divergence times using BEAST version 2. BEAST v2 estimates matched our previous results from MCMCTREE well and were within the estimated confidence intervals. Divergence times estimated by BEAST v2 were added to Figure 4 of the manuscript.

13 – Figure 1 - Doesn't need to be a main figure. Could be used as supplementary figure.
Response: The manuscript has been changed to Data Note and we believe that Figure 1 is important for this manuscript type.

14 – Figure 2 B - These sequences could be mentioned in the text and added as supplementary file. You can name these repeats if needed in figure 2 A.
Response: We agree and have moved the text to the supplementary data.

15 – Figure 3 - Figure is fine but needs to improve image quality. It is preferable to have a Venn diagram of the orthologs between these species.
Response: We have now replaced the figure with a non-proportional Venn diagram.

16 – Add a circus synteny plot figure between the new genome and the closest species genome available.
Response: We have re-submitted the manuscript as Data Note and we believe that in this case a synteny plot is not needed. Additionally, while we agree that a synteny plot would be valuable, generating a synteny plot would be problematic for the Paragonimus genome, as no close relative genome of high quality is available that would allow ordering of the scaffolds.

17 – Figure 4 - (optional) Try to make the same figure using Figtree. They have a nicer way to show the median of the mrca on each node.
Response: We have now improved the figure and aligned the numbers with the tree branches.

Reviewer 1 minor comments:
18 – Change the word faeces for stool. It's not wrong, but stool is more commonly used worldwide;
Response: We have changed "faeces" to "stool" as suggested.

19 – Line 148 - Data Sequencing: add the Illumina Platform used in the data generation (example: HiSeq2000)
Response: Done as suggested.

20 – Line 150 - Data Sequencing: add the PacBio Platform used in the data generation (example: PacBio Sequel or RSII)
Response: Done as suggested.

Reviewer: 2
The manuscript is currently written in my opinion as a data note rather than a research type manuscript. If this is the intention this should be made clearer by the authors as part of their submission. If the manuscript is intended to be submitted as a research paper, the authors should expand on their discussion and conclusions of their data.

Response: We have resubmitted the manuscript as a Data Note.


1 – Abstract, line 85 and Data description, line 157: The authors computationally determined the estimated size of the P.westermani genome, prior to assembly of the raw reads. The computationally determined estimated size was slightly larger than the assembled genome size. The authors should comment on the size difference. In addition, the authors interchange throughout the manuscript whether they compare the estimated size or the assembled genome size with other known published trematode genomes. Until it can be shown that the genome of P.westermani is actually 1.1 Gb, the authors should only refer to the assembled genome size particularly in the section around line 157, as these published trematode genomes describe only the assembled genome sizes.
Response: As suggested, we have added a comment regarding the genome size differences and now base the genome size comparison on the assembled genome sizes.


2 – Line 144 - at what point of infection were the parasites recovered - specifically how old were the parasites?
Response: The parasites were 30-40 days old, this information has been added to the manuscript.

3 – Lines 144-146 -Further information is required regarding the methodology of genomic DNA extraction. Was the extraction carried out on individual worms and then combined or were the worms combined for extraction? Was the genomic DNA quality checked?
Response: The following information has been added to the manuscript: "Genomic DNA was isolated from a pool of 50 worms (30 – 40 days of age), yielding 18 µg of DNA. DNA was quantified by Pico green, QUBIT and NanoDrop. Degradation was tested by Microplate Reader and Agarose Gel Electrophorese (concentration of agarose gel: 1%, electrophoresis time: 40 min, voltage 150 V).


4 – Line 150 - can the authors confirm that the PacBio sequencing was performed on the same sample of genomic DNA?
Response: We confirm that the same sample of genomic DNA was used for PacBio and Illumina sequencing.


5 – Line 255 - the authors should mention that the RNAseq data was from adult parasites only, not the various different lifecycle stages.
Response: This information has now been added to the manuscript.


6 – Line 221 - Related to point 3, as the authors extracted DNA from 50 individual worms, did they check the level of polymorphism at the individual worm level for this region?
Response: We agree with reviewer that this would be an interesting question. However, DNA was isolated from a pool of 50 individual worms. Moreover, only 5 reads spanned the region in full (anchored in non-repetitive sequence at both ends), which was sufficient to generate a consensus sequence for the region, but not to accurately quantify individual-level differences.


7 – Line 272-273 - If the authors are submitting a research themed manuscript, they could include some further discussion of the predicted protein coding genes, particularly those predicted proteins that have inferred homologs in other trematodes (Fig 3A) and the Paragonimus-specific predicted proteins.
Response: The manuscript has been re-submitted as Data Note.


8 – Figure 3A - the venn diagram is currently difficult to interpret, particularly given its current small size as a multi-panel figure. I suggest amending this figure to a classical

venn diagram or an Upset plot.
Response: The figure has been replaced by a non-proportional Venn diagram.

9 – The authors should include supplemental data detailing the functional annotation particularly the analysis of the functional domains, transmembrane domains and signal peptides, as well the data relating to the single copy predicted proteins used for the phylogenetic analysis.
Response: As requested, we have now uploaded our InterProScan results as well as the sequences for single copy proteins used for the phylogenetic analysis to the GigaScience ftp server.

10 – Minor corrections:
a.line 118 - develop into sporocysts
b.line 161 - 1.3 Gb
c.line 204, 291, 293 - BLAST
d.line 281 - predicted proteome
Response: We thank the reviewer for these comments. The manuscript has been modified as suggested.

| | |
|---|---|
| requested as detailed in our [Minimum Standards Reporting Checklist](#)? | |
| **Availability of data and materials**<br><br>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.<br><br>Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)? | Yes |

1   **Title Page**

2   Whole-genome sequence of the oriental lung fluke *Paragonimus*

3   *westermani*

4

5   Harald Oey[1, #], Martha Zakrzewski[2], Kanwar Narain[6], K. Rekha Devi[6], Takeshi Agatsuma[7],

6   Sujeevi Nawaratna[2], Geoffrey Gobart[2,4], Malcolm Jones[3], Mark A. Ragan[5], Donald P.

7   McManus[2,*], Lutz Krause[1,2,*,#]

8

9   [1] The University of Queensland Diamantina Institute, Brisbane, QLD, Australia

10  [2] QIMR Berghofer Medical Research Institute, Brisbane, QLD, Australia

11  [3] School of Veterinary Science, University of Queensland, Gatton, QLD, Australia

12  [4] School of Biological Sciences, Queen's University Belfast, UK

13  [5] Institute for Molecular Bioscience, The University of Queensland, Brisbane, QLD, Australia

14  [6] Regional Medical Research Centre for Northeastern region, Indian Council of Medical

15  Research Department of Health Research Ministry of H & FW, Govt. of India, Dibrugarh-786

16  001, India

17  [7] Department of Environmental Medicine, Kochi Medical School, Kochi University, Kohasu,

18  Oko, Nankoku City Kochi 783-8505

19

20  [*] authors contributed equally

21  [#] To whom correspondence should be addressed

22  Contact: l.krause@uq.edu.au and h.oey@uq.edu.au

23

## Author names, affiliations and emails:

**Dr Harald Oey**

The University of Queensland Diamantina Institute, The Faculty of Medicine, The University of Queensland, Translational Research Institute, Brisbane, QLD

Email: h.oey@uq.edu.au


**Dr Martha Zakrzewski**

QIMR Berghofer Medical Research Institute, Brisbane, QLD, Australia

Email: Martha.Zakrzewski@qimrberghofer.edu.au. ORCID: 0000-0003-1436-5070.


