

# An annotated draft genome for *Radix auricularia*

## (Gastropoda, Mollusca)

**Tilman Schell**<sup>1,2</sup> \*, **Barbara Feldmeyer**<sup>2</sup>, **Hanno Schmidt**<sup>2</sup>, **Bastian Greshake**<sup>3</sup>, **Oliver Tills**<sup>4</sup>,  
**Manuela Truebano**<sup>4</sup>, **Simon D. Rundle**<sup>4</sup>, **Juraj Paule**<sup>5</sup>, **Ingo Ebersberger**<sup>3,2</sup>,  
**Markus Pfenninger**<sup>1,2</sup>

<sup>1</sup> *Molecular Ecology Group, Institute for Ecology, Evolution and Diversity, Goethe-University, Frankfurt am Main, Germany*

<sup>2</sup> *Adaptation and Climate, Senckenberg Biodiversity and Climate Research Centre, Frankfurt am Main, Germany*

<sup>3</sup> *Department for Applied Bioinformatics, Institute for Cell Biology and Neuroscience, Goethe-University, Frankfurt am Main, Germany*

<sup>4</sup> *Marine Biology and Ecology Research Centre, Marine Institute, School of Marine Science and Engineering, Plymouth University, Plymouth, United Kingdom*

<sup>5</sup> *Department of Botany and Molecular Evolution, Senckenberg Research Institute, Frankfurt am Main, Germany*

\* Author for Correspondence: Senckenberg Biodiversity and Climate Research Centre,  
 Senckenberganlage 25, 60325 Frankfurt am Main, Germany. Tel.: +49 (0)69 75 42 18 30, E-mail:  
 Tilman.Schell@senckenberg.de

**Data deposition:** BioProject: PRJNA350764, SRA: SRP092167

## Abstract

Molluscs are the second most species rich phylum in the animal kingdom, yet only eleven genomes have been published so far. To complement the scarce genomic resources, we present the draft genome sequence of the pulmonate freshwater snail *Radix auricularia*. Six whole genome shotgun libraries with different layouts were sequenced to an overall read coverage of 73x. The resulting assembly comprises 4,823 scaffolds with a cumulative length of 910 Mb. The assembly contains 94.6 % of a metazoan core gene collection, indicating an almost complete coverage of the coding fraction. The discrepancy of ~690 Mb compared to the estimated genome size of *R. auricularia* (1.6 Gb) results from an extraordinary high repeat content of 70 % mainly comprising DNA transposons. The 17,338 annotated protein coding genes. The presented draft will serve as starting point for further genomic and population genetic research in this scientifically important phylum.

**Key words:** *De novo* assembly, Genome size, Repeats

## Introduction

Gastropods are the only eukaryotic taxon thriving in all ecosystems on earth. They occupy a maximally diverse set of habitats ranging from the deep sea to the highest mountains and from deserts to the Arctic, and have evolved to a range of specific adaptations (Romero et al. 2015, 2016). However, as for molluscs as a whole, whose species richness is second only to the arthropods (Dunn & Ryan 2015) they are strongly underrepresented among publicly available genomes (Figure 1). To date, only eleven mollusc genome sequences of which six originate from gastropods exist with varying qualities concerning contiguity and completeness (Table 1). Adding data to this white spot has the potential to substantially increase knowledge about molluscs in particular and animal genomics in general.

The pulmonate freshwater snail genus *Radix* has a holarctic distribution (Glöer, Meier-Brook 1998; Cordellier et al. 2012) and plays an important role in understanding climate change in freshwater ecosystems (Sommer et al. 2012). The number of European species as well as the precise evolutionary relationships within the genus are controversial, due to weak morphological differentiation and enormous environmental plasticity across species (Pfenninger et al. 2006). Members of the genus are simultaneously hermaphroditic (Yu et al. 2016; Jarne & Delay 1990) and both outcrossing and self-fertilisation occur (Jarne & Delay 1990; Jarne & Charlesworth 1993; Wiehn et al. 2002). The genus *Radix* is studied in many different fields, including parasitology (e.g. (Huňová et al. 2012)), evolutionary development (Tills et al. 2011), developmental plasticity (Rundle et al. 2011), ecotoxicology (e.g. (Hallgren et al. 2012)), climate change (Pfenninger et al. 2011), local adaptation (Quintela et al. 2014; Johansson et al. 2016), hybridisation (Patel et al. 2015) and biodiversity (Albrecht et al. 2012). Despite this broad range of interests, genomic resources are scarce and limited to transcriptomes (Feldmeyer et al. 2011, 2015; Tills et al. 2015) and mitochondrial genomes (Feldmeyer et al. 2010).

