

# 1 Supplement

2	Supplementary Note 1.	Individuals and Sequencing.....	2
3	Supplementary Table 1.	Sequenced raw data.....	2
4	Supplementary Table 2.	Summary statistics of different assembly steps.....	3
5	Supplementary Note 2.	Re-mapping .....	3
6	Supplementary Figure 1.	Number of partially mapped reads along continuous parts of scaffolds .....	4
7	Supplementary Figure 2.	Mapping quality frequency distribution .....	5
8	Supplementary Table 3.	Re-mapping statistics from mate pair libraries.....	5
9	Supplementary Note 3.	Results from flow cytometric measurements .....	6
10	Supplementary Note 4.	Genome size estimation from coverage .....	6
11	Supplementary Note 5.	Repeat content .....	6
12	Supplementary Figure 3.	Classified repeat families .....	7
13	Supplementary Figure 4.	Mollusc genome sizes .....	8
14	Supplementary Table 4.	Proteins similar to Swiss-Prot entries.....	9
15	Supplementary Figure 5.	Number of sequences from one species in orthogroups .....	9
16	Supplementary Figure 6.	Regression of protein sequences and orthogroups .....	10
17	Supplementary Table 5.	Protein sets used for ortholog clustering and GO-term enrichment .....	11
18	Supplementary Note 6.	Orthologous clustering and Gene Ontology enrichment .....	12
19	Supplementary Table 6.	Enriched GO-terms for <i>Radix</i> .....	13
20	Supplementary Table 7.	Enriched GO-terms not in <i>Radix</i> .....	13
21	Supplementary Note 7.	Preprocessing and trimming .....	13
22	Supplementary Note 8.	Contamination screening.....	14
23	Supplementary Figure 7.	Results of contamination screening.....	15
24	Supplementary Note 9.	Material and methods of flow cytometric analysis.....	15
25	Supplementary Note 10.	Transcriptome assemblies .....	16
26	Supplementary Note 11.	Genome assembly .....	17
27	Supplementary Figure 8.	Transcriptome filtering.....	20
28	Supplementary Table 8.	Genome scaffolding with transcriptomic data.....	21
29	Supplementary Note 12.	Repeat library .....	21
30	Supplementary Figure 9.	Read subsampling .....	23
31	Supplementary Table 9.	Summarized results of the CEGMA analysis.....	23
32	Supplementary Note 13.	Annotation.....	24
33			

34 **Supplementary Note 1.** Individuals and Sequencing

35 Three snails from the inbred line were used for DNA extraction. Pooling of tissue from whole snails,  
36 DNA extraction, library construction of paired libraries with insert sizes of 250, 500, 800, 2k, 5k and  
37 10k as well as sequencing using Illumina HiSeq 2500 technology with read length of 125bp for small  
38 insert libraries 200, 500 and 800bp and Illumina HiSeq 2000 with read length 100bp for mate pair  
39 libraries 2 kb, 5 kb and 10 kb was performed by BGI, Hong-Kong. In total more than one billion reads  
40 (1,000,372,010) containing more than 116Gb (116,162,940,950bp) raw data was produced.

41

**Supplementary Table 1.** Sequenced raw data. The paired-end libraries with insert sizes of 250 bp, 500 bp and 800 bp were each sequenced with read lengths of 125 bp (HiSeq 2500). The mate pair libraries with insert sizes of 2 kb, 5 kb und 10 kb were sequenced with read lengths of 100 bp (HiSeq2000).

Insert size	Number of sequences	Number of nucleotides	%GC	Coverage
250 bp	289,883,600	36,235,450,000	39	23
500 bp	197,642,448	24,705,306,000	36	15
800 bp	157,503,550	19,687,943,750	36	12
2,000 bp	168,299,780	16,829,978,000	39	11
5,000 bp	135,963,092	13,596,309,200	39	8
10,000 bp	51,079,540	5,107,954,000	41	3

42

**Supplementary Table 2.** Summary statistics of different assembly steps.

Assembly step	Number of Sequences	Total length [bp]	N50 [bp]	% Ns
Raw reads	1,000,372,010	116,162,940,950		0
Trimmed reads	994,535,287	115,378,553,461		0
Platanus				
Assembly	6,838,932	1,488,367,542	324	0
Scaffolding	193,639	966,366,534	259,302	15.27
Gap close	193,639	927,196,599	250,725	9.64
Length filter $\geq$ 500 bp	22,306	898,221,812	262,000	9.94
SSPACE	10,317	909,612,132	512,264	11.05
Removing mitochondrial scaffold	10,316	909,598,491	512,264	11.05
L_RNA_scaffolder				
MOTU4	10,268	909,604,080	518,249	11.05
MOTU5	10,036	909,629,612	555,879	11.06
MOTU2+3	9,965	909,636,872	575,006	11.06
Length filter $\geq$ 1 kb	4,825	906,300,918	576,630	11.09
GapFiller	4,825	909,751,983	578,730	6.42
Adding separate processed mitochondrial scaffold (cutting position)	4,826	909,765,727	578,730	6.42
Remove gap containing scaffolds $<$ 1 kb	4,823	909,764,068	578,730	6.42

43

44 **Supplementary Note 2.** Re-mapping

45 All trimmed genomic reads were mapped unpaired against the final genome assembly using BWA  
46 mem 0.7.12-r1039 (Li 2013) with the options -a -c 10000. All other parameters were kept as default.

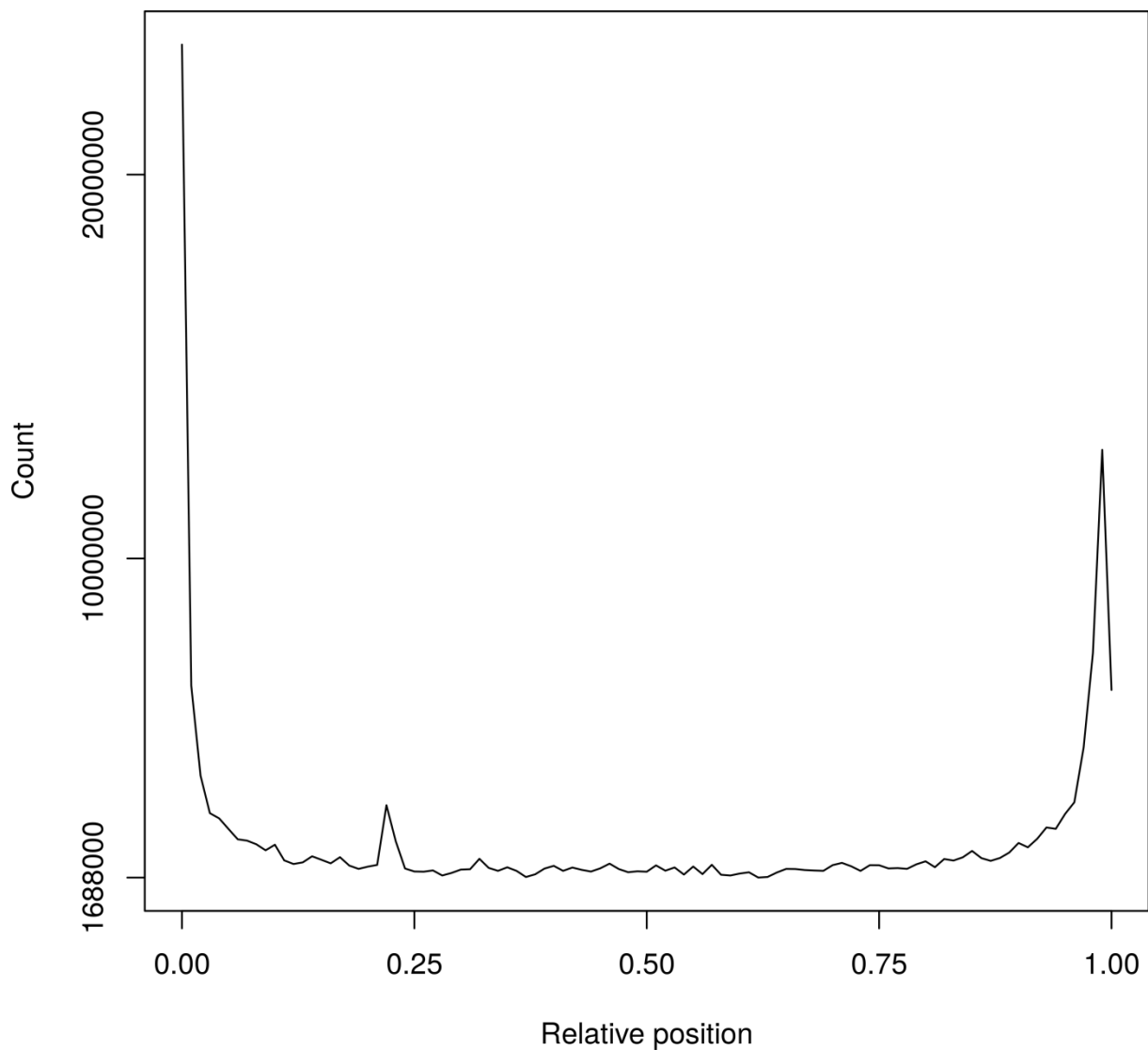
47 The coverage per position was calculated with samtools 1.1 mpileup and -A -C 50 -d 10000 except  
48 default parameters. Nearly all reads (97.64 %) mapped back but only 66.24 % of all nucleotides. The  
49 relatively low fraction of remapped nucleotides is due to the mapping quality cutoff of 50 during the  
50 mpileup step (94.9% without cutoff) and split mappings, which accumulate at the ends of contigs  
51 (Supplementary Figure 1). Split mappings are probably caused by incomplete assembled repeats or  
52 algorithmic problems to place a read correctly in a repeat. Per base coverage frequency distribution is  
53 shown in Figure 2A and the mapping quality frequency distribution in Supplementary Figure 2.

