

Manuscript Number:	GIGA-D-18-00193	
Full Title:	Complete genome sequence of the oriental lung fluke <i>Paragonimus westermani</i>	
Article Type:	Research	
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:	<p>Background Foodborne infections caused by lung flukes of the genus <i>Paragonimus</i> are a significant and widespread public health problem in tropical areas. Around 50 <i>Paragonimus</i> species have been reported to infect animals and humans, but <i>Paragonimus westermani</i> is responsible for the bulk of human disease. Despite their medical and economic importance, no genome sequence for any <i>Paragonimus</i> species is available.</p> <p>Results We sequenced and assembled the genome of <i>P. westermani</i>, which is among the largest of the known pathogen genomes with an estimated size of 1.1 Gb. A 924.5 Mb genome assembly was generated from Illumina and PacBio sequence data. 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 <i>Opisthorchis viverrini</i> with an average sequence identity of 64.1%. Assembly of the <i>P. westermani</i> 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 <i>P. westermani</i>, possibly even within single worm isolates.</p> <p>Conclusions The generated assembly represents the first <i>Paragonimus</i> genome sequence and will facilitate future molecular studies of this important, but neglected, parasite group.</p>	
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:	
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 deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p>	Yes

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

1 **Title Page**

2 Complete genome sequence of the oriental lung fluke *Paragonimus*
3 *westermani*

4
5 Harald Oey^{1, #}, Martha Zakrzewski², Kanwar Narain⁶, K. Rekha Devi⁶, Takeshi Agatsuma⁷,
6 Sujeevi Nawaratna², Geoffrey Gobart^{2,4}, Malcolm Jones³, Mark A. Ragan⁵, 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

24 **Author names, affiliations and emails:**

25 **Dr Harald Oey**

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

28 Email: h.oey@uq.edu.au

30 **Dr Martha Zakrzewski**

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

32 Email: Martha.Zakrzewski@qimrberghofer.edu.au

34 **Dr Kanwar Narain**

35 Regional Medical Research Centre for Northeastern region, Indian Council of Medical
36 Research Department of Health Research Ministry of H & FW, Govt. of India, Dibrugarh-786
37 001, India

38 Email: kanwar_narain@hotmail.com

40 **Dr K. Rekha Devi**

41 Regional Medical Research Centre for Northeastern region, Indian Council of Medical
42 Research Department of Health Research Ministry of H & FW, Govt. of India, Dibrugarh-786
43 001, India

44 Email: krekha75@yahoo.co.in

46 **Prof Takeshi Agatsuma**

47 Department of Environmental Medicine

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

1
2
3 50
4
5 51 **Dr Geoffrey Gobert**
6
7
8 52 School of Biological Sciences, Queen's University Belfast, UK
9
10 53 Email: g.gobert@qub.ac.uk
11
12
13 54
14
15 55 **Prof Mark A. Ragan**
16
17
18 56 Institute for Molecular Bioscience, The University of Queensland, Brisbane, QLD, Australia
19
20
21 57 Email: m.ragan@uq.edu.au
22
23
24 58
25
26 59 **Prof Malcolm Jones**
27
28
29 60 School of Veterinary Science, University of Queensland, Gatton, QLD, Australia
30
31
32 61 Email: m.jones@uq.edu.au
33
34 62
35
36
37 63 **Dr Sujeevi Nawaratna**
38
39
40 64 QIMR Berghofer Medical Research Institute, Brisbane, QLD, Australia
41
42 65 Email: sujeevi.nawaratna@qimrberghofer.edu.au
43
44
45 66
46
47 67 **Prof Donald P. McManus**
48
49
50 68 QIMR Berghofer Medical Research Institute, Brisbane, QLD, Australia
51
52
53 69 Email: Don.McManus@qimrberghofer.edu.au
54
55 70
56
57
58 71
59
60
61
62
63
64
65

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

72 **A/Prof Lutz Krause**

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

74 Queensland, Translational Research Institute, Brisbane, QLD

75 Email: L.Krause@uq.edu.au

76 **Abstract**

77 **Background**

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

83 **Results**

84 We sequenced and assembled the genome of *P. westermani*, which is among the largest of
85 the known pathogen genomes with an estimated size of 1.1 Gb. A 924.5 Mb genome assembly
86 was generated from Illumina and PacBio sequence data. The genome has a high proportion
87 (45%) of repeat-derived DNA, particularly of the LINE and LTR subtypes, and the expansion of
88 these elements may explain some of the large size. We predicted 12,852 protein coding
89 genes, showing a high level of conservation with related trematode species. The majority of
90 proteins (80%) had homologs in the human liver fluke *Opisthorchis viverrini* with an average
91 sequence identity of 64.1%. Assembly of the *P. westermani* mitochondrial genome from long
92 PacBio reads resulted in a single high-quality circularized 20.6 kb contig. The contig harboured
93 a 6.9 kb region of non-coding repetitive DNA comprised of three distinct repeat units. Our
94 results suggest that the region is highly polymorphic in *P. westermani*, possibly even within
95 single worm isolates.