**Dr Kanwar Narain**

Regional Medical Research Centre for Northeastern region, Indian Council of Medical Research Department of Health Research Ministry of H & FW, Govt. of India, Dibrugarh-786 001, India

Email: kanwar_narain@hotmail.com


**Dr K. Rekha Devi**

Regional Medical Research Centre for Northeastern region, Indian Council of Medical Research Department of Health Research Ministry of H & FW, Govt. of India, Dibrugarh-786 001, India

Email: krekha75@yahoo.co.in

49    **Prof Takeshi Agatsuma**

50    Department of Environmental Medicine

51    Kochi Medical School, Kochi University, Kohasu, Oko, Nankoku City Kochi 783-8505

52    Email: agatsuma@kochi-u.ac.jp

53

54    **Dr Geoffrey Gobert**

55    School of Biological Sciences, Queen's University Belfast, UK

56    Email: g.gobert@qub.ac.uk

57

58    **Prof Mark A. Ragan**

59    Institute for Molecular Bioscience, The University of Queensland, Brisbane, QLD, Australia

60    Email: m.ragan@uq.edu.au. ORCID: 0000-0003-1672-7020

61

62    **Prof Malcolm Jones**

63    School of Veterinary Science, University of Queensland, Gatton, QLD, Australia

64    Email: m.jones@uq.edu.au

65

66    **Dr Sujeevi Nawaratna**

67    QIMR Berghofer Medical Research Institute, Brisbane, QLD, Australia

68    Email: sujeevi.nawaratna@qimrberghofer.edu.au. ORCID: 0000-0001-8716-2296

69

70

71

72 **Prof Donald P. McManus**

73 QIMR Berghofer Medical Research Institute, Brisbane, QLD, Australia

74 Email: Don.McManus@qimrberghofer.edu.au

75

76 **A/Prof Lutz Krause**

77 The University of Queensland Diamantina Institute, The Faculty of Medicine, The University

78 of Queensland, Translational Research Institute, Brisbane, QLD

79 Email: L.Krause@uq.edu.au. ORCID: 0000-0003-3806-0845

## Abstract

### Background

Foodborne infections caused by lung flukes of the genus *Paragonimus* are a significant and widespread public health problem in tropical areas. Around 50 *Paragonimus* species have been reported to infect animals and humans, but *Paragonimus westermani* is responsible for the bulk of human disease. Despite their medical and economic importance, no genome sequence for any *Paragonimus* species is available.

### Results

We sequenced and assembled the genome of *P. westermani*, which is among the largest of the known pathogen genomes with an estimated size of 1.1 Gb. A 922.8 Mb genome assembly was generated from Illumina and PacBio sequence data, covering 84% of the estimated genome size. The genome has a high proportion (45%) of repeat-derived DNA, particularly of the LINE and LTR subtypes, and the expansion of these elements may explain some of the large size. We predicted 12,852 protein coding genes, showing a high level of conservation with related trematode species. The majority of proteins (80%) had homologs in the human liver fluke *Opisthorchis viverrini* with an average sequence identity of 64.1%. Assembly of the *P. westermani* mitochondrial genome from long PacBio reads resulted in a single high-quality circularized 20.6 kb contig. The contig harboured a 6.9 kb region of non-coding repetitive DNA comprised of three distinct repeat units. Our results suggest that the region is highly polymorphic in *P. westermani*, possibly even within single worm isolates.

### Conclusions

The generated assembly represents the first *Paragonimus* genome sequence and will facilitate future molecular studies of this important, but neglected, parasite group.

### Keywords

*Paragonimus westermani*, whole-genome sequence, genome assembly, paragonimiasis, food-borne disease, oriental lung fluke, parasitic infection, bioinformatics, high-throughput sequencing, comparative genomics, genome annotation, neglected tropical disease, flatworm

## Background

*Paragonimus* lung flukes represent a significant and widespread clinical problem with an estimated 23 million people infected worldwide (1). Around 50 species are described, with at least 7 being human pathogens (2). The majority of human *Paragonimus* infections can be attributed to the *P. westermani* species complex, mainly in Southeast Asia and Japan (1). *P. westermani* show considerable geographic genetic variability and human infections occur predominantly in East Asia and the Philippines. In India the incidence rates of paragonimiasis caused by *P. westermani* is currently unknown(2-4), however many cases of paragonimiasis are attributed to the related worm *Paragonimus heterotremus* (2). Paragonimiasis is a zoonotic disease and also pigs, dogs and other animals can harbour *P. westermani* (2).

*Paragonimus* spp. have a complex life cycle. Unembryonated eggs are expelled by coughing or passed with stool and develop in water. Miracidia hatch from the eggs and penetrate a freshwater snail, its first intermediate host. During several asexual developmental phases inside the snail, a miracidium develops into a sporocyst and then two redial generations occur, the second of which gives rise to microcercous cercariae that escape into fresh water. These crawling cercariae invade a species of crustacean, the second intermediate host, to encyst in muscles and other sites and develop into metacercariae. Humans and other definitive hosts become infected through consumption of raw or inadequately cooked freshwater crabs or crayfish (5). Ingested metacercariae excyst, penetrate through the gut and become encapsulated in the lungs where they mature into hermaphroditic adult worms (7.5 mm to 12 mm in length) in 6-10 weeks (5). Paragonimiasis can lead to a chronic inflammatory disease of the lung and can trigger asthma- or tuberculosis-like symptoms (6-8). In more severe cases *Paragonimus* can infect the brain or central nervous system of the definitive host, leading to headache, visual loss, and even death (1).

Paragonimiasis is commonly diagnosed by microscopic detection of parasite eggs in stool or sputum. The lack of sensitive and reliable diagnostic tests in conjunction with unspecific disease symptoms often leads to delayed treatment with the drug of choice, praziquantel (8). Despite their high medical, veterinary and economic importance, only limited information on the molecular biology of *Paragonimus* is currently available. Recent

transcriptome sequencing studies have provided some information on the gene content of *Paragonimus* (9), but until now no *Paragonimus* genome sequence has been available. Here we present a 922.8 Mb assembly of the *P. westermani* genome which provides new insights into the genomic composition of the *Paragonimus* genus and represents an invaluable resource for future studies of the neglected tropical disease paragonimiasis.

## Data Description

**Sequencing**

Diploid *Paragonimus westermani* metacercariae (NCBI:txid34504) were collected from the freshwater crab *Maydelliathelphusa lugubris* in 2009 in the East Siang district of Arunachal Pradesh, Northeast India, and fed to Wistar rats as experimental hosts. Genomic DNA was isolated from a pool of 50 worms (30 − 40 days of age), yielding 18 μg of DNA. DNA was quantified by Pico green, Qubit and NanoDrop and degradation was tested by Microplate Reader and Agarose Gel Electrophorese (concentration of agarose gel: 1%, electrophoresis time: 40 min, voltage 150 V). The *Paragonimus westermani* genome was then sequenced from 2 μg of the isolated DNA using a whole-genome shotgun approach. Paired-end short-insert (200 bp and 450 bp) and mate-pair (5 Kb and 10 Kb) genomic DNA libraries were sequenced on the Illumina HiSeq 2000 platform, yielding 58 Gb of sequence data (**Table 1**). For genome scaffolding and quality evaluation of the assembled sequence, additional long-read data were generated from the same genomic DNA sample using the PacBio RSII platform, yielding 1.7 Gb of information (**Table 1**). The genome size was estimated from the K-mer coverage of the 450 bp insert library. K-mer frequencies were calculated by the program Jellyfish (10), version 2.2.6 , using a K-mer size of 17bp. The 17-mer distribution in the 450 bp library had a single peak at 26x (**Figure 1**), demonstrating low sequence heterozygosity. The genome size (G) was deduced from the K-mer distribution via the formula $G = N * (L − K + 1) / K\_depth(11)$, where N is the total number of reads, L is the read length, K is the K-mer size and K_depth is the peak frequency. The *P. westermani* genome size was estimated to be 1.1 Gb.