Here, we present the annotated draft genome sequence for *Radix auricularia* L. This serves as an important foundation for future genomic and applied research in this scientifically important genus.

## Results and Discussion

### *Genome assembly*

A total of 1,000,372,010 raw reads (Supplementary Note 1; Supplementary Table 1) were generated and assembled into 4,823 scaffolds (Table 2; Supplementary Table 2). The mitochondrial genome (13,744 bp) was fully reconstructed, evidenced by comparison to the previously published sequence (Feldmeyer et al. 2015). Re-mapping the reads used for assembly revealed that 97.6 % could be unambiguously placed, resulting in a per position coverage distribution with its peak at 72 x (Supplementary Note 2; Supplementary Figure 1).

The cumulative length of all scaffolds sums up to 910 Mb, which is about 665 Mb below the genome length estimates resulting from flow cytometric analyses (1,575 Mb, Supplementary Note 3) and from a read-mapping analysis (1,603 Mb; Supplementary Note 4). This difference in length is most likely caused by a high repeat content in the *Radix* genome. Already within the scaffolds, 40.4 % of the sequence was annotated as repeats. Analysing the read coverage distribution along the scaffolds revealed furthermore a pronounced increase towards contig ends (Supplementary Figure 3). Such a pattern is typical for collapsed repeat stretches, and a repeat annotation in these regions supported this hypothesis (Supplementary Figure 4). The genome size estimates from the read mapping coverage are consistent with estimates from flow cytometry. This indicates an approximately uniform coverage of the nuclear genome in shotgun libraries without substantial bias introduced during library generation. The convex read coverage distribution along the continuous parts in the scaffolds as revealed by the re-mapping of sequence reads to the assembly (Supplementary Figure 3) with coverages achieving extreme values of 50,000 and more, indicating the presence of extensive collapsed repeat stretches in the present assembly. Furthermore the majority of reads could be mapped back to the assembly, which indicates the faithful representation of especially repeat sequences. The collapsed repeats at the end of contigs (Supplementary Figures 3 + 4) explain the majority of gaps and thus the lacking part of the assembly to the full estimated genome length. Overall, the repeat content of the genome was estimated to be approximately 70 % (Supplementary Note 5). The majority of repeats were either classified as Transposable Elements or as ‘unknown’ (Supplementary Figure 5). Compared to other published mollusc genomes, genome size and assembly length of this *R. auricularia* draft assembly are typical (Supplementary Figure 6), but when considering contiguity, N50 and completeness regarding metazoan core orthologues it ranks among the top mollusc genomes (Table 1; Table 2), with 94.6 % of

all metazoan core orthologues present, indicating that there are no conspicuous gaps in the assembled gene space (Simão et al. 2015).

### Annotation

The annotation resulted in 17,338 protein coding genes (Table 2), 70.4% of these show high sequence similarity to entries in the Swiss-Prot database (e-value  $< 10^{-10}$ , accessed on May 11<sup>th</sup> 2016). This number of genes is at the lower end compared to other annotated mollusc genomes (Min: *Lottia gigantea* 23,822; Max: *Crassostrea gigas* 45,406; Table 1; Supplementary Table 4). A substantial number of *Radix* genes (2,690) lack an ortholog in other molluscs.

We tested for overrepresentation of functional categories in *Radix* versus all other molluscs with available annotated genomes and vice versa. This identified nine Gene Ontology (GO) terms to be significantly enriched amongst the 2,690 proteins private to *Radix* compared to all other published and annotated mollusc genomes. These include “carbohydrate metabolic process”, “chitin catabolic process” and “nucleoside transmembrane transport” (Supplementary Table 5). Among the categories found in all annotated molluscs but *Radix*, the “G-protein coupled receptor signalling pathway” was significantly enriched (Supplementary Table 6). G-protein receptors are involved in reactions to “hormones, neurotransmitters and environmental stimulants” (Rosenbaum et al. 2009). It thus seems that *Radix auricularia* might have lost its sensitivity to some of these stimuli, however which ones still needs to be investigated. On the other hand membrane proteins are generally more diverse than water-soluble proteins in the tree of life (Sojo et al. 2016), so we hypothesise that its proteins could be highly modified and were thus not detected in *Radix*. A third possibility is the lack of those genes in the sequencing data or the annotation.