96 **Conclusions**

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

99 **Keywords**

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

105 **Background**

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

115 *Paragonimus* spp. have a complex life cycle. Unembryonated eggs are expelled by coughing
116 or passed with faeces and develop in water. Miracidia hatch from the eggs and penetrate a
117 freshwater snail, its first intermediate host. During several asexual developmental phases
118 inside the snail, miracidia develop into a sporocysts and then two redial generations occur,
119 the second of which giving rise to microcercous cercariae that escape into fresh water. These
120 crawling cercariae invade a species of crustacean, the second intermediate host, to encyst in
121 muscles and other sites and develop into metacercariae. Humans and other definitive hosts
122 become infected through consumption of raw or inadequately cooked freshwater crabs or
123 crayfish [5]. Ingested metacercariae excyst, penetrate through the gut and become
124 encapsulated in the lungs where they mature into hermaphroditic adult worms (7.5 mm to
125 12 mm in length) in 6-10 weeks [5]. Paragonimiasis can lead to a chronic inflammatory disease
126 of the lung and can trigger asthma- or tuberculosis-like symptoms [6-8]. In more severe cases
127 *Paragonimus* can infect the brain or central nervous system of the definitive host, leading to
128 headache, visual loss, and even death [1].

129 Paragonimiasis is commonly diagnosed by microscopic detection of parasite eggs in faeces or
130 sputum. The lack of sensitive and reliable diagnostic tests in conjunction with unspecific
131 disease symptoms often leads to delayed treatment with the drug of choice, praziquantel [8].
132 Despite their high medical, veterinary and economic importance, only limited information on
133 the molecular biology of *Paragonimus* is currently available. Recent transcriptome
134 sequencing studies have provided some information on the gene content of *Paragonimus* [9],

135 but until now no *Paragonimus* genome sequence has been available. Here we present a 924.5
136 Mb assembly of the *P. westermani* genome which provides new insights into the genomic
137 composition of the *Paragonimus* genus and represents an invaluable resource for future
138 studies of the neglected tropical disease paragonimiasis.

140 Data Description

141 Sequencing

142 Diploid *Paragonimus westermani* metacercariae were collected from the freshwater crab
143 *Maydelliathelphusa lugubris* in 2009 in the East Siang district of Arunachal Pradesh, Northeast
144 India, and fed to Wistar rats as experimental hosts. Genomic DNA was isolated from 50
145 individual worms, yielding 18 µg of DNA, and 2 µg of DNA was used to sequence the
146 *Paragonimus westermani* genome using a whole-genome shotgun approach. Paired-end
147 short-insert (200 bp and 450 bp) and mate-pair (5 Kb and 10 Kb) genomic DNA libraries were
148 sequenced on the Illumina HiSeq platform, yielding 58 Gb of sequence data (**Table 1**). For
149 genome scaffolding and quality evaluation of the assembled sequence, additional long-read
150 data were generated using the PacBio platform, yielding 1.7 Gb of information (**Table 1**). The
151 genome size was estimated from the K-mer coverage of the 450 bp insert library. K-mer
152 frequencies were calculated by the program Jellyfish[10], version 2.2.6, using a K-mer size of
153 17bp. The 17-mer distribution in the 450 bp library had a single peak at 26x (**Figure 1**),
154 demonstrating low sequence heterozygosity. The genome size (G) was deduced from the K-
155 mer distribution via the formula $G = N * (L - K + 1) / K_depth$ [11], where N is the total number
156 of reads, L is the read length, K is the K-mer size and K_depth is the peak frequency. With an
157 estimated size of 1.1 Gb the *P. westermani* genome is among the largest known pathogen
158 genomes and one of the largest parasite genomes sequenced to date. The genome is
159 considerably larger than the published genomes of the related trematodes *Clonorchis sinensis*
160 (644 Mb)[12], *Opisthorchis viverrini* (assembly size of 634.5 Mb) [13], and *Schistosoma* spp
161 (381-403 Mb) [14-16] and comparable to the 1.3 Gp genome of *Fasciola hepatica* [17].

164 **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

165
166
167 **Genome assembly**

168 PacBio sequence data were error corrected by proofread version 2.13.13 [18], using Illumina
 169 short reads from the 200bp and 450bp libraries as input , and assembled into contigs by Mira
 170 v4.0.2 [19]. Short-read Illumina sequence data were trimmed using Trimmomatic v0.36 and
 171 subsequently error corrected by KmerFreq_HA (part of SoapDenovo2 [20]) with a K-mer size
 172 of 23. The 10 kb mate-pair library showed a high proportion of PCR duplicates and was
 173 subjected to PCR de-duplication prior to genome assembly. Illumina paired-end sequence
 174 data were assembled using the ABYSS assembly pipeline [21], version 2.0.2, with options n=5
 175 s=200 N=36 S=500 k=33 and including the PacBio contigs via the re-scaffolding feature. The
 176 resulting assembly was de-gapped using the SoapDenovo2 GapCloser program [20]. Mate-
 177 pair libraries were then used to scaffold the assembly with SSPACE v3.0[22] (with options -x
 178 0 -a 0.60 -n 30 -z 200 -g 0) and gaps were again filled with GapCloser. To detect and resolve
 179 scaffolding errors, the resulting assembly was processed by the program REAPER [23] using
 180 the 5kb mate-pair library as input, breaking the assembly at sites with poor evidence for
 181 contiguity.