54

55 *Re-mapping mate pairs*

56 The trimmed mate pair reads of each library were mapped separately in paired mode against the final  
57 genome assembly using BWA mem with all default parameters. The mappings were sorted by position  
58 with samtools sort 1.1 (Li et al. 2009) and computation of statistics was realized in QualiMap bamqc  
59 2.2 (Okonechnikov et al. 2015; Figure 2B). Accumulation of insert sizes on the lower end of the  
60 distribution and higher mean coverage than expected (Supplementary Table 3) can be explained by  
61 mate pairs that cannot span repetitive regions.

62

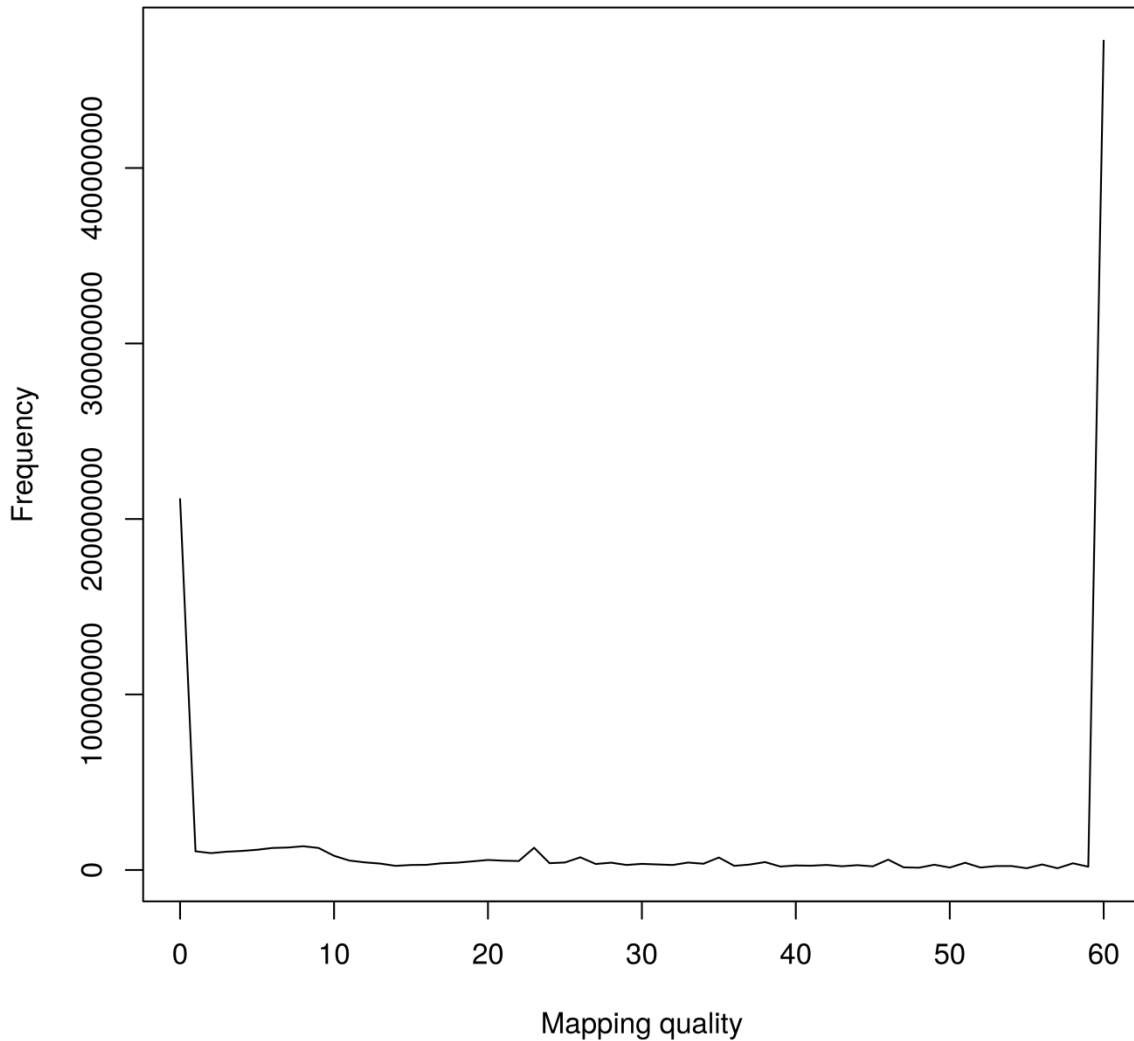


63

64 **Supplementary Figure 1.** Number of partially mapped reads along continuous parts of scaffolds.

65

## Mapping quality of primary alignments



66

67 **Supplementary Figure 2.** Mapping quality frequency distribution. Secondary and supplementary  
68 alignments are excluded.

**Supplementary Table 3.** Re-mapping statistics from mate pair libraries. Expected coverages were calculated from the number of nucleotides from each trimmed library and the estimated genome size of 1.6 Gb.

Library	Reads mapped [%]	Both in pair [%]	Mean / Expected Coverage
2k	98.42	97.88	17.27 / 10.51
5k	98.62	98.04	14.01 / 8.49
10k	98.28	97.67	5.24 / 3.19

69

70 **Supplementary Note 3.** Results from flow cytometric measurements.

71 Genome size estimations using different standards yielded comparable 2C-values. Mean 2C-value  
72 estimated using the standard *Glycine max* was 3.22 pg ( $\pm$  0.02 s.d.), using *Lycopersicon esculentum*  
73 3.19 pg ( $\pm$  0.01 s.d.). These values correspond to 3149.16 Mb and 3119.82 Mb, respectively. The CVs  
74 for the G0/G1 peak of the analysed samples ranged from 1.74 to 3.92% (mean 2.56).

75

76 **Supplementary Note 4.** Genome size estimation from coverage

77 The number of total trimmed nucleotides which were used in the assembly divided by the maximum of  
78 the per-position coverage frequency distribution (Figure 2A) is an estimate for the genome size,  
79 assuming even sequencing coverage throughout the genome.

$$\frac{\text{Total trimmed nucleotides}}{\text{peak coverage}} = \frac{115.35\text{Gb}}{72} = 1.6025\text{Gb}$$

80

81 **Supplementary Note 5.** Repeat content

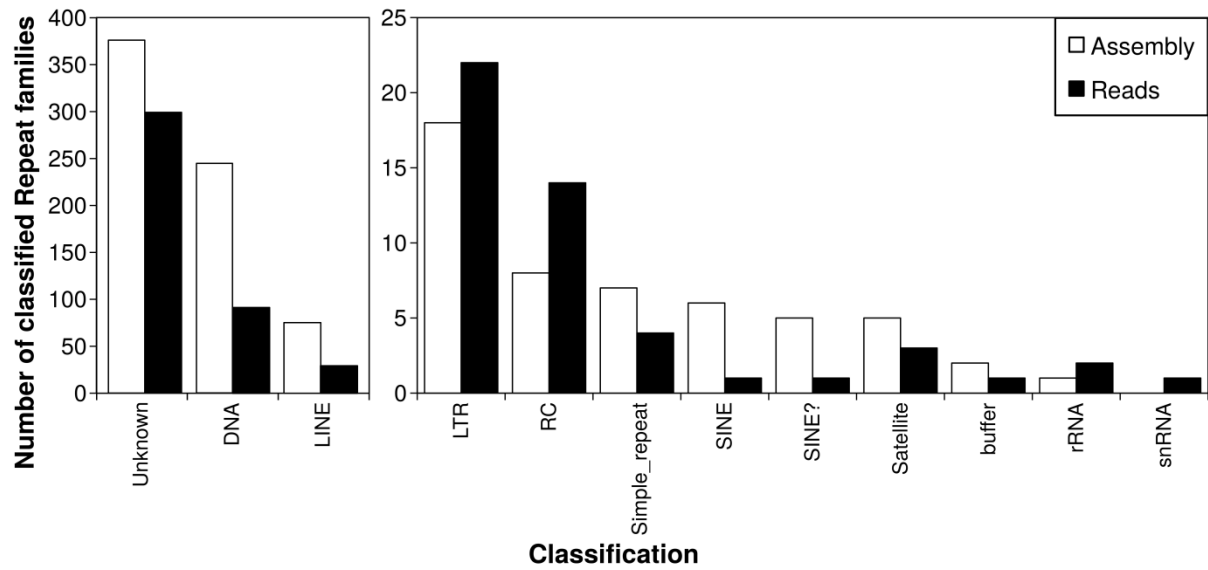
82 Re-mapping of all reads to the repeat library (same procedure as re-mapping on the genome assembly)  
83 revealed 244,177,765 reads (23.58 Gb) to be lying within repetitive sequences. A size estimation on  
84 the peak coverage results in 327.54Mb (20.5% of the 1.6 Gb estimated genome size) which is  
85 comparable to the length of annotated repeats in the assembly by the MAKER2 pipeline respectively  
86 RepeatMasker.

$$\frac{\text{Total nucleotides in repeats}}{\text{peak coverage}} = \frac{23.58\text{Gb}}{72} = 327.54\text{Mb}$$

367.14Mb annotated repeats + 58.45Mb gaps = 425.59Mb repeats in assembly

425.59Mb repeats in assembly + 692.74Mb missing = 1118.33Mb (69.9% possible repeats)