**Table 1. *Paragonimus westermani* sequencing libraries.**

| Library | Platform | Library type | Insert size (bp) | Read length (bp) | Read count (raw) |
|---|---|---|---|---|---|
| 200bp | HiSeq | Paired-end | 200 | 2 x 120 | 140,542,299 |
| 450bp | HiSeq | Paired-end | 450 | 2 x 100 | 171,954,230 |
| 5kb | HiSeq | Mate-pair | 5,000 | 2 x 49 | 232,630,904 |
| 10kb | HiSeq | Mate-pair | 10,000 | 2 x 49 | 266,480,540 |
| PacBio | PacBio | Long read | - | - | 1,731,327 |

169

170

## Genome assembly

172 PacBio sequence data were error corrected by proovread version 2.13.13 (12), using

173 Illumina short reads from the 200bp and 450bp libraries as input , and assembled into

174 contigs by Mira v4.0.2 (MIRA, RRID:SCR_010731)(13). Short-read Illumina sequence data

175 were trimmed using Trimmomatic v0.36 (Trimmomatic , RRID:SCR_011848) and

176 subsequently error corrected by KmerFreq_HA (part of SoapDenovo2 (14)) with a K-mer size

177 of 23. The 10 kb mate-pair library showed a high proportion of PCR duplicates and was

178 subjected to PCR de-duplication prior to genome assembly. For assembly of short read data,

179 several assembly programs were evaluated. ABYSS performed best for this particular

180 genome with its large size, high percentage of repetitive regions and some low-level

181 sequence heterogeneity resulting from pooling genomic DNA from the 50 individual worms.

182 ABYSS is also one of the few assemblers that allow inclusion of long-read data to guide

183 scaffolding. Illumina paired-end sequence data were assembled using the ABYSS assembly

184 pipeline (ABySS , RRID:SCR_010709)(15), version 2.0.2, with options n=5 s=200 N=36 S=500

185 k=33 and including the PacBio contigs via the re-scaffolding feature.

186 The resulting assembly was de-gapped using the SoapDenovo2 GapCloser program

187 (GapCloser, RRID:SCR_015026)(14) which is well suited for closing gaps larger than 1kb and

188 it performed particularly well on this genome. Mate-pair libraries were then used to scaffold

189 the assembly with SSPACE v3.0(16) (with options -x 0 -a 0.60 -n 30 -z 200 -g 0) and gaps

190 were again filled with GapCloser. Un-closed gaps are represented by N's spanning the

191 estimated sizes of the gaps. To detect and resolve scaffolding errors, the resulting assembly

192 was processed by the program REAPR (17) using the 5kb mate-pair library as input, breaking

the assembly at sites with poor evidence for contiguity. Contamination due to the experimental rat host and the bacterium *Delftia* sp was detected based on a comparison of predicted proteins with the NCBI protein database using BLAST and, additionally, via the NCBI Genome Submission Portal quality control pipeline. A targeted comparison of all scaffolds with the genomes of the rat and *Delftia* using BLAT identified 531 short scaffolds with high similarity (>90%) to these genomes. These sequences were manually scrutinized, with 529 of the affected scaffolds found to be completely derived from rat or *Delftia*, and these were removed from the assembly. The remaining two contaminated sequences represented rat ribosomal DNA that had been erroneously incorporated into *Paragonimus* scaffolds and were also removed from the final assembly by cutting and trimming the affected scaffolds.

The final assembly resulted in a 922.8 Mb genome sequence (30,466 scaffolds with N50 of 135 Kb) (**Table 2**), covering 84.0% of the estimated genome size. The discrepancies in genome size can potentially be the result of problematic DNA regions that are difficult to sequence or assemble (e.g. regions with strong secondary structures, highly repetitive regions or long homopolymeric runs) or the result of low-level sequence heterogeneity which can lead to an overestimation of genome size by k-mer approaches. The *P. westermani* genome sequence is among the largest known pathogen genomes and one of the largest parasite genomes sequenced to date. The assembled genome sequence is considerably larger than the published genomes of the related trematodes *Clonorchis sinensis* (assembly size of 546.9 Mb) (18), *Opisthorchis viverrini* (606.0 Mb) (19)*,* and *Schistosoma* spp (364.5-397.7 Mb) (20-22) and comparable to the 1.3 Gb genome of *Fasciola hepatica* (23).

The GC content of the genome was 43.3%, comparable to genomes of other related trematodes (**Table 2**). Genome assembly completeness was evaluated by BUSCO (BUSCO , RRID:SCR_015008)(24) using the metazoan lineage data, resulting in similar scores to those obtained for the genomes of comparable trematode species (**Table 2**). The proportion of duplicated genes reported by BUSCO was also similar to that of comparable trematodes suggesting that the relatively large size of the *P. westermani* genome is not the result of genome duplication events.

224

225

**Table 2. Assembly statistics for *P. westermani* and comparable trematode genomes of similar size**.

| | *P. westermani* | *F. hepatica* | *O. viverrini* | *C. sinensis* |
|---|---|---|---|---|
| Assembly size (Mb)[a] | 922.8 | 1,275.0 | 606.0 | 546.9 |
| Ungapped size (Mb)[b] | 877.7 | 1,183.5 | 558.0 | 547.1 |
| Contig N50 (Kb) | 7.0 (>100bp) | 9.7 | NA | 14.7 |
| Scaffold N50 (Kb) | 135 (>1Kb) | 204 | 1,324 | 30.2 |
| Scaffold L50 | 1,943 | 1,799 | 135 | 408 |
| Scaffold count | 30,466 (>1Kb) | 45,354 (>1kb) | 4,919 (>1kb) | 31,822 |
| GC content (%) | 43.3 | 44.1 | 43.8 | 44.1 |
| Repeat content (%) | 45.2 | 57.1 | 28.9 | 32.6 |
| Protein coding genes | 12,852 | 15,740[c] | 16,356 | 13,634 |
| Longest scaffold (Kb) | 809 | 1,565 | 9,657 | 2,050 |
| BUSCOS - Complete | 65.3% | 65.8% | 71.4% | 70.8% |
| BUSCOS - Duplicated | 1.4% | 0.8% | 1.1% | 1.5% |
| BUSCOS - Missing | 25.8% | 25.4% | 23.0% | 23.1% |

228 [a]Combined length of all scaffolds in Mb. [b]Combined length of all scaffolds without gaps (N's)

229 in Mb. [c]Non-overlapping RNAseq-supported gene models(23)

230

**Mitochondrial genome**

232 The mitochondrial genome of *P. westermani* is present at a much higher copy number than

233 the nuclear genome and we were able to assemble the full mitochondrial genome at high

234 coverage from error corrected long PacBio reads using the Mira assembler(13), version

235 4.0.2. This resulted in a single mitochondrial contig of 20.3 Kb (**Figure 2**). The accuracy of the

236 contig was confirmed by mapping short insert paired-end sequences directly onto the contig

237 revealing single nucleotide discrepancies at only 4 positions. The mitochondrial genome was

238 found to closely match previously published *Paragonimus* mitochondrial genomes with the

239 best match from a BLAST search against the Nucleotide collection at the NCBI

240 (https://blast.ncbi.nlm.nih.gov/Blast.cgi) being accession NC_027673.1, a *P. westermani*

241 complex sp. type 1 mitochondrial genome isolated in India (97% sequence identity across

242 13.4 Kb of NC_027673.1). This sequence was used as reference for mitochondrial gene

243 identification and annotation, supplemented by mitochondrial gene predictions by Mitos

244 (25) and tRNA prediction by Aragorn (Aragorn, RRID:SCR_015974)(26). The mitochondrial

245 genomes of flatworms are known to harbour a region of non-coding repetitive DNA,

246 generally comprised of a long noncoding region (LNR) and a short non-coding region (SNR)

247 with a single tRNA gene separating them(27). Reconstructing this region from short-read

248 data proved challenging, but our long-read PacBio data allowed complete assembly of the

249 repetitive region and circularization of the genome. Interestingly, our assembled

250 mitochondrial genome sequence had a much longer non-coding region (6.9 Kb) than the

251 previously published NC_027673.1 (0.7 Kb) and the non-coding regions of both genomes

252 showed only partial homology, but with close homology of the intervening tRNA gene. We

253 found the LNR to be comprised of two distinct repeat units with 8 and 13 copies while the

254 SNR was comprised of another distinct repeat unit with 3 copies (**Figure 2 and Additional**

255 **File 1**). Strikingly, five independent PacBio reads spanned the entirety of the non-coding

256 region but with slight differences in length (6.3 – 6.9 Kb), suggesting that the region is

257 polymorphic, possibly even within individual worms.