182 The assembly resulted in a 924.5 Mb genome sequence (30,997 scaffolds with N50 of 135 Kb)
 183 (**Table 2**), covering 84.0% of the estimated genome size. The GC content of the genome was
 184 43.3%, comparable to genomes of other related trematodes (**Table 2**). Genome assembly
 185 completeness was evaluated by BUSCO [24] using the metazoan lineage data, resulting in
 186 similar scores to those obtained for the genomes of comparable trematode species (**Table 2**).
 187 The proportion of duplicated genes reported by BUSCO was also similar to that of comparable

188 trematodes suggesting that the relatively large size of the *P. westermani* genome is not the
 189 result of genome duplication events.

190

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

	<i>P. westermani</i>	<i>F. hepatica</i>	<i>O. viverrini</i>	<i>C. sinensis</i>
Assembly size (Mb) ^a	924.5	1,275.0	606.0	546.9
Total base pairs (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 count	30,997 (>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%

193 ^aCombined length of all scaffolds. ^bCombined length of all scaffolds without gaps (N's). ^cNon-
 194 overlapping RNAseq-supported gene models[17]

196 Mitochondrial genome

197 The mitochondrial genome of *P. westermani* is present at a much higher copy number than
 198 the nuclear genome and we were able to assemble the full mitochondrial genome at high
 199 coverage from error corrected long PacBio reads using the Mira assembler[19], version 4.0.2.
 200 This resulted in a single mitochondrial contig of 20.3 Kb. The accuracy of the contig was
 201 confirmed by mapping short insert paired-end sequences directly onto the contig revealing
 202 single nucleotide discrepancies at only 4 positions. The mitochondrial genome was found to
 203 closely match previously published *Paragonimus* mitochondrial genomes with the best match
 204 from a blast search against the Nucleotide collection at the NCBI
 205 (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) being accession NC_027673.1, a *P. westermani*

206 complex sp. type 1 mitochondrial genome isolated in India (97% sequence identity across 13.4
207 Kb of NC_027673.1). This sequence was used as reference for mitochondrial gene
208 identification and annotation, supplemented by mitochondrial gene predictions by Mitos
209 [25]. The mitochondrial genomes of flatworms are known to harbour a region of non-coding
210 repetitive DNA, generally comprised of a long noncoding region (LNR) and a short non-coding
211 region (SNR) with a single tRNA gene separating them[26]. Reconstructing this region from
212 short-read data proved challenging, but our long-read PacBio data allowed complete
213 assembly of the repetitive region and circularization of the genome. Interestingly, our
214 assembled mitochondrial genome sequence had a much longer non-coding region (6.9 Kb)
215 than the previously published NC_027673.1 (0.7 Kb) and the non-coding regions of both
216 genomes showed only partial homology, but with close homology of the intervening tRNA
217 gene. We found the LNR to be comprised of two distinct repeat units with 8 and 13 copies
218 while the SNR was comprised of another distinct repeat unit with 3 copies. Strikingly, five
219 independent PacBio reads spanned the entirety of the non-coding region but with slight
220 differences in length (6.3 – 6.9 Kb), suggesting that the region is polymorphic, possibly even
221 within individual worms.

223 **Repeat annotation**

224 RepBase repeat consensus sequences did not adequately represent the repeats found in the
225 *P. westermani* assembly, consistent with the distant evolutionary relationship of lung flukes
226 with previously sequenced worm genomes. We therefore carried out *de-novo* repeat
227 characterization using the RepeatModeller package, version 1.0.9, and used the generated
228 consensus sequences to identify repetitive regions by RepeatMasker, version 4.0.7 (both
229 available at <http://www.repeatmasker.org>). To enable direct comparison with related
230 trematode species we also ran RepeatModeller and RepeatMasker separately on the *F.*
231 *hepatica*, *O. viverrini* and *C. sinensis* genomes with the same program parameters as those
232 used for *P. westermani*.

233 A relatively high percentage (45.2%) of the *P. westermani* genome sequence was repeat-
234 derived, similar to the rate reported for *Schistosoma* spp. (40.1-47.5%) [14-16] and *F. hepatica*
235 (57.1%), but considerably higher than the rate observed for the closer relatives *O. viverrini*

(28.9%) and *C. sinensis* (32.6%) (Table 3). Retrotransposons of the long interspersed nuclear element (LINE) subtype were found to be the greatest contributors of repetitive DNA (21.6%) (Table 3), consistent with reports for other trematode genomes [17]. In *P. westermani* and *F. hepatica*, the two largest of the four included trematode genomes, long terminal repeat (LTR) retrotransposons were also highly abundant contributing 7.7% and 10.1% of the genomes, respectively. Additionally, all four genomes had considerable amounts of repetitive DNA (10.7-17.1%) that did not match repeat consensus sequences of any of the known repeat classes modelled by RepeatModeler. The relatively high proportion of repeat-derived sequences in *P. westermani* may explain some of the increased size observed for this genome compared to the genomes of related flatworm species.