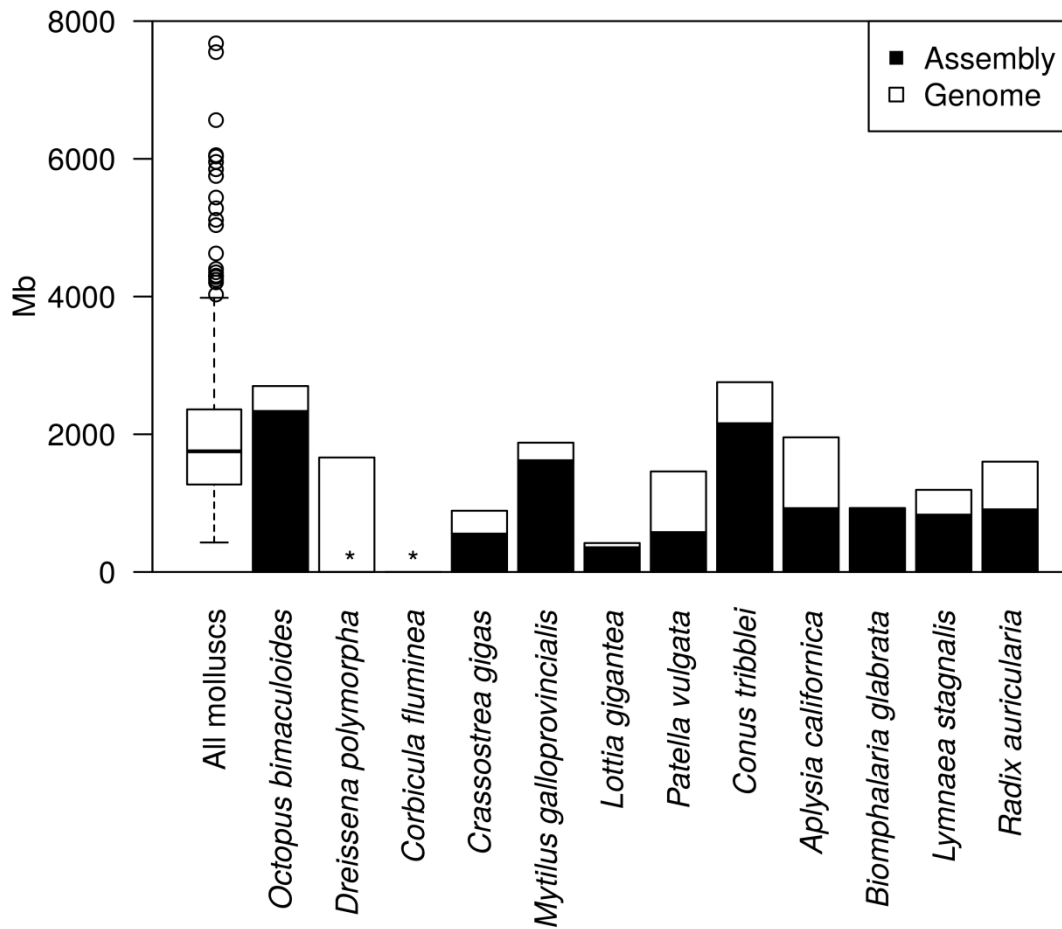
87



88

89 **Supplementary Figure 3.** Classified repeat families. RepeatModeler was executed on the  
 90 assembly and on the contigs originating from randomly drawing reads at certain coverage and their  
 91 assembly. Please note the different scales.

92



93

94 **Supplementary Figure 4.** Mollusc genome sizes. The first column shows the distribution of  
 95 genome sizes of all molluscan records found in the Animal Genome Size Database (Gregory 2016;  
 96 accessed on 28<sup>th</sup> of October 2016; N=263; circles are outliers). The total height of bars from mollusc  
 97 species with available genome assemblies shows the estimated genome size and the black filled part  
 98 the fraction represented in the corresponding assembly. \*The total assembly length of *Dreissena* and  
 99 *Corbicula* is below 1 Mb and therefore no black bar is visible in the graph. There is no estimated  
 100 genome size for *Corbicula* and therefore no white bar is displayed. Citations from assemblies and  
 101 genome sizes are given in Table 1.

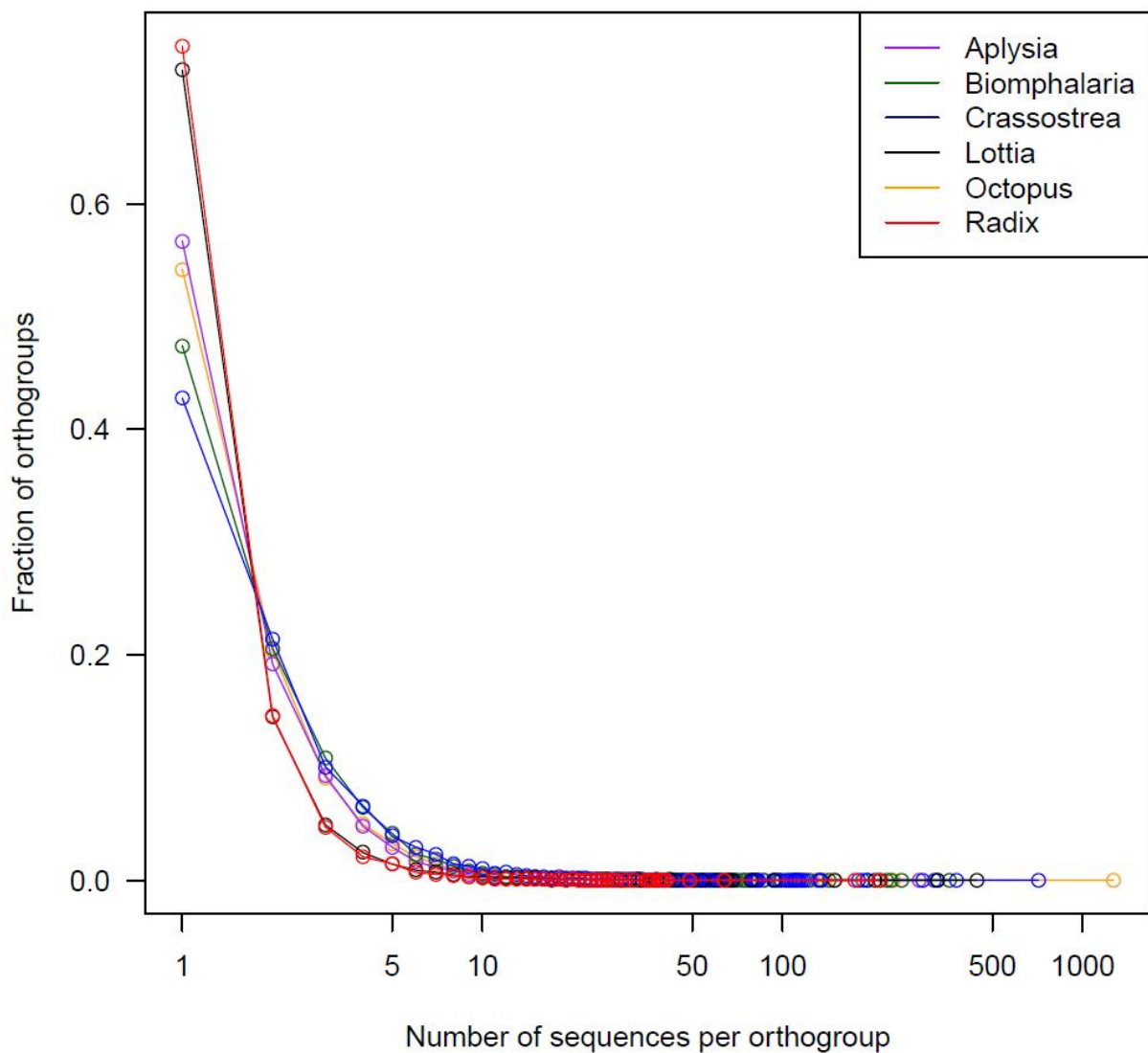
102



**Supplementary Table 4.** Proteins similar to Swiss-Prot entries (accessed May 11<sup>th</sup> 2016; blastp e-value < 10<sup>-10</sup>).