**Repeat annotation**

259 RepBase repeat consensus sequences did not adequately represent the repeats found in the

260 *P. westermani* assembly, consistent with the distant evolutionary relationship of lung flukes

261 with previously sequenced worm genomes. We therefore carried out *de-novo* repeat

262 characterization using the RepeatModeller package, version 1.0.9 (RepeatModeler,

263 RRID:SCR_015027), and used the generated consensus sequences to identify repetitive

264 regions by RepeatMasker (RepeatMasker , RRID:SCR_012954), version 4.0.7 (both available

265 at http://www.repeatmasker.org). To enable direct comparison with related trematode

266 species we also ran RepeatModeller and RepeatMasker separately on the *F. hepatica*, *O.*

267 *viverrini* and *C. sinensis* genomes with the same program parameters as those used for *P.*

268 *westermani*.

269 A relatively high percentage (45.2%) of the *P. westermani* genome sequence was repeat-

270 derived, similar to the rate reported for *Schistosoma* spp. (40.1-47.5%) (20-22) and *F.*

271 *hepatica* (57.1%), but considerably higher than the rate observed for the closer relatives *O.*

272 *viverrini* (28.9%) and *C. sinensis* (32.6%) (**Table 3**). Retrotransposons of the long interspersed

273 nuclear element (LINE) subtype were found to be the greatest contributors of repetitive

274 DNA (21.6%) (**Table 3**), consistent with reports for other trematode genomes (23). In *P.*

275 *westermani* and *F. hepatica*, the two largest of the four included trematode genomes, long

276 terminal repeat (LTR) retrotransposons were also highly abundant contributing 7.7% and

277 10.1% of the genomes, respectively. Additionally, all four genomes had considerable

278 amounts of repetitive DNA (10.7-17.1%) that did not match repeat consensus sequences of

279 any of the known repeat classes modelled by RepeatModeler. The relatively high proportion

280 of repeat-derived sequences in *P. westermani* may explain some of the increased size

281 observed for this genome compared to the genomes of related flatworm species.

282

283 **Table 3. Repeat content percentage of *P. westermani* and related trematode genome**

284 **sequences.**

| Repeat class | *P. westermani* | *F. hepatica* | *O. viverrini* | *C. sinensis* |
|---|---|---|---|---|
| LINE | 21.57 | 26.17 | 12.76 | 14.85 |
| LTR | 7.71 | 10.06 | 2.82 | 1.97 |
| DNA elements | 1.76 | 2.14 | 0.94 | 1.04 |
| SINE | 0.96 | 1.06 | 1.26 | 1.22 |
| Simple repeats | 0.18 | 0.63 | 0.43 | 0.36 |
| Unclassified | 12.97 | 17.06 | 10.69 | 13.15 |
| **Total** | **45.15** | **57.12** | **28.9** | **32.59** |

285

286 **Gene prediction and functional annotation**

287 Genes were predicted by the Maker pipeline, version 2.31.9, using Augustus (28), version

288 3.2.3, and GeneMark-ES (29), version 4.32, for *ab-initio* gene prediction. To accurately

289 model the sequence properties of the *P. westermani* genome, both gene finders were

290 initially trained by BRAKER1 (30), version 1.9, which makes use of mapped transcriptome

291 sequence data. Previously published RNA-seq data from adult *P. westermani* (9) were

292 obtained from the short read archive and mapped to our genome assembly using the Star

293 aligner (31), version 2.5, with the option --twopassMode Basic. BRAKER1 was then run with

294 default parameters. The RNA-seq data was further assembled into transcripts using cufflinks

295 (32), version 2.2.1, with the options --frag-bias-correct <p.westermani assembly> --multi-

296 read-correct. The resulting transcripts were provided as input for Maker via the "est_gff"

297 option. For homology based searches Maker was provided with the following wormbase v8

298 protein datasets: *Clonorchis sinensis* (PRJDA72781), *Opisthorchis viverrini* (PRJNA222628),

299 *Schistosoma mansoni* (PRJEA36577), *Caenorhabditis elegans* (PRJNA13758), *Echinococcus*

300 *granulosus* (PRJEB121), *Hymenolepis diminuta* (PRJEB507) and *Schistosoma haematobium*

301 (PRJNA78265). Additionally, the Swiss-Prot dataset from UniProt was included. Maker was

302 allowed to report single exon genes, and otherwise run with default parameters.

303 Proteins were functionally annotated based on a BLASTp search against the NCBI non-

304 redundant protein database (obtained on 25.10.17) requiring an e-value <1e-15 and the

305 best hit spanning at least 40% of the query sequence. KEGG annotations were identified

306 using the BlastKoala server with the option "genus_eukaryotes" (33). Additionally,

307 functional domains, GO annotations, transmembrane proteins and signal peptides were

308 identified with InterProScan (InterProScan , RRID:SCR_005829)(34), version 5.25-64.0. GO

309 annotations were then visualized using WEGO (35). In total, 12,852 protein encoding genes

310 were predicted in the *P. westermani* genome and functionally annotated (**Table 2**).

311

**Genome comparison**

313 Predicted *P. westermani* coding genes were mapped to the genomes of related trematode

314 species using Exonerate, version 2.4.0, requiring a minimal sequence identity of 30% and

315 excluding matches spanning less than 40% of the query protein. The majority of predicted

316 proteins (86.2%) had inferred homologs in the related trematode species (**Figure 3A**) and

317 showed a similar distribution of protein functional categories (**Figure 3C**). The *P. westermani*

318 predicted proteome was most similar to *O. viverrini* and *C. sinensis*. Of the 12,852 predicted

319 proteins, 10,350 (80%) had inferred homologs in *O. viverrini* with an average sequence

320 identity of 64.1%, and 10,227 (79.6%) had homologs in *C. sinensis* with an average sequence

321 identity of 63.8% (**Figures 3A and 3B**).

322

**Phylogenetic analysis and estimation of divergence time**

A protein-based phylogenetic tree was inferred from 14 worm genomes, including *P. westermani*, 12 related trematode/cestode species and *Schmidtea mediterranea*, a free-living turbellarian flatworm, as outgroup (**Figure 4**). We first identified single-copy proteins shared across all 14 included worm species. Single-copy proteins were identified based on BLASTp searches of a species proteins against the species own proteome using a sequence-identity cut-off of 30% and requiring hits to cover >50% of the query sequence. Single-copy proteins shared across all 14 species were then identified using a less stringent BLASTp search with a 30% sequence identity cut-off but requiring only >40% coverage of the query sequence. We identified 104 single-copy proteins shared across the 14 worm species that were then aligned using MUSCLE (36). The resulting multiple sequence alignment was de-gapped with trimAI (37) and a phylogenetic tree was reconstructed by PhyML (PhyML , RRID:SCR_014629)(38). Model selection in PhyML (39) identified the LG model (40) with decorations +G+I+F as optimal. PHYLIP v3.696 (41) using the maximum likelihood method and the Jones-Taylor-Thornton (JTT) probability model (42) resulted in the same tree topology, demonstrating the robustness of the inferred phylogenetic relationships.

The multiple alignment and the inferred phylogenetic tree were then used to estimate species divergence by a Bayesian model with relaxed molecular clock using MCMCTREE in PAML 4.9e (**Figure 4**)(PAML , RRID:SCR_014932). The model was calibrated based on previously published divergence times and ages of fossil records. Evidence for trematode infestation have been reported from the Eocene (56 to 33.9 million years (myr) ago) and preserved trematode eggs have been found in dinosaur coprolites from the Early Cretaceous (146 to 100 myr ago); however, fossil records indicate that trematodes may have already existed more than 400 myr ago (43, 44). The trematode split from other neodermatan lineages was therefore fixed at >56 myr. The origin of schistosomes has been estimated somewhere in the Miocene around 15-20 myr ago (45, 46). It has further been estimated that the divergence of *S. mansoni* did likely not occur before 2-5 myr ago, based on fossil records of its intermediate host *Biomphalaria* (47). From these data, the split of Plagiorchiida (including *P. westermani*) and Opistorchiida (including *O. viverrini* and *C. sinensis*) was estimated to have occurred 38.9 myr ago (95% confidence interval of 28.0-58.6 myr) (**Figure 4**). To estimate the robustness of the inferred divergence times the analysis was repeated using BEAST 2 version 2.5.0 (48), based on the JTT substitution matrix,

355 gamma category count of 4, estimated substitution rate, relaxed clock log normal model,

356 and a chain length of 6M (49, 50). A maximum clade credibility tree using median node

357 heights was generated by the BEAST 2 treeannotator tool. Divergence times inferred by

358 BEAST 2 well matched the MCMCTREE results and were within the estimated confidence

359 intervals (Figure 4). The split of the Plagiorchiida and the Opistorchiida was estimated to

360 have occurred 31.5 myr ago.