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

Repeat class	<i>P. westermani</i>	<i>F. hepatica</i>	<i>O. viverrini</i>	<i>C. sinensis</i>
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

Gene prediction and functional annotation

Genes were predicted by the Maker pipeline, version 2.31.9, using Augustus [27], version 3.2.3, and GeneMark-ES [28], version 4.32, for *ab-initio* gene prediction. To accurately model the sequence properties of the *P. westermani* genome, both gene finders were initially trained by BRAKER1 [29], version 1.9, which makes use of mapped transcriptome sequence data. Previously published *P. westermani* RNA-seq data [9] were obtained from the short read archive and mapped to our genome assembly using the Star aligner [30], version 2.5, with the option `--twopassMode Basic`. BRAKER1 was then run with default parameters. The RNA-seq data was further assembled into transcripts using cufflinks [31], version 2.2.1, with the

options --frag-bias-correct <p.westermani assembly> --multi-read-correct. The resulting transcripts were provided as input for Maker via the "est_gff" option. For homology based searches Maker was provided with the following wormbase v8 protein datasets: *Clonorchis sinensis* (PRJDA72781), *Opisthorchis viverrini* (PRJNA222628), *Schistosoma mansoni* (PRJEA36577), *Caenorhabditis elegans* (PRJNA13758), *Echinococcus granulosus* (PRJEB121), *Hymenolepis diminuta* (PRJEB507) and *Schistosoma haematobium* (PRJNA78265). Additionally, the Swiss-Prot dataset from UniProt was included. Maker was allowed to report single exon genes, and otherwise run with default parameters.

Proteins were functionally annotated based on a BLASTp search against the NCBI non-redundant protein database (obtained on 25.10.17) requiring an e-value <1e-15 and the best hit spanning at least 40% of the query sequence. Additionally, functional domains, GO annotations, transmembrane proteins and signal peptides were identified with InterProScan [32], version 5.25-64.0. GO annotations were then visualized using WEGO [33]. In total, 12,852 protein encoding genes were predicted in the *P. westermani* genome and functionally annotated (**Table 2**).

Genome comparison

Predicted *P. westermani* coding genes were mapped to the genomes of related trematode species using Exonerate, version 2.4.0, requiring a minimal sequence identity of 30% and excluding matches spanning less than 40% of the query protein. The majority of predicted proteins (86.2%) had inferred homologs in the related trematode species (**Figure 3A**) and showed a similar distribution of protein functional categories (**Figure 3C**). The *P. westermani* proteome was most similar to *O. viverrini* and *C. sinensis*. Of the 12,852 predicted proteins, 10,350 (80%) had inferred homologs in *O. viverrini* with an average sequence identity of 64.1%, and 10,227 (79.6%) had homologs in *C. sinensis* with an average sequence identity of 63.8% (**Figures 3A and 3B**).

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-

289 living turbellarian flatworm, as outgroup (**Figure 4**). We first identified single-copy proteins
290 shared across all 14 included worm species. Single-copy proteins were identified based on
291 blastp searches of a species proteins against the species own proteome using a sequence-
292 identity cut-off of 30% and requiring hits to cover >50% of the query sequence. Single-copy
293 proteins shared across all 14 species were then identified using a less stringent blastp search
294 with a 30% sequence identity cut-off but requiring only >40% coverage of the query sequence.
295 We identified 104 single-copy proteins shared across the 14 worm species that were then
296 aligned using MUSCLE [34]. The resulting multiple sequence alignment was de-gapped with
297 trimAl [35] and a phylogenetic tree was reconstructed by PHYLIP v3.696 using the maximum
298 likelihood method and the Jones-Taylor-Thornton probability model.

299 The multiple alignment and the inferred phylogenetic tree were then used to estimate species
300 divergence by a Bayesian model with relaxed molecular clock using MCMCTREE in PAML 4.9e
301 (**Figure 4**). The model was calibrated based on previously published divergence times and
302 ages of fossil records. Evidence for trematode infestation have been reported from the
303 Eocene (56 to 33.9 million years (myr) ago) and preserved trematode eggs have been found
304 in dinosaur coprolites from the Early Cretaceous (146 to 100 myr ago); however, fossil records
305 indicate that trematodes may have already existed more than 400 myr ago [36, 37]. The
306 trematode split from other neodermatan lineages was therefore fixed at >56 myr. The origin
307 of schistosomes has been estimated somewhere in the Miocene around 15-20 myr ago [38,
308 39]. It has further been estimated that the divergence of *S. mansoni* did likely not occur before
309 2-5 myr ago, based on fossil records of its intermediate host *Biomphalaria* [40]. From these
310 data, the split of Plagiorchiida (including *P. westermani*) and Opistorchiida (including *O.*
311 *viverrini* and *C. sinensis*) was estimated to have occurred 38.9 myr ago (95% confidence
312 interval of 28.0-58.6 myr) (**Figure 4**).

313

314 **Conclusion**

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

318

1
2
3 319 **Availability of supporting data**
4

5 320 The nuclear and mitochondrial genomes are available from NCBI under accession number

6 321 PRJNA454344
7

8 322
9

10 323 **Abbreviations**
11

12
13 324 LTR Long terminal repeat
14

15
16 325 SINE Short interspersed nuclear elements
17

18 326 LINE Long interspersed nuclear elements
19

20
21 327 MYR Million years
22

23
24 328
25

26
27 329 **Competing interests**
28

29
30 330 All authors declare that they have no competing interests.
31

32
33 331
34

35 332 **Funding**
36

37
38 333 This work has been supported by grants from the QIMR Berghofer Medical Research Institute
39

40 334 (Chenhall Estate) and the Australian Infectious Diseases Research Centre. DPM is a National
41

42 335 Health and Medical Research Council Senior Principal Research Fellow.
43
44

45 336
46

47
48 337 **Author contributions**
49

50
51 338 LK and DPM conceived and managed the project; KN and KRD provided *P. westermani*
52

53 339 material. TA and SN isolated genomic DNA. MZ and GG managed DNA sequencing. HO carried
54

55 340 out genome assembly, gene prediction and functional genome annotation. HO and LK carried
56

57 341 out comparative genomics. LK, DPM, MKJ and MAR attracted funding and designed the study.
58
59
60
61
62
63
64
65

342 LK and HO drafted the manuscript and all authors read, edited and approved the final
343 manuscript.

344

345 **Acknowledgements**

346 We wish to acknowledge the QIMR Berghofer Medical Research Institute (Chenhall Estate
347 donation) and the Australian Infectious Diseases Research Centre for funding this work.