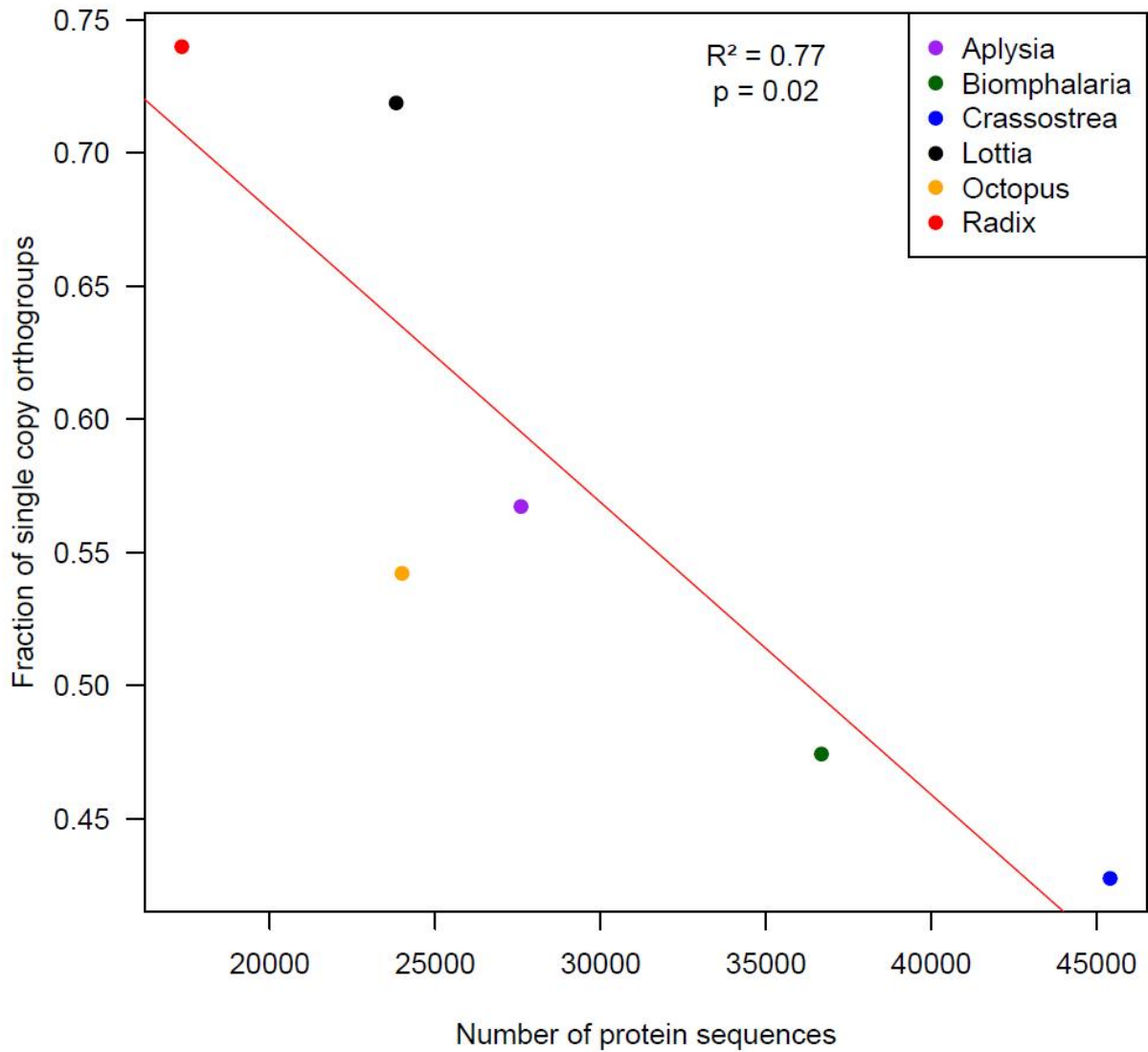
Species	Annotated proteins	Proteins similar to Swiss-Prot	Percentage
<i>Octopus bimaculoides</i>	23,994	19,217	80.1
<i>Aplysia californica</i>	27,591	19,920	72.2
<b><i>Radix auricularia</i></b>	<b>17,338</b>	<b>12,207</b>	<b>70.4</b>
<i>Crassostrea gigas</i>	45,406	30,917	68.1
<i>Biomphalaria glabrata</i>	36,675	24,422	66.6
<i>Lottia gigantean</i>	23,822	13,950	58.6

103



104

105 **Supplementary Figure 5.** Number of sequences from one species in orthogroups. The fraction of  
 106 orthogroups (see Material and Methods) is the number of orthogroups containing this number of  
 107 different protein sequences of one species relative to the number of orthogroups containing this  
 108 species. Note the logarithmic x-axis.



110

111 **Supplementary Figure 6.** Correlation of the number of annotated protein sequences per species

112 and the number of orthogroups containing only one sequence of this species. The fraction of

113 orthogroups is relative to the number of orthogroups containing this species.

114

**Supplementary Table 5.** Protein sets used for ortholog clustering and GO-term enrichment. Protein sets from six different mollusc species and 16 different spiralian species from outside the Mollusca, which were used for orthologous clustering and Gene Ontology enrichment analysis. (# annotated OGs = number of ortholog groups, which were functionally annotated by GO-terms using a sequence of the according species)

Species	NCBI Genome ID	# sequences	# sequences in orthogroups	# unassigned sequences	# orthogroups containing species	# species specific orthogroups	# annotated sequences from InterProScan	# annotated OGs from species in reference list
<i>Radix auricularia</i>		17,338	15,880	1,458	9,128	5	10,436	6,105
<i>Biomphalaria glabrata</i>	357	36,675	33,044	3,631	10,594	43	21,813	1,977
<i>Aplysia californica</i>	443	27,591	25,463	2,128	10,208	44	17,495	863
<i>Lottia gigantea</i>	15113	23,822	20,444	3,378	9,696	56	12,387	886
<i>Crassostrea gigas</i>	10758	45,406	38,769	6,637	9,710	115	28,107	1,576
<i>Octopus bimaculoides</i>	41501	23,994	22,211	1,783	8,139	72	16,396	500
<i>Caenorhabditis elegans</i>	41	28,137	25,780	2,357	11,481	50	14,052	2,237
<i>Caenorhabditis remanei</i>	253	31,476	27,504	3,972	12,515	57	13,795	1,015
<i>Caenorhabditis brenneri</i>	254	30,670	25,398	5,272	11,387	94	13,154	646
<i>Caenorhabditis briggsae</i>	40	21,959	18,977	2,982	11,674	45	9,446	322
<i>Wuchereria bancrofti</i>	2616	19,323	15,942	3,381	8,641	14	7,768	683
<i>Pristionchus pacificus</i>	246	16,763	11,488	5,275	6,022	65	7,054	659
<i>Loa loa</i>	2686	16,281	12,623	3,658	8,439	14	6,646	238
<i>Clonorchis sinensis</i>	2651	13,634	12,479	1,155	6,750	2	6,590	808
<i>Schistosoma mansoni</i>	236	11,713	10,603	1,110	6,677	8	6,175	256
<i>Opisthorchis viverrini</i>	32471	16,356	13,425	2,931	6,475	6	6,120	129
<i>Hymenolepis microstoma</i>	24432	12,371	9,822	2,549	5,906	35	5,839	402
<i>Schistosoma haematobium</i>	10705	11,140	10,496	644	6,686	13	5,689	129
<i>Echinococcus multilocularis</i>	22333	10,656	9,368	1,288	6,559	1	5,619	162
<i>Echinococcus granulosus</i>	10706	11,319	8,611	2,708	6,094	5	5,274	110
<i>Capitella teleta</i>	15118	31,978	25,701	6,277	9,353	107	0	0
<i>Helobdella robusta</i>	15112	23,426	17,475	5,951	6,908	53	0	0

117 **Supplementary Note 6.** Ortholog clustering and Gene Ontology enrichment  
118 Protein sets from six (including *Radix*) different mollusc species and 16 different spiralian species  
119 from outside the Mollusca (Supplementary Table 4) were used to predict ortholog cluster by  
120 OrthoFinder 0.7.1 (Emms & Kelly 2015) with default parameters. Functional annotation of all protein  
121 sets with Gene Ontology (GO) terms was performed with InterProScan 5 (Quevillon et al. 2005;  
122 Zdobnov & Apweiler 2001), using default parameters.

123 To construct the reference list for the GO enrichment analysis a step wise procedure was implemented  
124 to obtain GO annotations for as many orthogroups as possible. As *R. auricularia* is our focal species,  
125 all GO-annotations obtained for *R. auricularia* were assigned to orthogroups. In the following steps  
126 GO-annotations from one species after another were added in a phylogenetic order until all  
127 orthogroups had assigned GO-terms, where possible (see order of species from top to bottom and the  
128 corresponding number of orthogroups annotated with each corresponding species in Supplementary  
129 Table 5). The GO-term enrichment analysis was based on this reference list of annotated orthogroups  
130 compared to a test set containing orthogroups of certain sets of species. The analysis were performed  
131 using TopGO (Alexa & Rahnenfuhrer 2016).

132

**Supplementary Table 6.** Significantly enriched GO-terms for *Radix* specific and unassigned orthogroups (FDR < 5%).

Term	Annotated	Significant	Expected	q-value
Nucleoside transmembrane transport	15	6	0.29	2.10E-07
Glycolytic process	41	6	0.79	0.00012
Microtubule-based movement	138	10	2.66	0.00033
Carbohydrate metabolic process	399	18	7.7	0.00172
Tubulin complex assembly	4	2	0.08	0.00216
Post-chaperonin tubulin folding pathway	4	2	0.08	0.00216
Cilium or flagellum-dependent cell motility	15	3	0.29	0.00271
Neurotransmitter transport	59	5	1.14	0.00546
transport	1,617	42	31.19	0.00594
Sulfate transport	13	2	0.25	0.0251
cGMP biosynthetic process	14	2	0.27	0.02891
Peptidyl-glutamic acid carboxylation	2	1	0.04	0.0382
Inositol trisphosphate metabolic process	2	1	0.04	0.0382
Clathrin coat assembly	2	1	0.04	0.0382
glutamine biosynthetic process	2	1	0.04	0.0382
Clathrin-mediated endocytosis	2	1	0.04	0.0382
Chitin metabolic process	134	6	2.58	0.045

133

**Supplementary Table 7.** Significantly enriched GO-terms in orthogroups containing all five mollusc species but *Radix* (FDR < 5%).