361

## Discussion

363 We have presented the first whole-genome sequence of a *Paragonimus* spp. worm,

364 providing a valuable resource to the field that will aid our understanding of this group of

365 clinically important parasites. The genome was found to be unusually large for a worm, a

366 feature that at least in part appears attributable to an expansion of retrotransposable

367 elements, rather than genome duplication events.

368 The mitochondrial genome was also found to be very large comprising 20.3 Kb. Such a large

369 size appears to be a common feature of worms and results from a long repetitive region of

370 unknown function. However, while this region appears to be a feature of most flatworms it

371 is rarely sequenced in full due to the technical challenges of sequencing long tandemly

372 repeated sequences.

373 *P. westermani* has been described as a species complex with considerable genetic

374 differences across geographic regions (2). The genome presented herein is of an Indian

375 isolate and it will be of considerable interest to compare this and the genomes of isolates

376 from other regions where *P. westermani* is endemic to elucidate the region specific genetic

377 features. This would be particularly informative as not all endemic regions are associated

378 with paragonimiasis in humans (2).

379 Phylogenetic analyses of *P. westermani* shows that it has diverged considerably from its

380 closest relatives, *Clonorchis sinensis* and *Opisthorchis viverrini* with a split estimated to have

381 occurred 28-59 myr ago. Subsequent to that split the species spread out across a vast

382 geographical range, acquiring distinct local traits in what may eventually be considered

383 speciation events. This time-span has also seen an expansion of two repeat families, in

384 particular the LINE and LTR elements. In mammals these elements are known to
385 occasionally become exapted and gain novel regulatory functions (51), and they are
386 therefore likely to add to the diversity of the *P. westermani* species complex.

387

388

389

## Conclusion

391 The presented *P. westermani* genome assembly provides new insights into the molecular
392 biology of *Paragonimus* and provides an unprecedented resource for functional studies of
393 lung flukes and for the design of new disease interventions and diagnostics tests.

394

## Availability of supporting data

396 The nuclear and mitochondrial genomes are available from NCBI under accession number
397 PRJNA454344. Annotation and tree data is available from the *GigaScience* GigaDB repository
398 (52).

399

## Abbreviations

401 bp        base pair

402 BUSCO Benchmarking Universal Single-Copy Orthologs

403 Kb        Kilo base pair

404 LNR        Long noncoding region

405 LINE        Long interspersed nuclear elements

406 LTR        Long terminal repeat

407 Mb        Mega base pair

408    MYR    Million years

409    SINE    Short interspersed nuclear elements

410

## Competing interests

412    All authors declare that they have no competing interests.

413

## Funding

418

## Author contributions

420    LK and DPM conceived and managed the project; KN and KRD provided *P. westermani*
421    material. TA and SN isolated genomic DNA. MZ and GG managed DNA sequencing. HO
422    carried out genome assembly, gene prediction and functional genome annotation. HO and
423    LK carried out comparative genomics. LK, DPM, MKJ and MAR attracted funding and
424    designed the study. LK and HO drafted the manuscript and all authors read, edited and
425    approved the final manuscript.

426

## Acknowledgements

430

## References

432  1.      Furst T, Keiser J, Utzinger J. Global burden of human food-borne trematodiasis: a systematic
433  review and meta-analysis. Lancet Infect Dis. 2012;12(3):210-21.
434  2.      Blair D. Paragonimiasis. Adv Exp Med Biol. 2014;766:115-52.
435  3.      Roy JS, Das PP, Borah AK, Das JK. Paragonimiasis in a Child from Assam, India. J Clin Diagn
436  Res. 2016;10(4):DD06-7.
437  4.      Singh TS, Hiromu S, Devi KR, Singh WA. First case of Paragonimus westermani infection in a
438  female patient in India. Indian J Med Microbiol. 2015;33 Suppl:156-9.
439  5.      Jones MK, Keiser J, McManus DP. Trematodes. In: Jorgensen JH, Pfaller MA, Carroll KC,
440  Funke G, Landry ML, Richter SS, et al., editors. Manual of Clinical Microbiology, Eleventh Edition:
441  American Society of Microbiology; 2015.
442  6.      Luo J, Wang MY, Liu D, Zhu H, Yang S, Liang BM, et al. Pulmonary Paragonimiasis Mimicking
443  Tuberculous Pleuritis: A Case Report. Medicine (Baltimore). 2016;95(15):e3436.
444  7.      Zhou R, Zhang M, Cheng N, Zhou Y. Paragonimiasis mimicking chest cancer and abdominal
445  wall metastaisis: A case report. Oncol Lett. 2016;11(6):3769-71.
446  8.      Kalhan S, Sharma P, Sharma S, Kakria N, Dudani S, Gupta A. Paragonimus westermani
447  infection in lung: A confounding diagnostic entity. Lung India. 2015;32(3):265-7.
448  9.      Li BW, McNulty SN, Rosa BA, Tyagi R, Zeng QR, Gu KZ, et al. Conservation and diversification
449  of the transcriptomes of adult Paragonimus westermani and P. skrjabini. Parasit Vectors.
450  2016;9:497.
451  10.     Marcais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of
452  occurrences of k-mers. Bioinformatics. 2011;27(6):764-70.
453  11.     Song L, Bian C, Luo Y, Wang L, You X, Li J, et al. Draft genome of the Chinese mitten crab,
454  Eriocheir sinensis. Gigascience. 2016;5:5.
455  12.     Hackl T, Hedrich R, Schultz J, Forster F. proovread: large-scale high-accuracy PacBio
456  correction through iterative short read consensus. Bioinformatics. 2014;30(21):3004-11.
457  13.     Chevreux B, Pfisterer T, Drescher B, Driesel AJ, Muller WE, Wetter T, et al. Using the miraEST
458  assembler for reliable and automated mRNA transcript assembly and SNP detection in sequenced
459  ESTs. Genome research. 2004;14(6):1147-59.
460  14.     Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved
461  memory-efficient short-read de novo assembler. Gigascience. 2012;1(1):18.
462  15.     Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I. ABySS: a parallel assembler for
463  short read sequence data. Genome research. 2009;19(6):1117-23.
464  16.     Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-assembled contigs
465  using SSPACE. Bioinformatics. 2011;27(4):578-9.
466  17.     Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD. REAPR: a universal tool for
467  genome assembly evaluation. Genome Biol. 2013;14(5):R47.
468  18.     Wang X, Chen W, Huang Y, Sun J, Men J, Liu H, et al. The draft genome of the carcinogenic
469  human liver fluke Clonorchis sinensis. Genome Biol. 2011;12(10):R107.
470  19.     Young ND, Nagarajan N, Lin SJ, Korhonen PK, Jex AR, Hall RS, et al. The Opisthorchis viverrini
471  genome provides insights into life in the bile duct. Nat Commun. 2014;5:4378.
472  20.     Young ND, Jex AR, Li B, Liu S, Yang L, Xiong Z, et al. Whole-genome sequence of Schistosoma
473  haematobium. Nat Genet. 2012;44(2):221-5.
474  21.     Berriman M, Haas BJ, LoVerde PT, Wilson RA, Dillon GP, Cerqueira GC, et al. The genome of
475  the blood fluke Schistosoma mansoni. Nature. 2009;460(7253):352-8.
476  22.     Schistosoma japonicum Genome S, Functional Analysis C. The Schistosoma japonicum
477  genome reveals features of host-parasite interplay. Nature. 2009;460(7253):345-51.
478  23.     Cwiklinski K, Dalton JP, Dufresne PJ, La Course J, Williams DJ, Hodgkinson J, et al. The
479  Fasciola hepatica genome: gene duplication and polymorphism reveals adaptation to the host
480  environment and the capacity for rapid evolution. Genome Biol. 2015;16:71.