348

349 **References**

- 350 1. Furst T, Keiser J and Utzinger J. Global burden of human food-borne trematodiasis: a
351 systematic review and meta-analysis. *Lancet Infect Dis.* 2012;12 3:210-21.
352 doi:10.1016/S1473-3099(11)70294-8.
- 353 2. Blair D. Paragonimiasis. *Adv Exp Med Biol.* 2014;766:115-52. doi:10.1007/978-1-4939-0915-
354 5_5.
- 355 3. Roy JS, Das PP, Borah AK and Das JK. Paragonimiasis in a Child from Assam, India. *J Clin Diagn*
356 *Res.* 2016;10 4:DD06-7. doi:10.7860/JCDR/2016/18160.7616.
- 357 4. Singh TS, Hiromu S, Devi KR and Singh WA. First case of *Paragonimus westermani* infection in
358 a female patient in India. *Indian J Med Microbiol.* 2015;33 Suppl:156-9. doi:10.4103/0255-
359 0857.150950.
- 360 5. Jones MK, Keiser J and McManus DP. Trematodes. In: Jorgensen JH, Pfaller MA, Carroll KC,
361 Funke G, Landry ML, Richter SS, et al., editors. *Manual of Clinical Microbiology*, Eleventh
362 Edition. American Society of Microbiology; 2015.
- 363 6. Luo J, Wang MY, Liu D, Zhu H, Yang S, Liang BM, et al. Pulmonary Paragonimiasis Mimicking
364 Tuberculous Pleuritis: A Case Report. *Medicine (Baltimore).* 2016;95 15:e3436.
365 doi:10.1097/MD.0000000000003436.
- 366 7. Zhou R, Zhang M, Cheng N and Zhou Y. Paragonimiasis mimicking chest cancer and
367 abdominal wall metastasis: A case report. *Oncol Lett.* 2016;11 6:3769-71.
368 doi:10.3892/ol.2016.4434.
- 369 8. Kalhan S, Sharma P, Sharma S, Kakria N, Dudani S and Gupta A. *Paragonimus westermani*
370 infection in lung: A confounding diagnostic entity. *Lung India.* 2015;32 3:265-7.
371 doi:10.4103/0970-2113.156248.
- 372 9. Li BW, McNulty SN, Rosa BA, Tyagi R, Zeng QR, Gu KZ, et al. Conservation and diversification
373 of the transcriptomes of adult *Paragonimus westermani* and *P. skrjabini*. *Parasit Vectors.*
374 2016;9:497. doi:10.1186/s13071-016-1785-x.
- 375 10. Marçais G and Kingsford C. A fast, lock-free approach for efficient parallel counting of
376 occurrences of k-mers. *Bioinformatics.* 2011;27 6:764-70.
377 doi:10.1093/bioinformatics/btr011.
- 378 11. Song L, Bian C, Luo Y, Wang L, You X, Li J, et al. Draft genome of the Chinese mitten crab,
379 *Eriocheir sinensis*. *Gigascience.* 2016;5:5. doi:10.1186/s13742-016-0112-y.
- 380 12. Wang X, Chen W, Huang Y, Sun J, Men J, Liu H, et al. The draft genome of the carcinogenic
381 human liver fluke *Clonorchis sinensis*. *Genome Biol.* 2011;12 10:R107. doi:10.1186/gb-2011-
382 12-10-r107.