Term	Annotated	Significant	Expected	q-value
Meiotic prophase I	1	1	0	0.0013
Synaptonemal complex assembly	2	1	0	0.0027
Actin ubiquitination	3	1	0	0.004
Carbohydrate transport	9	1	0.01	0.0119
G-protein coupled receptor signaling pathway	457	3	0.61	0.0206

134

135 **Supplementary Note 7.** Preprocessing and trimming

136 All read files were quality checked using FastQC 0.10.1 (Andrews 2010). Reads from small insert  
 137 libraries were preprocessed with Trimmomatic 0.33 (Bolger et al. 2014) using the adapter trimming  
 138 along with a custom adapter file (ILLUMINACLIP:<adapter.fasta>:2:30:10). FastQC reports from

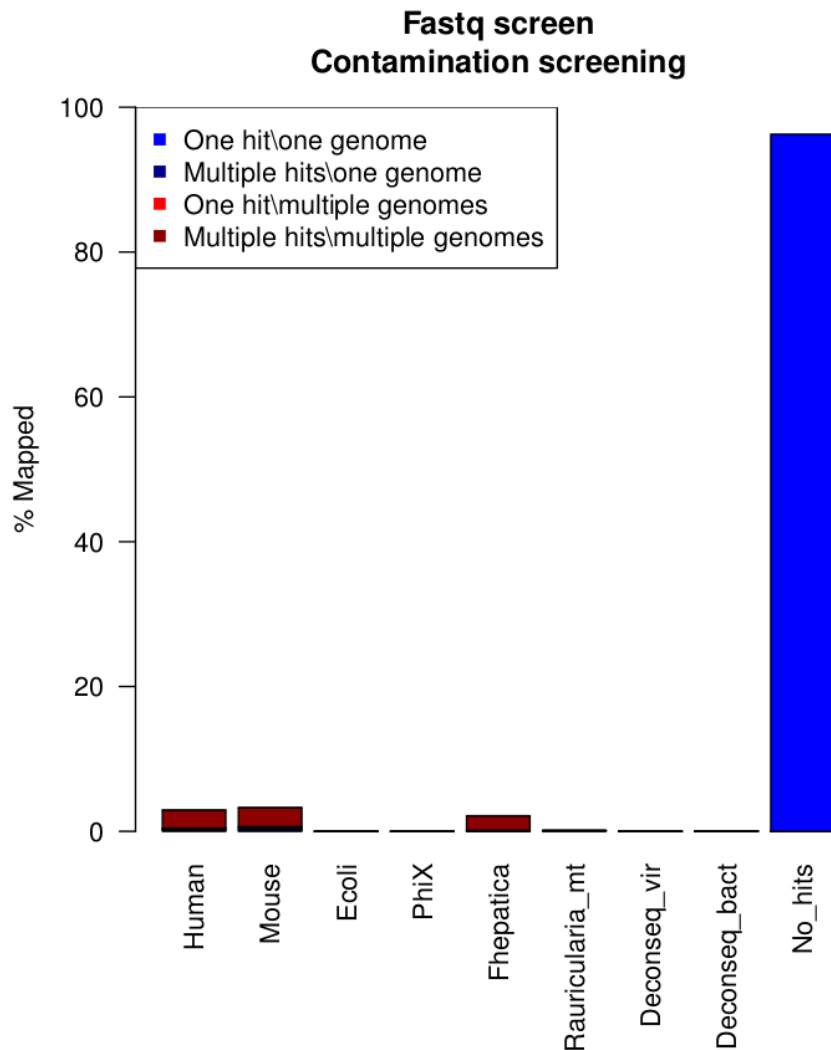
139 mate pair libraries showed overrepresented k-mers at 3' ends of reads that could be assembled into  
140 fragments of external adapters which is only possible if selection for fragments containing junction  
141 adapter was not completely free from other, smaller fragments. Mate pair reads with (partial) external  
142 adapter sequence do indeed overlap. Therefore, motifs from external adapter contamination only were  
143 used during adapter trimming additional to a minimum length threshold of 100 bp to discard false  
144 mate pair reads (ILLUMINACLIP:<motifs.fasta>:2:7:7 MINLEN:100).

145 99.4 % (994,535,287) of all reads and 99.3 % (115,378,553,461 bp) of all nucleotides survived  
146 trimming.

147

#### 148 **Supplementary Note 8.** Contamination screening

149 All adapter trimmed reads were checked for possible contamination using FastqScreen 0.5.2 (Andrews  
150 et al. 2015) with libraries from human (*Homo sapiens* GRCh38), mouse (*Mus musculus* GRCm38), *E.*  
151 *coli* (U00096.3), Enterobacteria phage phiX174 (NC\_001422.1), the known *Radix* parasite *Fasciola*  
152 *hepatica* (GCA\_000947175.1), *R. auricularia* mt genome as positive control (NC\_026538.1) and  
153 simulated bacterial and viral Databases from DeconSeq (Schmieder & Edwards 2011). Most  
154 sequences could not be mapped to the provided libraries (Supplementary Figure 5). Small fractions of  
155 reads could be mapped multiple times to multiple genomes which hint towards similar repeats in  
156 provided libraries and *Radix*.



157

158 **Supplementary Figure 7.** Results of contamination screening. All trimmed reads were mapped  
 159 against different libraries containing possible contamination sources using FastqScreen.

160

161 **Supplementary Note 9.** Material and methods of flow cytometric analysis

162 Genome size (2C-values; Greilhuber et al. 2005) was estimated by flow cytometry using fresh foot  
 163 tissue of six individuals of the same *Radix auricularia* inbred line as used for sequencing  
 164 (Supplementary Note S3) and the Partec CyFlow Space (Partec, Münster, Germany) equipped with a  
 165 green solid-state laser (Partec, 532 nm, 30 mW). Sample preparation followed two-step Otto protocol  
 166 (Otto 1990), with an internal standard *Glycine max* (L.) Merr. cv. Polanka (2C = 2.50 pg; Doležel et al.  
 167 1994) and *Lycopersicon esculentum* Mill. cv. Stupické polní tyčkové rané (2C = 1.96 pg; Doležel et al.  
 168 1992). The tissue of *Radix auricularia* and the internal reference standard were mixed and chopped  
 169 with a razor blade in a Petri dish containing 1 ml of ice-cold Otto I buffer (0.1 M citric acid, 0.5%

170 Tween 20). The suspension was filtered through a 42- $\mu$ m nylon mesh and incubated for approximately  
171 15 min at room temperature. The staining solution consisted of 1 ml of Otto II buffer (0.4 M  
172  $\text{Na}_2\text{HPO}_4 \cdot 12 \text{H}_2\text{O}$ ),  $\beta$ -mercaptoethanol (final concentration of 2  $\mu\text{l/ml}$ ), intercalating fluorochrome  
173 propidium iodide (PI) and RNase IIA (both at final concentrations of 50  $\mu\text{g/ml}$ ). Fluorescence  
174 intensities of 5000 particles (nuclei) were recorded. Sample/standard ratios were calculated from the  
175 means of the sample and standard fluorescence histograms, and only histograms with coefficients of  
176 variation (CVs) < 5% for the G0/G1 sample peak were considered. Four to five replicates were  
177 measured on three different days in order to minimize potential random instrumental error. If the  
178 between-day variation in fluorescence intensity was above 4% the most extreme value was discarded.  
179 The 2C-values estimated using different standards were inferred separately.

180

#### 181 **Supplementary Note 10.** Transcriptome assemblies

182 In order to generate de novo transcriptomes (contigs) of available RNA-seq data, raw Roche 454 reads  
183 from *R. auricularia* and MOTU3 (Feldmeyer et al. 2015; accession numbers SRR1926149 and  
184 SRR1926204) were analysed in FastQC and trimmed with Trimmomatic 0.33 (“HEADCROP:35  
185 TRAILING:20 MINLEN:200” and “ILLUMINACLIP:<overrepresented\_k-mers.fa>:2:1:1:1  
186 TRAILING:20 SLIDINGWINDOW:10:20 HEADCROP:35 CROP:540 MINLEN:50” respectively).  
187 The *R. auricularia* and MOTU3 transcriptomes were assembled using the Overlap Layout Consensus  
188 algorithms of MIRA 4.0.2 (Chevreux et al. 1999) (job = est,denovo,accurate; parameters =  
189 454\_SETTINGS -ALIGN:min\_relative\_score=75; technology = 454). Contiguous sequences were  
190 constructed from the above mentioned contigs and contigs from Feldmeyer et al. (2015) via meta-  
191 assembly using MIRA (job = genome,denovo,accurate; parameters = TEXT\_SETTINGS --noclipping  
192 -AS:epoq=no -AS:mrpc=1 -OUT:sssip=yes; technology = text). Meta-contigs and unassembled  
193 contigs from the meta-assembly were merged to obtain the final transcriptomes.

194 Raw Illumina sequences for MOTU5 were obtained from Feldmeyer et al (2015) (SRR1926203),  
195 and *R. balthica* (MOTU2) from (Feldmeyer et al. 2011; Tills et al. 2015) and so far unpublished 50bp  
196 Illumina reads. All embryonic *R. balthica* reads are deposit in NCBI’s BioProject PRJEB9533.  
197 Analyses in FastQC provided the overrepresented k-mers. Trimming using Trimmomatic 0.33



198 (“ILLUMINACLIP:<adapter.fasta>:2:30:10 ILLUMINACLIP:<overrepresented.fasta>:2:1:1:1  
199 SLIDINGWINDOW:5:20 HEADCROP:15 CROP:50 MINLEN:50”;  
200 “ILLUMINACLIP:<adapter.fasta>:2:30:10 ILLUMINACLIP:<overrepresented\_k.mers.fa>:0:1:1:1  
201 SLIDINGWINDOW:5:30 MINLEN:50” and “ILLUMINACLIP<adapter.fasta>:2:30:10  
202 HEADCROP:10 MINLEN:35” respectively) survived 75.1 % for MOTU5 and 84.1 % for *R. balthica*  
203 respectively. The two datasets from MOTU5 and *R. balthica* were assembled unpaired with Trinity  
204 2.0.6 (Grabherr et al. 2011) and Bridger 2014-12-01 (Chang et al. 2015) each. The Trinity and Bridger  
205 contigs from one species were meta-assembled in MIRA with the same parameters as the meta  
206 assembly from *R. auricularia* (MOTU4) and MOTU3.

207

#### 208 **Supplementary Note 11.** Genome assembly

209 All trimmed reads were assembled using the Platanus 1.2.1 pipeline (Kajitani et al. 2014). The  
210 assembly was performed with default parameters except the initial k-mer size of 63 and a stepsize of 2.  
211 The automatically detected maximum k-mer size was 88, implying that the last stepsize was 1 (from  
212 87 to 88). Scaffolding within the Platanus pipeline was computed with default parameters (e.g.  
213 minimum 3 links) using all six libraries in ascending order along with their average insert sizes  
214 according library preparation. Gapclose as last step of the Platanus pipeline was performed with  
215 standard parameters using all six libraries in ascending order regarding insert size. All sequences  
216 smaller 500 bp were excluded from following assembly steps.