481    24.    Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing
482    genome assembly and annotation completeness with single-copy orthologs. Bioinformatics.
483    2015;31(19):3210-2.
484    25.    Bernt M, Donath A, Juhling F, Externbrink F, Florentz C, Fritzsch G, et al. MITOS: improved de
485    novo metazoan mitochondrial genome annotation. Mol Phylogenet Evol. 2013;69(2):313-9.
486    26.    Laslett D, Canback B. ARAGORN, a program to detect tRNA genes and tmRNA genes in
487    nucleotide sequences. Nucleic Acids Res. 2004;32(1):11-6.
488    27.    Le TH, Blair D, McManus DP. Mitochondrial genomes of parasitic flatworms. Trends
489    Parasitol. 2002;18(5):206-13.
490    28.    Stanke M, Steinkamp R, Waack S, Morgenstern B. AUGUSTUS: a web server for gene finding
491    in eukaryotes. Nucleic Acids Res. 2004;32(Web Server issue):W309-12.
492    29.    Lomsadze A, Ter-Hovhannisyan V, Chernoff YO, Borodovsky M. Gene identification in novel
493    eukaryotic genomes by self-training algorithm. Nucleic Acids Res. 2005;33(20):6494-506.
494    30.    Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: Unsupervised RNA-Seq-
495    Based Genome Annotation with GeneMark-ET and AUGUSTUS. Bioinformatics. 2016;32(5):767-9.
496    31.    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal
497    RNA-seq aligner. Bioinformatics. 2013;29(1):15-21.
498    32.    Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript
499    assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching
500    during cell differentiation. Nat Biotechnol. 2010;28(5):511-5.
501    33.    Kanehisa M, Sato Y, Morishima K. BlastKOALA and GhostKOALA: KEGG Tools for Functional
502    Characterization of Genome and Metagenome Sequences. J Mol Biol. 2016;428(4):726-31.
503    34.    Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale
504    protein function classification. Bioinformatics. 2014;30(9):1236-40.
505    35.    Ye J, Fang L, Zheng H, Zhang Y, Chen J, Zhang Z, et al. WEGO: a web tool for plotting GO
506    annotations. Nucleic Acids Res. 2006;34(Web Server issue):W293-7.
507    36.    Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
508    Nucleic Acids Res. 2004;32(5):1792-7.
509    37.    Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. trimAl: a tool for automated alignment
510    trimming in large-scale phylogenetic analyses. Bioinformatics. 2009;25(15):1972-3.
511    38.    Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New Algorithms and
512    Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0.
513    Systematic Biology. 2010;59(3):307-21.
514    39.    Lefort V, Longueville JE, Gascuel O. SMS: Smart Model Selection in PhyML. Mol Biol Evol.
515    2017;34(9):2422-4.
516    40.    Le SQ, Gascuel O. An improved general amino acid replacement matrix. Mol Biol Evol.
517    2008;25(7):1307-20.
518    41.    Felsenstein J. Confidence-Limits on Phylogenies - an Approach Using the Bootstrap.
519    Evolution. 1985;39(4):783-91.
520    42.    Jones DT, Taylor WR, Thornton JM. The Rapid Generation of Mutation Data Matrices from
521    Protein Sequences. Computer Applications in the Biosciences. 1992;8(3):275-82.
522    43.    Poinar G, Jr., Boucot AJ. Evidence of intestinal parasites of dinosaurs. Parasitology.
523    2006;133(Pt 2):245-9.
524    44.    Huntley JW, De Baets K. Trace Fossil Evidence of Trematode-Bivalve Parasite-Host
525    Interactions in Deep Time. Adv Parasitol. 2015;90:201-31.
526    45.    Littlewood DTJe, Baets Kde. Fossil parasites.
527    46.    Snyder SD, Loker ES. Evolutionary relationships among the Schistosomatidae
528    (Platyhelminthes:Digenea) and an Asian origin for Schistosoma. J Parasitol. 2000;86(2):283-8.
529    47.    Morgan JA, Dejong RJ, Snyder SD, Mkoji GM, Loker ES. Schistosoma mansoni and
530    Biomphalaria: past history and future trends. Parasitology. 2001;123 Suppl:S211-28.

531 48.     Bouckaert R, Heled J, Kuhnert D, Vaughan T, Wu CH, Xie D, et al. BEAST 2: A Software
532 Platform for Bayesian Evolutionary Analysis. Plos Computational Biology. 2014;10(4).
533 49.     Heled J, Drummond AJ. Calibrated tree priors for relaxed phylogenetics and divergence time
534 estimation. Syst Biol. 2012;61(1):138-49.
535 50.     Drummond AJ, Ho SY, Phillips MJ, Rambaut A. Relaxed phylogenetics and dating with
536 confidence. PLoS Biol. 2006;4(5):e88.
537 51.     Mita P, Boeke JD. How retrotransposons shape genome regulation. Curr Opin Genet Dev.
538 2016;37:90-100.
539 52.     Oey H; Zakrzewski M; Narain K; Devi KR; Agatsuma T; Nawaratna S; Gobart G; Jones M;
540 Ragan MA; McManus DP; Krause L (2018): Supporting data for "Complete genome sequence of the
541 oriental lung fluke Paragonimus westermani" GigaScience Database.
542 http://dx.doi.org/10.5524/100524

543

544

545

546

547

548

549

550

551 **Figure legends:**

552 **Figure 1. K-mer frequencies for the 450bp library.** Distribution of 17-mers in the 450bp

553 short-insert library demonstrated low sequence heterozygosity. We observed a single peak

554 at 26x and the *P. westermani* genome size was estimated to be 1.1 Gb.

555

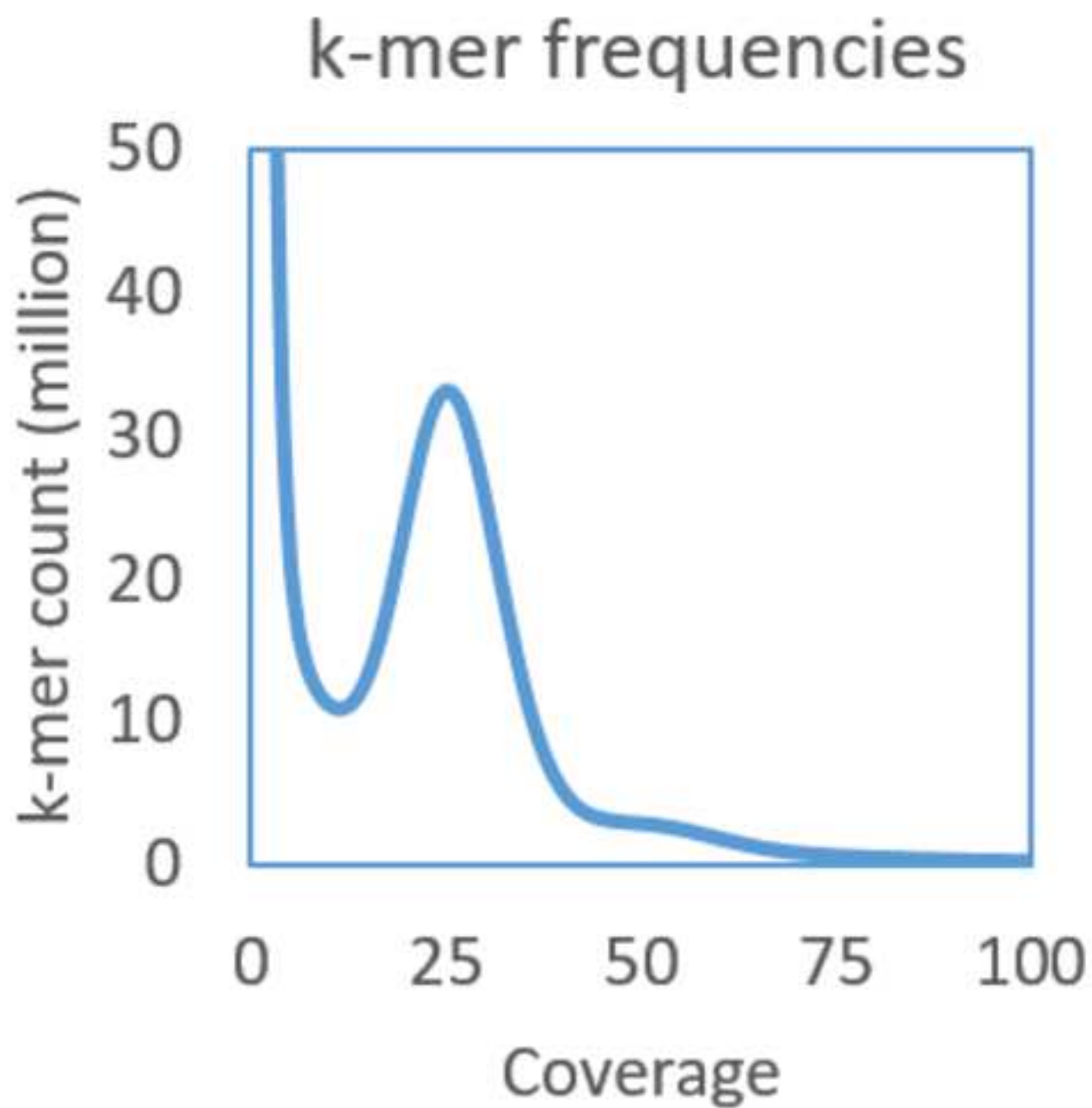556 **Figure 2. The complete *P. westermani* mitochondrial genome.**