- 383 13. Young ND, Nagarajan N, Lin SJ, Korhonen PK, Jex AR, Hall RS, et al. The *Opisthorchis viverrini*
1 384 genome provides insights into life in the bile duct. *Nat Commun.* 2014;5:4378.
2 385 doi:10.1038/ncomms5378.
- 3 386 14. Young ND, Jex AR, Li B, Liu S, Yang L, Xiong Z, et al. Whole-genome sequence of *Schistosoma*
4 387 *haematobium*. *Nat Genet.* 2012;44 2:221-5. doi:10.1038/ng.1065.
- 5 388 15. Berriman M, Haas BJ, LoVerde PT, Wilson RA, Dillon GP, Cerqueira GC, et al. The genome of
6 389 the blood fluke *Schistosoma mansoni*. *Nature.* 2009;460 7253:352-8.
7 390 doi:10.1038/nature08160.
- 8 391 16. *Schistosoma japonicum* Genome S and Functional Analysis C. The *Schistosoma japonicum*
9 392 genome reveals features of host-parasite interplay. *Nature.* 2009;460 7253:345-51.
10 393 doi:10.1038/nature08140.
- 11 394 17. Cwiklinski K, Dalton JP, Dufresne PJ, La Course J, Williams DJ, Hodgkinson J, et al. The
12 395 *Fasciola hepatica* genome: gene duplication and polymorphism reveals adaptation to the
13 396 host environment and the capacity for rapid evolution. *Genome Biol.* 2015;16:71.
14 397 doi:10.1186/s13059-015-0632-2.
- 15 398 18. Hackl T, Hedrich R, Schultz J and Forster F. proovread: large-scale high-accuracy PacBio
16 399 correction through iterative short read consensus. *Bioinformatics.* 2014;30 21:3004-11.
17 400 doi:10.1093/bioinformatics/btu392.
- 18 401 19. Chevreux B, Pfisterer T, Drescher B, Driesel AJ, Muller WE, Wetter T, et al. Using the miraEST
19 402 assembler for reliable and automated mRNA transcript assembly and SNP detection in
20 403 sequenced ESTs. *Genome research.* 2004;14 6:1147-59. doi:10.1101/gr.1917404.
- 21 404 20. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved
22 405 memory-efficient short-read de novo assembler. *Gigascience.* 2012;1 1:18.
23 406 doi:10.1186/2047-217X-1-18.
- 24 407 21. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ and Birol I. ABySS: a parallel assembler
25 408 for short read sequence data. *Genome research.* 2009;19 6:1117-23.
26 409 doi:10.1101/gr.089532.108.
- 27 410 22. Boetzer M, Henkel CV, Jansen HJ, Butler D and Pirovano W. Scaffolding pre-assembled
28 411 contigs using SSPACE. *Bioinformatics.* 2011;27 4:578-9. doi:10.1093/bioinformatics/btq683.
- 29 412 23. Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M and Otto TD. REAPR: a universal tool
30 413 for genome assembly evaluation. *Genome Biol.* 2013;14 5:R47. doi:10.1186/gb-2013-14-5-
31 414 r47.
- 32 415 24. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing
33 416 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.*
34 417 2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.
- 35 418 25. Bernt M, Donath A, Juhling F, Externbrink F, Florentz C, Fritzsche G, et al. MITOS: improved de
36 419 novo metazoan mitochondrial genome annotation. *Mol Phylogenet Evol.* 2013;69 2:313-9.
37 420 doi:10.1016/j.ympev.2012.08.023.
- 38 421 26. Le TH, Blair D and McManus DP. Mitochondrial genomes of parasitic flatworms. *Trends*
39 422 *Parasitol.* 2002;18 5:206-13.
- 40 423 27. Stanke M, Steinkamp R, Waack S and Morgenstern B. AUGUSTUS: a web server for gene
41 424 finding in eukaryotes. *Nucleic Acids Res.* 2004;32 Web Server issue:W309-12.
42 425 doi:10.1093/nar/gkh379.
- 43 426 28. Lomsadze A, Ter-Hovhannisyanyan V, Chernoff YO and Borodovsky M. Gene identification in
44 427 novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res.* 2005;33 20:6494-
45 428 506. doi:10.1093/nar/gki937.
- 46 429 29. Hoff KJ, Lange S, Lomsadze A, Borodovsky M and Stanke M. BRAKER1: Unsupervised RNA-
47 430 Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics.* 2016;32
48 431 5:767-9. doi:10.1093/bioinformatics/btv661.
- 49 432 30. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal
50 433 RNA-seq aligner. *Bioinformatics.* 2013;29 1:15-21. doi:10.1093/bioinformatics/bts635.

434 31. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript
1 435 assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform
2 436 switching during cell differentiation. *Nat Biotechnol.* 2010;28 5:511-5. doi:10.1038/nbt.1621.
3 437 32. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale
4 438 protein function classification. *Bioinformatics.* 2014;30 9:1236-40.
5 439 doi:10.1093/bioinformatics/btu031.
6 440 33. Ye J, Fang L, Zheng H, Zhang Y, Chen J, Zhang Z, et al. WEGO: a web tool for plotting GO
7 441 annotations. *Nucleic Acids Res.* 2006;34 Web Server issue:W293-7. doi:10.1093/nar/gkl031.
8 442 34. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
9 443 *Nucleic Acids Res.* 2004;32 5:1792-7. doi:10.1093/nar/gkh340.
10 444 35. Capella-Gutierrez S, Silla-Martinez JM and Gabaldon T. trimAl: a tool for automated
11 445 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics.* 2009;25 15:1972-3.
12 446 doi:10.1093/bioinformatics/btp348.
13 447 36. Poinar G, Jr. and Boucrot AJ. Evidence of intestinal parasites of dinosaurs. *Parasitology.*
14 448 2006;133 Pt 2:245-9. doi:10.1017/S0031182006000138.
15 449 37. Huntley JW and De Baets K. Trace Fossil Evidence of Trematode-Bivalve Parasite-Host
16 450 Interactions in Deep Time. *Adv Parasitol.* 2015;90:201-31. doi:10.1016/bs.apar.2015.05.004.
17 451 38. Littlewood DTJe and Baets Kde. Fossil parasites.
18 452 39. Snyder SD and Loker ES. Evolutionary relationships among the Schistosomatidae
19 453 (Platyhelminthes:Digenea) and an Asian origin for *Schistosoma*. *J Parasitol.* 2000;86 2:283-8.
20 454 doi:10.1645/0022-3395(2000)086[0283:ERATSP]2.0.CO;2.
21 455 40. Morgan JA, Dejong RJ, Snyder SD, Mkoji GM and Loker ES. *Schistosoma mansoni* and
22 456 *Biomphalaria*: past history and future trends. *Parasitology.* 2001;123 Suppl:S211-28.
23
24
25
26
27
28 457
29
30
31 458
32