217 The filtered and gapclosed Platanus scaffolds were scaffolded again using SSPACE Standard  
218 3.0 (Boetzer et al. 2011). Since SSPACE/our hardware was not able to run the scaffolding successfully  
219 until the end, separately all six libraries were mapped unpaired using bowtie2 2.2.5 (Langmead &  
220 Salzberg 2012) against the filtered and gapclosed Platanus scaffolds. The Mappings were sorted by  
221 read name with samtools 1.1 (Li et al. 2009) and converted to the SSPACE readable TAB-format  
222 using the script sam\_bam2tab.pl provided by SSPACE. All trimmed reads from all six libraries were  
223 used during the scaffolding via SSPACE using default parameters apart from contig extension  
224 switched on. The insert sizes and corresponding errors from paired reads were calculated from the  
225 re-mappings to the filtered and gapclosed Platanus scaffolds (250 bp: 236/0.17; 500 bp: 488/0.08;

226 800 bp: 771/0.06; 2 kb: 2147/0.12; 5 kb 4909/0.09; 10 kb 10010/0.10; “Library”: “mean insert  
227 size”/”error” respectively). Unpaired reads for contig extension only were aligned by SSPACE using  
228 bowtie.

229 The SSPACE scaffolds were scaffolded again with L\_RNA\_scaffolder (Xue et al. 2013) using  
230 the transcriptomes of four *Radix* species (Supplementary Note 10). Therefore raw transcriptomic data  
231 from *R. auricularia* (MOTU4), *R. balthica* (MOTU2), MOTU3 and MOTU5 was preprocessed and  
232 assembled to obtain ESTs. Afterwards the ESTs were mapped with BLAT 35 (Kent 2002) and filtered  
233 for order and orientation on one scaffold using a custom perl script. Only ESTs that survived filtering  
234 were used as input in BLAT again to obtain the psl-file which is used as input for L\_RNA\_scaffolder.  
235 All four transcriptomic meta assemblies were mapped with BLAT 35 against the SSPACE scaffolds  
236 using standard parameters except *-extendThroughN* and *-out=blast8*. Only alignments with correct  
237 order and orientation of the parts of the split alignment on one scaffold, covering at least 20 bp and  
238 80 % of the transcriptomic contig length were kept (Supplementary Figure 8).

239 Before transcriptomic scaffolding the mitochondrial scaffold was excluded to avoid  
240 misscaffolding caused by Mitochondrial DNA-like sequences in the nucleus (NUMTs). Afterwards the  
241 filtered transcriptomic contigs were used to sequential scaffold the mitochondrial-free SSPACE  
242 scaffolds in three steps using L\_RNA\_scaffolder by first mapping the filtered transcriptomic contigs  
243 again with BLAT 35 (*-extendThroughN -noHead*) to create the correct input for L\_RNA\_scaffolder.  
244 All scaffoldings with L\_RNA\_scaffolder were performed with default parameters. First the  
245 mitochondrial-free SSPACE scaffolds were scaffolded with the filtered transcriptomic contigs from  
246 the same species *R. auricularia* (MOTU4). Second the output from the first scaffolding was used as  
247 input along with the filtered transcriptomic contigs from *Radix sp.* MOTU5. Third the output from the  
248 second scaffolding was used as input along with the filtered transcriptomic contigs from *R. balthica*  
249 (MOTU2) and *Radix sp.* MOTU3 together since hybridization is observed between the two MOTUs  
250 (Patel et al. 2015). All results from transcript filtering and scaffolding are summarized in  
251 Supplementary Table 8.

252 The excluded mitochondrial scaffold was cut at the same site as by Feldmeyer et al. (2015)  
253 and scaffolded without the shortest (250 bp) insert library in SSPACE (using the same parameters and

254 pipeline as before). Omitting the library with the smallest insert size ensures to keep the cutting  
255 position. Afterwards a gapfilling with GapFiller 1-10 (Boetzer et al. 2012) was performed using all  
256 genomic trimmed paired Illumina reads along with the insert sizes and errors used during genomic  
257 scaffolding with SSPACE except from default parameters. Mapping within GapFiller was executed  
258 via BWA for all libraries and closed all eight gaps and 750 missing nucleotide positions within one  
259 iteration.

260         After scaffolding with L\_RNA\_scaffolder all scaffolds smaller than 1 kb where excluded  
261 (5,140 scaffolds [51.58 %] / 3,335,954 bp total length [0.37 %]) and the mitochondrial scaffold was  
262 added.

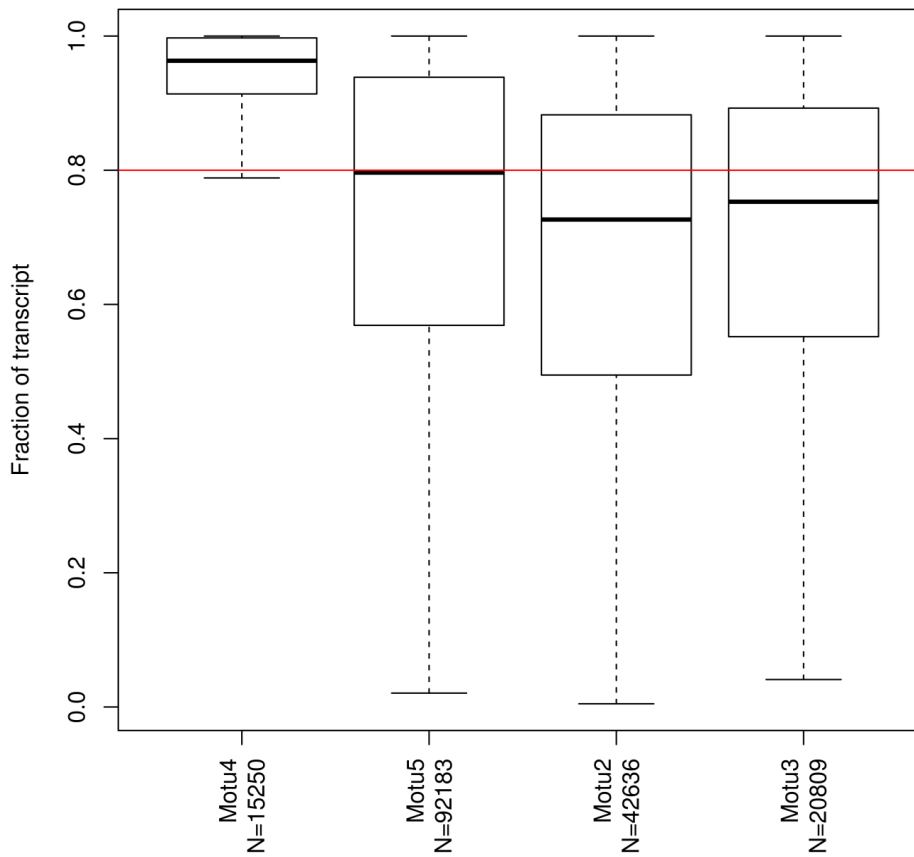
263         All remaining scaffolds were used as input to GapFiller together with all trimmed paired reads  
264 from all libraries and the parameters used before during gapfilling of the mitochondrial scaffold.

265         Because of gap resizing ten scaffolds smaller than 1 kb remained in the assembly after gapfilling.  
266 Three of them still contained N's and were removed from the assembly.

267         A summary of the results of all assembly steps is in Supplementary Table 2 together with  
268 corresponding intermediate results.

269

### Relative representation of correct transcripts in SSPACE scaffolds



270

271 **Supplementary Figure 8.** Transcriptome filtering. Fraction of mapped positions of each contig  
272 surviving order and orientation filtering. The red line marks the cut-off for contigs used in scaffolding  
273 with L\_RNA\_scaffolder at 80 %.

274

**Supplementary Table 8.** Genome scaffolding with transcriptomic data. Four *Radix* species (MOTUs) were used. The first line of each cell shows the number of sequences first and the total length in Mb second. In the second line the fraction of sequences and the fraction of total length relative to the unfiltered/raw meta transcriptome is shown respectively.

	MOTU4	MOTU5	MOTU2	MOTU3	Total
Raw meta transcriptome	22,798 / 16.3	145,687 / 75.2	91,728 / 101.4	34,418 / 27.9	294,631 / 220.8
Mapped with BLAT	22,108 / 16.0 97.0 / 98.3	116,566 / 66.6 80.0 / 88.5	74,411 / 95.6 81.1 / 94.3	30,686 / 26.4 92.9 / 94.3	243,771 / 204.5 82.7 / 92.6
Order / orientation	15,250 / 10.6 66.9 / 65.3	92,183 / 50.2 63.3 / 66.7	42,636 / 45.3 46.5 / 44.7	20,809 / 17.2 60.5 / 61.3	170,878 / 123.3 58.0 / 55.8
Relative representation $\geq 80\%$	14,687 / 10.4 64.4 / 63.6	45,672 / 31.2 31.3 / 41.5	16,584 / 22.2 18.1 / 21.9	8,864 / 7.6 26.8 / 27.1	85807 / 71.4 29.1 / 32.3
Scaffolded sites	49	232	71		352

275

276 **Supplementary Note 12.** Repeat library

277 The repeat library was created using dnaPipeTE 1.2 (Goubert et al. 2015) and RepeatModeler 1.0.4  
278 (Simit & Hubley 2015).