557 A graphical representation of the *P. westermani* circular mitochondrial genome is shown,

558 including a ~6.9 Kb repetitive region. Three distinct repeat units were identified in this

559 region, as well as an intervening tRNA gene (tRNA-Glu). All genes are transcribed in the
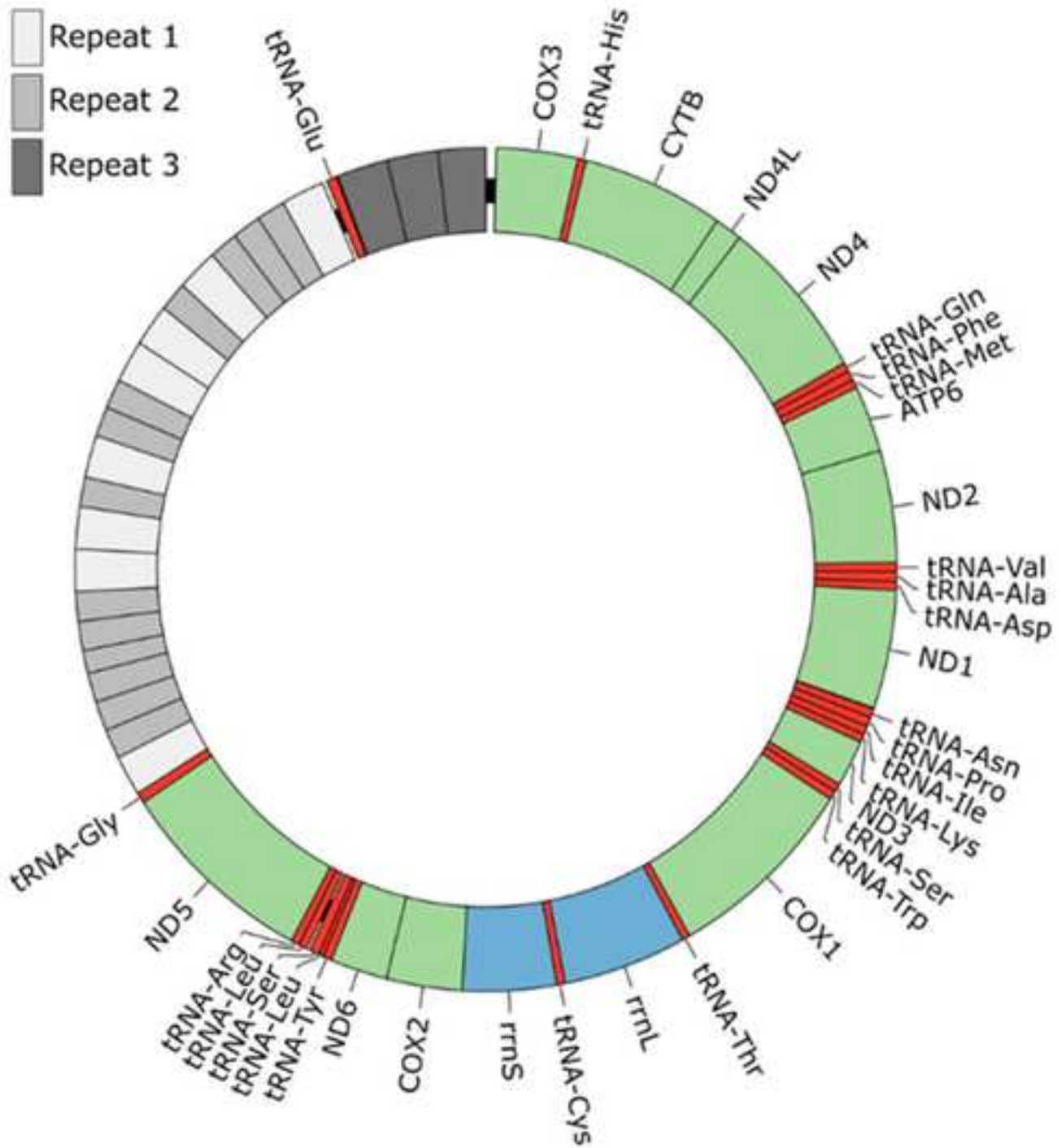
560 clock-wise direction.

561

562 **Figure 3. Conservation of the *P. westermani* proteome across four related trematode**

563 **species.** *P. westermani* proteins were mapped to the genome sequences of *O. viverrini, C.*

564 *sinensis, F. hepatica* and *S. mansoni* using Exonerate. A) *P. westermani* centred Venn

565 diagram of 12,852 predicted proteins. The four included trematode species shared a core

566 set of 7,599 proteins. B) Sequence identity of *P. westermani* proteins and orthologues

567 inferred in genomes of related trematodes. Average sequence identity is given in brackets.

568 C) Distribution of identified functional GO categories across three trematode species. GO

569 annotations were assigned by InterProScan and visualized using WEGO.

570

571 **Figure 4. Phylogenetic tree and estimated divergence times.** A phylogenetic tree of

572 selected trematodes and cestodes and *S. mediterranea* as outgroup was reconstructed from

573 104 shared single-copy proteins using the maximum likelihood method. Species divergence

574 was estimated by a Bayesian model using MCMCTREE with relaxed molecular clock and is

575 given in million years with 95% confidence intervals in round brackets. The split of *P.*

576 *westermani* was estimated to have occurred somewhere around 38.9 myr ago (28.0-58.6

577 myr). The analysis was repeated using BEAST 2 and estimated divergence times are shown in