33 **Figure legends:**
34
35

36 460
37
38
39 461 **Figure 1. K-mer frequencies for the 450bp library.** Distribution of 17-mers in the 450bp
40 462 short-insert library demonstrated low sequence heterozygosity. We observed a single peak
41 463 at 26x and the *P. westermani* genome size was estimated to be 1.1 Gb.
42
43
44
45 464
46
47

48 465 **Figure 2. The complete *P. westermani* mitochondrial genome.**
49

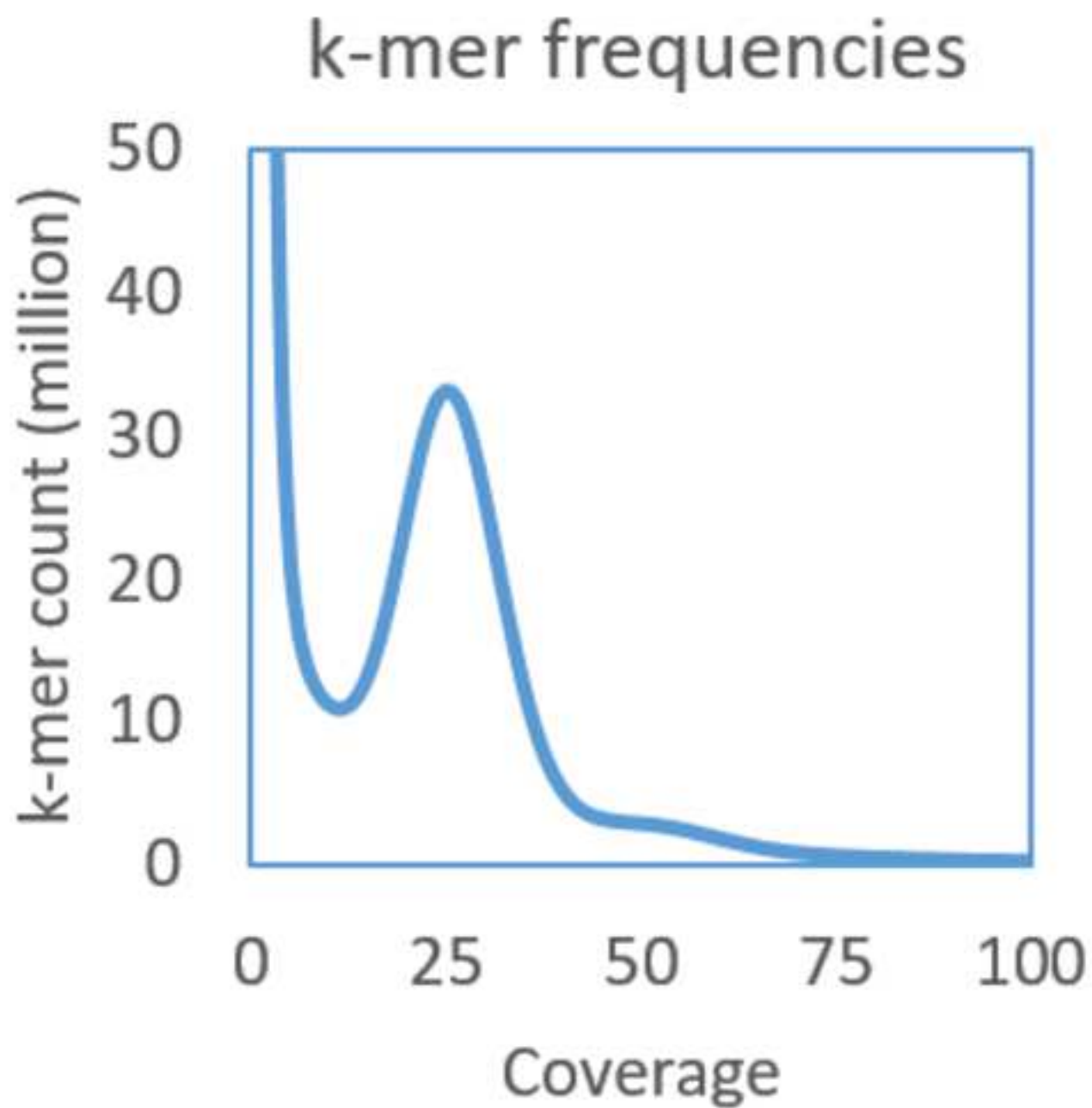
50 466 A) A graphical representation of the *P. westermani* circular mitochondrial genome is shown,
51 467 including a ~6.9 Kb repetitive region. Three distinct repeat units were identified in this region,
52 468 as well as an intervening tRNA gene (tRNA-Glu). All genes are transcribed in the clock-wise
53 469 direction. B) The consensus sequence for three repeat units identified in the ~6.9 Kb repetitive
54 470 region.
55
56
57
58
59