279 *dnaPipeTE*

280 As suggested from the first author, mitochondrial and reverse reads were excluded from the trimmed  
281 reads as input for dnaPipeTE, thereby producing an input file containing only forward and unpaired  
282 reads from all libraries (499,530,440 reads / 57,943,393,315 bp). For 30 coverages from 0.0001 to 0.7  
283 at an estimated genome size of 1.3 Gb two samples each were drawn. All parameters were set as  
284 default except the minimal contig length was lowered from 200 to 50 bp. After analysing the N50  
285 distribution (Supplementary Figure 9) 0.025 was determined as the optimal sampling coverage and 50  
286 repetitions were executed for this coverage with same parameters as above, resulting in 734,889  
287 contigs with total length of 110,811,062 bp. The maximized N50 at 0.0075x coverage was not viewed  
288 to be the optimal sampling coverage because the samples are very small and a huge variation  
289 regarding N50 is expected, furthermore filter steps are induced afterwards to discard non-repeats.

290 *Coverage of dnaPipeTE contigs*

291 All trimmed reads were mapped unpaired against all contigs from the 50 repetitions using BWA mem  
292 with the options -t 80 -k 25 -a -y 26 -c 1000000000 apart from default settings. Using samtools 1.3, the

293 coverage per position was calculated running *mpileup* with options -A -C 50 -d 1000000 apart from  
294 standard parameters. Contigs with a median coverage smaller than the 90 % quantile (94x) of the per  
295 position coverage distribution from the re-mapping of all genomic reads used in the assembly to the  
296 final genome assembly were filtered out (669,284 contigs / 78,314,604 bp). The remaining 65,605  
297 contigs (32,496,458bp total length) with sufficient coverage were used in the next steps.

298 RepeatModeler was then executed with default parameters on final genome assembly and on  
299 the high covered contigs from dnaPipeTE. The resulting fasta files containing the repeat families were  
300 concatenated into 1,216 sequences with a total length of 1,111,662bp.

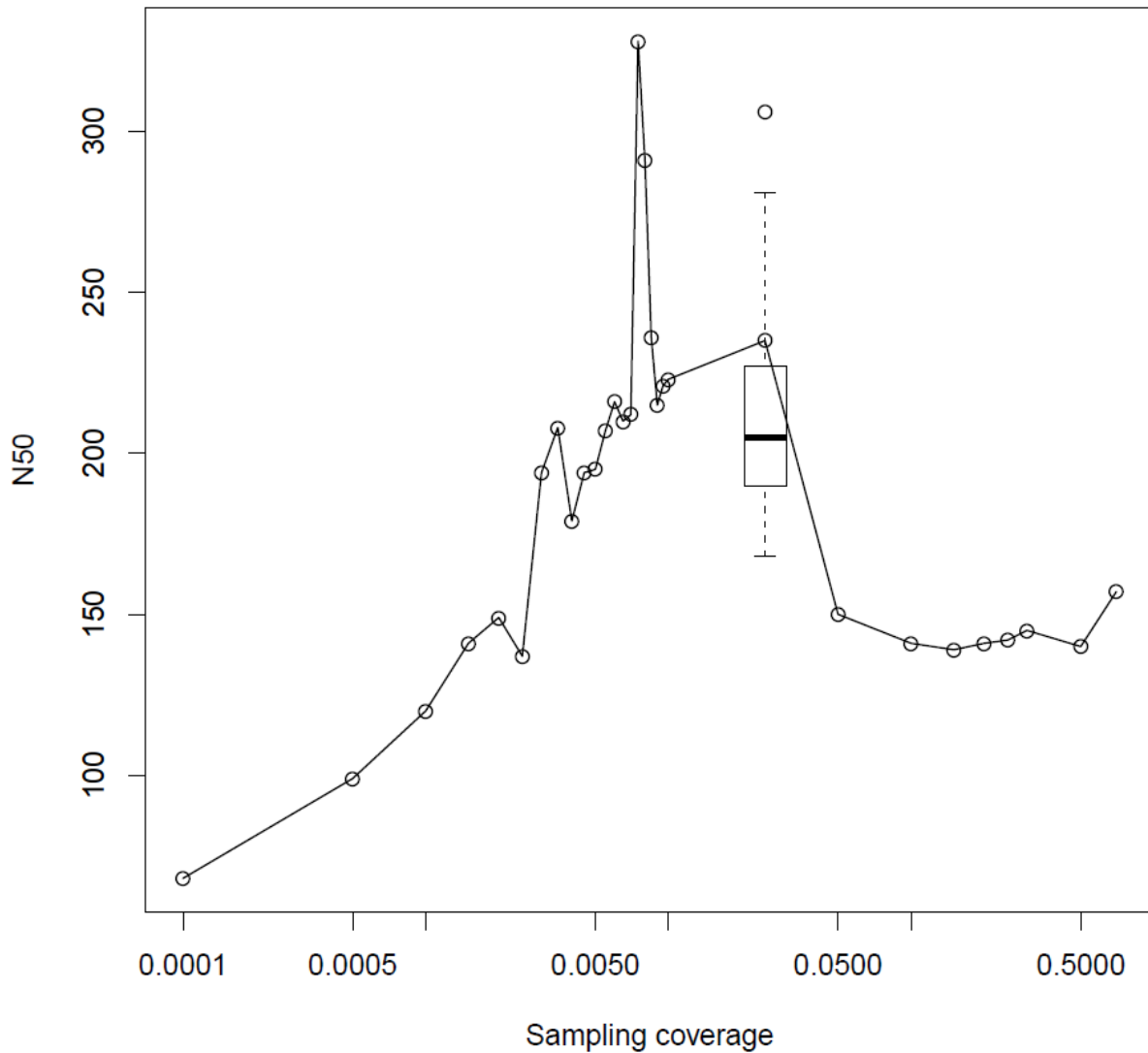
#### 301 *Filtering out possible proteins from the repeat library*

302 First a repeat-protein free protein database was built from Swiss-Prot database (Bateman et al. 2015)  
303 (accessed on May 11<sup>th</sup> 2016) and all repeat sequences from the Repbase database (Bao et al. 2015)  
304 (accessed on May 11<sup>th</sup> 2016). The repeat sequences were searched in the Swiss-Prot database via  
305 BlastX 2.3.0+ (Camacho et al. 2009) with an e-value cutoff of  $10^{-11}$ . Protein sequences with hits from  
306 repeat sequences containing at least 20 bp in one hit were removed to obtain a repeat-protein free  
307 protein database.

308 The concatenated repeat families from RepeatModeler runs on the genome assembly and the  
309 contigs from read subsampling were blasted (BlastX; e-value cutoff  $10^{-11}$ ) against the repeat-protein  
310 free protein database. For six different families a hit was reported. These sequences were removed  
311 from final repeat library containing 1,210 sequences totalling a length of 1,101,332bp and an N50 of  
312 1,492bp.

313

**dnaPipeTE subsampling**  
**Minimum Contig length 50bp**



314 **Supplementary Figure 9.** Read subsampling. N50 distribution of 30 repeat sequence assemblies  
 315 at different sampling coverages is shown by the connected circles. The boxplot shows the N50  
 316 distribution from the 50 repetitions at final sampling coverage at 0.025. Note the logarithmic x-axis.  
 317

**Supplementary Table 9.** Summarized results of the CEGMA analysis.

Statistics of the completeness of the genome based on 248 CEGs					
	#Prots	%Completeness	#Total	Average	%Ortho
Complete	152	61.29	165	1.09	8.55
Group 1	36	54.55	40	1.11	11.11
Group 2	36	64.29	40	1.11	11.11
Group 3	33	54.10	36	1.09	9.09
Group 4	47	72.31	49	1.04	4.26
Partial	233	93.95	328	1.41	32.62
Group 1	58	87.88	72	1.24	22.41
Group 2	54	96.43	82	1.52	37.04
Group 3	58	95.08	84	1.45	41.38
Group 4	63	96.92	90	1.43	30.16

These results are based on the set of genes selected by Genis Parra

Key:

Prots = number of 248 ultra-conserved CEGs present in genome

%Completeness = percentage of 248 ultra-conserved CEGs present

Total = total number of CEGs present including putative orthologs

Average = average number of orthologs per CEG

%Ortho = percentage of detected CEGS that have more than 1 ortholog

318

319 **Supplementary Note 13.**      Annotation

320 For annotation of the assembly we used the MAKER2 2.31.8 pipeline (Cantarel et al. 2008; Holt &  
321 Yandell 2011) combined with MPICH2 (<http://www.mpich.org/>) in three iterations combined with  
322 retraining of the species model in between.

323         Firstly an Augustus species model was computed on the Augustus webserver (Stanke et al.  
324 2004; <http://bioinf.uni-greifswald.de/webaugustus/training/create>). The assembly, the species own  
325 ESTs and annotations in gff format from BUSCO 1.2 (Simão et al. 2015) were used as input. BUSCO  
326 was run on the assembly using the metazoan dataset together with the option --long apart from  
327 standard parameters. From a CEGMA 2.5 (Parra et al. 2007) run on the assembly a SNAP 2006-07-28  
328 (Korf 2004) model was built using the script cegma2zff from the MAKER2 distribution and the SNAP  
329 scripts fathom (fathom genome.ann genome.dna -categorize 1000 && fathom -export 1000 -plus  
330 uni.ann uni.dna), forge (export.ann export.dna) and hmm-assembler.pl. A Genemark (GeneMark-ES  
331 suite 4.32; Lomsadze et al. 2005) model was built from a self-training (--ES) on the assembly.

332 The assembly, the Augustus species model, the ESTs, the ESTs from the other *Radix* species as  
333 alternative ESTs, the complete Swiss-Prot database (Accessed May 23<sup>rd</sup> 2016), the custom repeat  
334 library and the HMM models from SNAP and GeneMark were used as input for the first MAKER  
335 iteration. The options est2genome and protein2genome were switched off. Furthermore the minimum  
336 protein length being reported was set to 10 amino acids.