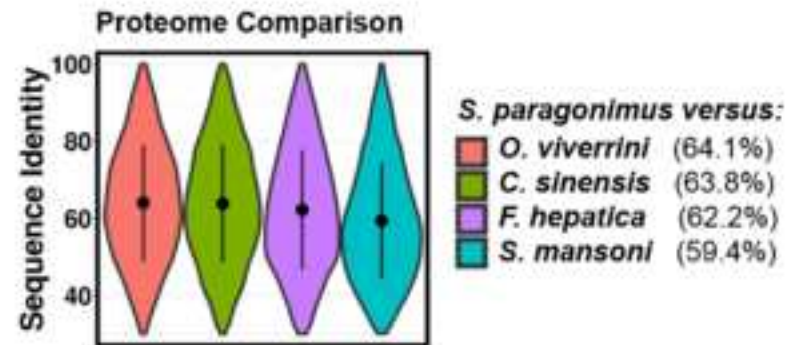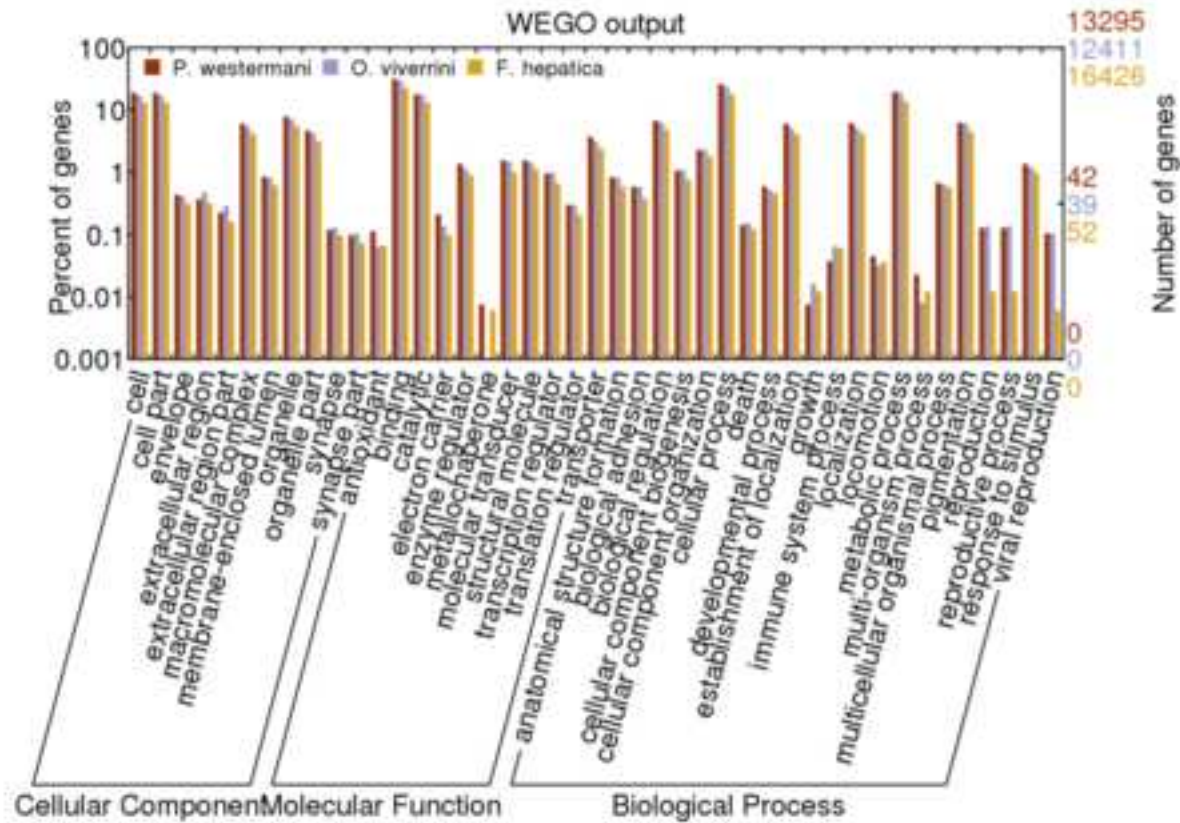578 square brackets. BEAST 2 estimated the split of *P. westermani* to have occurred 31.5 myr

579 ago.

21

Figure1
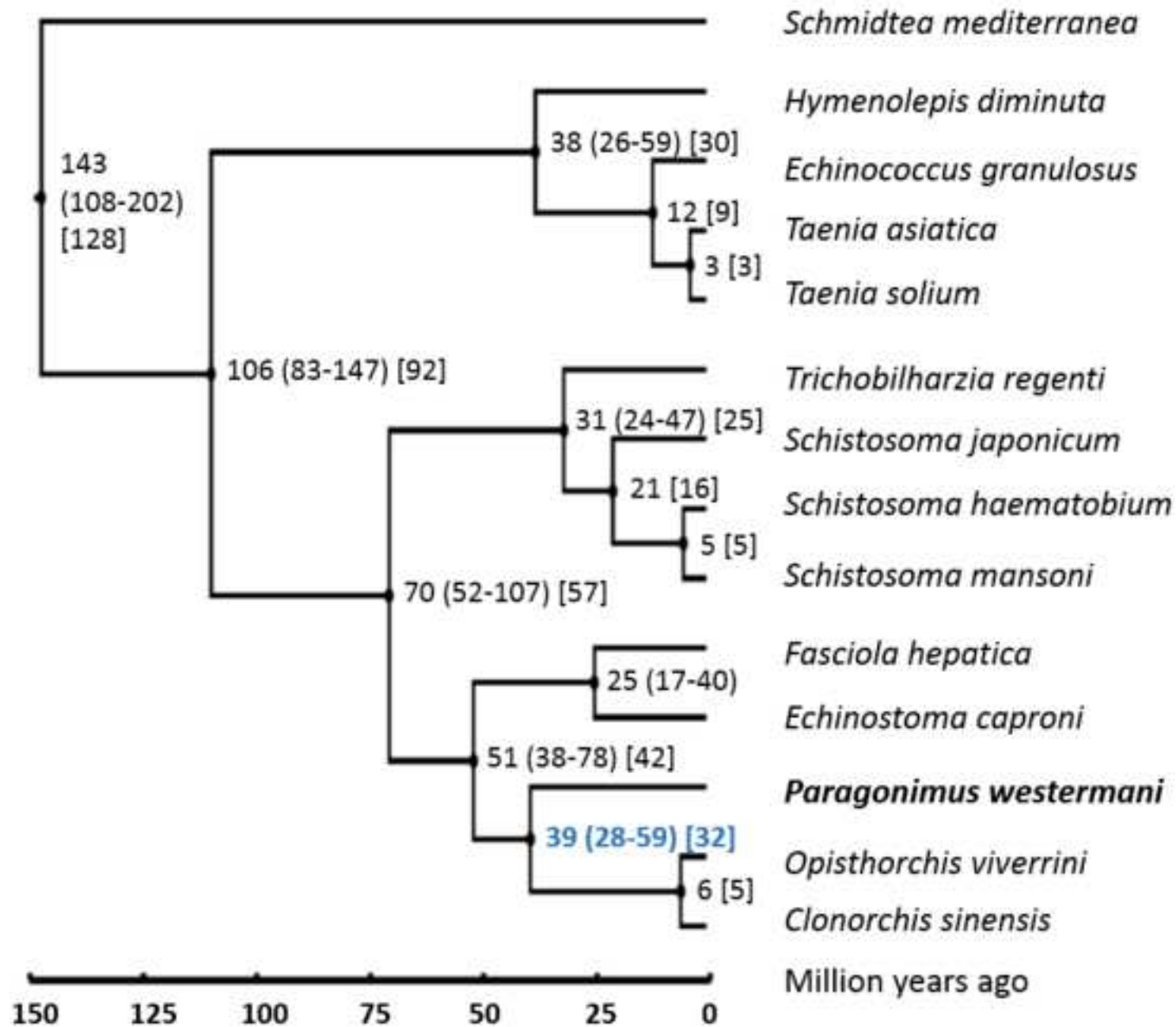
Figure2
Click here to access/download;Figure;Figure 2.jpg



Repeat 1
Repeat 2
Repeat 3

Figure3                    Click here to access/download;Figure;Figure 3.png ⬇

**A.**



**B.**

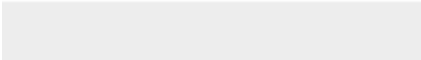

**C.**

Figure4

Figure 4

Click here to access/download
**Supplementary Material**
Additional file 1.pdf

Dear Scott Edmunds (Executive Editor),

Thank you for considering our manuscript "Complete genome sequence of the oriental lung fluke *Paragonimus westermani*" (GIGA-D-18-00193) for publication in GigaScience.

The reviewers were generally positive in their comments about the manuscript. However they did raise valuable points that we have addressed. We have prepared a detailed response to the points. Furthermore, we have modified sections of the manuscript and figures to address reviewer's questions. The changes have been highlighted (yellow) in the main text of the manuscript.

Both reviewers questioned the manuscript type. We have resubmitted the manuscript as a Data Note.

Yours sincerely,

Lutz Krause, PhD
Principal Research Fellow / Associate Professor
Head, Computational Medical Genomics Group
The University of Queensland Diamantina Institute