471

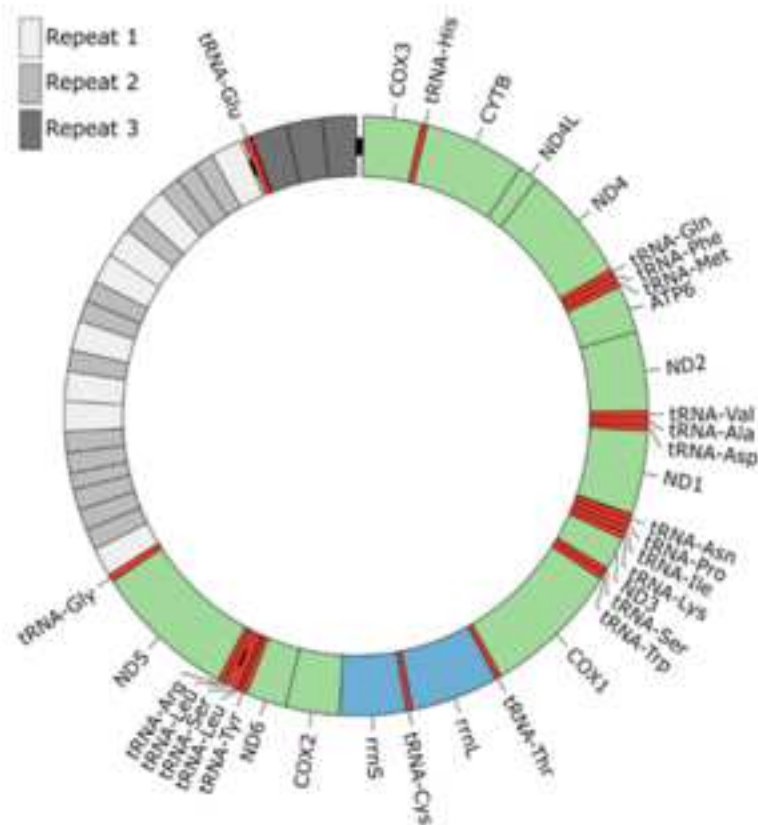
1
2
3 472 **Figure 3. Conservation of the *P. westermani* proteome across four related trematode**
4 473 **species.** *P. westermani* proteins were mapped to the genome sequences of *O. viverrini*, *C.*
5
6 474 *sinensis*, *F. hepatica* and *S. mansoni* using Exonerate. A) *P. westermani* centred Venn diagram
7
8 475 of 12,852 predicted proteins. The four included trematode species shared a core set of 7,599
9
10 476 proteins. B) Sequence identity of *P. westermani* proteins and orthologues inferred in genomes
11
12 477 of related trematodes. Average sequence identity is given in brackets. C) Distribution of
13
14 478 identified functional GO categories across three trematode species. GO annotations were
15
16 479 assigned by InterProScan and visualized using WEGO.

17
18
19 480

20
21 481 **Figure 4. Phylogenetic tree and estimated divergence times.** A phylogenetic tree of selected
22
23 482 trematodes and cestodes and *S. mediterranea* as outgroup was reconstructed from 104
24
25 483 shared single-copy proteins using the maximum likelihood method. Species divergence was
26
27 484 estimated by a Bayesian model with relaxed molecular clock and is given in million years with
28
29 485 95% confidence intervals in brackets. The split of *P. westermani* was estimated to have
30
31 486 occurred somewhere around 38.9 myr ago (28.0-58.6 myr).



A)



B)

Repeat 1 (328bp ~8 copies)

```
TGTC AAGTTTGAAGGGACCGATTTCGATTC CAATGGGTGTAGAGGTTTGGAGTTGGCGTTGCCGTGTGATTTTCTGTGTCAA
GGGGGGTCTGAAACTATGCGCGAAAAGGGTGC AAAAAAATTCGTAAGGGGGGGGCATTGCAAACTTTTTCCTTTTAAAAATTTAC
AGCTTAATTCAGGTC TAGTCGAAGAGTGAAGTGGTTTTTATCTCCOCTTAATTTGACTGTGCATTAAAAATTTTCGTTACTTTTGTG
TCAAAATTACATCATAGCTTTTTTCAGGGGAGTTCGGAGGTGAAAAGTTGGATTTTTGAAGGGTTTG
```

Repeat 2 (229bp ~13 copies)

```
TGTC AAGTTTGAAGGGACCGATTTCGATTC CAATGGGTGTAGAGGTTTGGAGTTGGCGTTGCCGTGTGATTTTCTGTGTCAA
GGGGGGTTTTAAACTATGCTGTGCAGGGGTGTTACCGTAGCTTTTTTCAGGGGAGTTCGGAGGTGAAAAGTTCGGTTTTTTCGATGA
GCTGGTACGAAATGCTATTATGTTAATCATAAGTAGAGTTATAATTAGAGGTCTC
```

Repeat 3 (406bp 3 copies)

```
AAAAAAGATATCATCGCTAAAAGAGAATAATTGGAAATGACTGTGTTGTCGTAAAGGATTGGGATTAAGTGTAGTCCGAGCAGTGT
TGTTTGTGCGAGGAAGAGATGAGGGGGCTTAATAATAATATGAGGAGGCCOCTACAATGAAAAGTGGTTAAGAAGGTCCATTGTAAA
GGATTAGGATTAAGTGTAGTCCGAGCAGTGTGTTGGTTGTGOGAGGAAGAGATGAGGGGGGCTTAATAATAATATGAGGAGGCCOCTAC
AATGAAAAGTGGTTAAGAAGGTCCATTGTAAAAGATTAGGATTAAGTGTAGTCCGAGCAGTGTGTTGGTTGTGCGAGGAAGAGATGAGGG
GATCTTTGGTTAATAATACAGGAAGTCATCTGTAATGGGGGGAAAAGGGGGGGCTGACC
```