337         After the first iteration the gff file for the whole assembly was extracted using the MAKER  
338 gff3\_merge, converted with maker2zff and a new HMM model was built for SNAP the same way as  
339 above. The Augustus species model was retrained locally by first converting with the SNAP script  
340 zff2gff3.pl (zff2gff3.pl genome.ann | perl -plne 's/\t(\S+)\\$/\t.\t\$1/') and second the autoAug.pl script  
341 from Augustus 3.2.2 (Stanke et al. 2006). The input for the autoAug.pl was the draft genome



342 assembly, the trained Augustus species model, the ESTs and the gff3 file created from the first  
343 MAKER iteration. Apart from standard parameters -v --useexisting were used.

344 For the second MAKER iteration the SNAP HMM model from CEGMA was exchanged to that  
345 created from the output of the first iteration as well as the updated Augustus species model were used  
346 as input. The minimum protein length was raised to 30 amino acids. Afterwards a retraining as above  
347 and a third MAKER iteration was realized.

348         The annotation pipeline resulted in 17,338 protein coding genes from lengths between 31 AA  
349 and 9,660 AA with a median of 332 AA, containing 7,968,643 AA in total. Gene lengths are between  
350 141 and 127,541 bp with a mean of 11,570 bp and sum up to a total of 201 Mb which corresponds to  
351 12.5% of the estimated genome size of 1.6 Gb. These genes contain 147,195 exons with an average of  
352 8.5 exons per gene. Exonlengths reach from 3 bp to 10,740 bp with a mean of 171.9 bp and sum up to  
353 a total of 25 Mb which corresponds to 1.6% of estimated genome size of 1.6 Gb.

354

355 **Literature cited**

- 356 Alexa A, Rahnenfuhrer J. 2016. topGO: Enrichment Analysis for Gene Ontology. R package version  
 357 2.26.0.
- 358 Andrews S. 2010. FastQC: a quality control tool for high throughput sequence data.  
 359 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>.
- 360 Andrews S, Fiers M, Wingett S. 2015. FastQ Screen.  
 361 [http://www.bioinformatics.babraham.ac.uk/projects/fastq\\_screen](http://www.bioinformatics.babraham.ac.uk/projects/fastq_screen).
- 362 Bao W, Kojima KK, Kohany O. 2015. Repbase Update, a database of repetitive elements in eukaryotic  
 363 genomes. *Mob. DNA*. 6:11. doi: 10.1186/s13100-015-0041-9.
- 364 Bateman A et al. 2015. UniProt: A hub for protein information. *Nucleic Acids Res*. 43:D204–D212.  
 365 doi: 10.1093/nar/gku989.
- 366 Boetzer M et al. 2012. Toward almost closed genomes with GapFiller. *Genome Biol*. 13:R56. doi:  
 367 10.1186/gb-2012-13-6-r56.
- 368 Camacho C et al. 2009. BLAST plus: architecture and applications. *BMC Bioinformatics*. 10:1. doi:  
 369 Artn 421\nDoi 10.1186/1471-2105-10-421.
- 370 Cantarel BL et al. 2008. MAKER: An easy-to-use annotation pipeline designed for emerging model  
 371 organism genomes. *Genome Res*. 18:188–196. doi: 10.1101/gr.6743907.
- 372 Chang Z et al. 2015. Bridger: a new framework for de novo transcriptome assembly using RNA-seq  
 373 data. *Genome Biol*. 16:1–10. doi: 10.1186/s13059-015-0596-2.
- 374 Chevreux B, Wetter T, Suhai S. 1999. Genome Sequence Assembly Using Trace Signals and  
 375 Additional Sequence Information. *Comput. Sci. Biol. Proc. Ger. Conf. Bioinforma. '99, GCB,*  
 376 *Hann. Ger.* 45–56. doi: 10.1.1.23/7465.
- 377 Doležel J, Doleželová M, Novák FJ. 1994. Flow cytometric estimation of nuclear DNA amount in  
 378 diploid bananas (*Musa acuminata* and *M. balbisiana*). *Biol. Plant*. 36:351–357. doi:  
 379 10.1007/BF02920930.
- 380 Doležel J, Sgorbati S, Lucretti S. 1992. Comparison of three DNA fluorochromes for flow cytometric  
 381 estimation of nuclear DNA content in plants. *Physiol. Plant*. 85:625–631. doi: 10.1111/j.1399-  
 382 3054.1992.tb04764.x.
- 383 Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome comparisons  
 384 dramatically improves orthogroup inference accuracy. *Genome Biol*. 16:157. doi:  
 385 10.1186/s13059-015-0721-2.
- 386 Feldmeyer B, Greshake B, Funke E, Ebersberger I, Pfenninger M. 2015. Positive selection in  
 387 development and growth rate regulation genes involved in species divergence of the genus  
 388 *Radix*. *BMC Evol. Biol*. 15:164. doi: 10.1186/s12862-015-0434-x.
- 389 Feldmeyer B, Wheat CW, Krezdorn N, Rotter B, Pfenninger M. 2011. Short read Illumina data for the  
 390 de novo assembly of a non-model snail species transcriptome (*Radix balthica*,  
 391 *Basommatophora*, *Pulmonata*), and a comparison of assembler performance. *BMC Genomics*.  
 392 12:317. doi: 10.1186/1471-2164-12-317.
- 393 Goubert C et al. 2015. De novo assembly and annotation of the Asian tiger mosquito  
 394 (*Aedes albopictus*) repeatome with dnaPipeTE from raw genomic reads and comparative  
 395 analysis with the yellow fever mosquito (*Aedes aegypti*). *Genome Biol. Evol*. 7:1192–1205.  
 396 doi: 10.1093/gbe/evv050.
- 397 Grabherr MG et al. 2011. Full-length transcriptome assembly from RNA-Seq data without a reference  
 398 genome. *Nat. Biotechnol*. 29:644–52. doi: 10.1038/nbt.1883.
- 399 Gregory TR. 2016. Animal Genome Size Database.
- 400 Greilhuber J, Doležel J, Lysák MA, Bennett MD. 2005. The origin, evolution and proposed  
 401 stabilization of the terms ‘genome size’ and ‘C-value’ to describe nuclear DNA contents. *Ann.*  
 402 *Bot*. 95:255–260. doi: 10.1093/aob/mci019.
- 403 Holt C, Yandell M. 2011. MAKER2 : an annotation pipeline and genome- database management tool  
 404 for second- generation genome projects. *BMC Bioinformatics*. 12:491. doi: 10.1186/1471-2105-  
 405 12-491.
- 406 Kajitani R et al. 2014. Efficient de novo assembly of highly heterozygous genomes from whole-  
 407 genome shotgun short reads. *Genome Res*. 24:1384–1395. doi: 10.1101/gr.170720.113.
- 408 Korf I. 2004. Gene finding in novel genomes. *BMC Bioinformatics*. 5:59. doi: 10.1186/1471-2105-5-  
 409 59.

410 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv  
411 Prepr. doi: arXiv:1303.3997 [q-bio.GN].  
412 Li H et al. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. 25:2078–2079.  
413 doi: 10.1093/bioinformatics/btp352.  
414 Lomsadze A, Ter-Hovhannisyanyan V, Chernoff YO, Borodovsky M. 2005. Gene identification in novel  
415 eukaryotic genomes by self-training algorithm. *Nucleic Acids Res*. 33:6494–6506. doi:  
416 10.1093/nar/gki937.  
417 Okonechnikov K, Conesa A, García-alcalde F. 2015. Qualimap 2: advanced multi-sample quality  
418 control for high- throughput sequencing data. *Bioinformatics*. 1–3.  
419 Otto F. 1990. DAPI Staining of Fixed Cells for High-Resolution Flow Cytometry of Nuclear DNA.  
420 *Methods Cell Biol*. 33:105–110. doi: 10.1016/S0091-679X(08)60516-6.  
421 Parra G, Bradnam K, Korf I. 2007. CEGMA: A pipeline to accurately annotate core genes in  
422 eukaryotic genomes. *Bioinformatics*. 23:1061–1067. doi: 10.1093/bioinformatics/btm071.  
423 Patel S, Schell T, Eifert C, Feldmeyer B, Pfenninger M. 2015. Characterizing a hybrid zone between a  
424 cryptic species pair of freshwater snails. *Mol. Ecol*. 24:643–655. doi: 10.1111/mec.13049.  
425 Quevillon E et al. 2005. InterProScan: Protein domains identifier. *Nucleic Acids Res*. 33:116–120.  
426 doi: 10.1093/nar/gki442.  
427 Schmieder R, Edwards R. 2011. Fast identification and removal of sequence contamination from  
428 genomic and metagenomic datasets. *PLoS One*. 6. doi: 10.1371/journal.pone.0017288.  
429 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. 2015. BUSCO: Assessing  
430 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.  
431 31:3210–3212. doi: 10.1093/bioinformatics/btv351.  
432 Simit A, Hubley R. 2015. RepeatModeler Open-1.0.  
433 Stanke M, Steinkamp R, Waack S, Morgenstern B. 2004. AUGUSTUS: A web server for gene finding  
434 in eukaryotes. *Nucleic Acids Res*. 32:309–312. doi: 10.1093/nar/gkh379.  
435 Stanke M, Tzvetkova A, Morgenstern B. 2006. AUGUSTUS at EGASP: using EST, protein and  
436 genomic alignments for improved gene prediction in the human genome. *Genome Biol*. 7 Suppl  
437 1:S11.1-8. doi: 10.1186/gb-2006-7-s1-s11.  
438 Tills O, Truebano M, Rundle S. 2015. An embryonic transcriptome of the pulmonate snail *Radix*  
439 *balthica*. *Mar. Genomics*. 24:259–260. doi: 10.1016/j.margen.2015.07.014.  
440 Zdobnov EM, Apweiler R. 2001. InterProScan - an integration platform for the signature-recognition  
441 methods in InterPro. *Bioinformatics*. 17:847–848. doi: 10.1093/bioinformatics/17.9.847.  
442