### Conclusion

Here we present a draft genome of the snail *Radix auricularia*. The genome is comparable in size to other mollusc genomes and like them, rich in repeats. This new genomic resource will allow conducting future studies on genome evolution, population genomics and gene evolution within this genus and higher gastropod and mollusc taxa.

## Material and Methods

### *Sample collection and sequencing*

Snails were collected from a pond in the Taunus, Germany, identified with COI barcoding (Pfenninger et al. 2006) and kept under laboratory conditions for several generations. Three specimens of *R. auricularia* (Figure 2) were used for DNA extraction. Pooled DNA was used for preparation of three paired end and three mate pair (2 kbp, 5 kbp, 10 kbp insert size) sequencing libraries, that were sequenced on an Illumina HiSeq 2000 and 2500 at Beijing Genomics Institute, Hong Kong (Supplementary Note 1; Supplementary Table 1). Reads were cleaned of adapter sequences using Trimmomatic 0.33 (Bolger et al. 2014) (Supplementary Note 6) and screened for contaminations with FastqScreen 0.5.2 ([http://www.bioinformatics.babraham.ac.uk/projects/fastq\\_screen/](http://www.bioinformatics.babraham.ac.uk/projects/fastq_screen/); Supplementary Note 7; Supplementary Figure 7). Raw reads have been deposited under NCBI BioProject PRJNA350764.

### *Genome size estimation*

Genome size was estimated by flow cytometry based on a modified protocol of (Otto 1990) (Supplementary Note 8). Additionally, we estimated the genome size from our sequence data by dividing the total sum of nucleotides used for the assembly by the peak coverage from mapping back the assembly reads with the *bwa mem* algorithm from BWA 0.5.10 (Li 2013) (Supplementary Note 4). Backmappings were also used to estimate the repeat content of the genome (Supplementary Note 5).

### *Assembly strategy*

Reads were assembled using the Platanus 1.2.1 pipeline (Kajitani et al. 2014) with kmer sizes between ranging from 63 to 88 and a step size of 2. All other assembly parameter were kept at the default value. From the output of assembly, scaffolding and gap closure with Platanus, all sequences with at least 500 bp were used for scaffolding with SSPACE 3.0 (Boetzer et al. 2011) with ‘contig extension’ turned on. To increase further the contiguity of the draft genome we applied a third scaffolding step, making use of the cDNA sequence data. Transcriptome contig sequences of *R. auricularia* and three closely related species (Supplementary Note 9) were mapped sequential according phylogeny (Feldmeyer et al. 2015) with BLAT 35 (Kent 2002), with *-extendThroughN* enabled apart from default settings, onto the scaffolds; the gapped alignments were then used for joining of sequences with L\_RNA\_scaffolder (Xue et al. 2013). Finally, all sequences with at least 1,000 bp were used as input

for GapFiller 1.10 (Boetzer et al. 2012) to close extant gaps in the draft genome (Supplementary Note 10).

### *Annotation strategy*

Metazoan core orthologous genes were searched in the *R. auricularia* assembly and all other available mollusc genomes using BUSCO 1.2b (Simão et al. 2015).

The whole annotation process was performed using the MAKER2 2.31.8 (Cantarel et al. 2008; Holt & Yandell 2011) pipeline and affiliated programs. Initially, we built a custom repeat library from the assembly using RepeatModeler 1.0.4 (Smit & Hubley 2015) and read data using dnaPipeTE 1.2 (Goubert et al. 2015) with 30 upstream trials on varying coverage depths and then 50 parallel runs on the best-fitting coverage of 0.025 (Supplementary Note 11). The draft genome and transcriptome of *R. auricularia* (Supplementary Note S9) in addition to the BUSCO 1.2b (Simão et al. 2015) annotations of core metazoan genes on the draft genome were used as input for the initial training at the Augustus webserver (Stanke et al. 2004) (<http://bioinf.uni-greifswald.de/webaugustus/>). As additional input for MAKER2 we created two hidden Markov models on the gene structure of *R. auricularia*. One was generated by GeneMark 4.32 (Lomsadze et al. 2005) and another by SNAP 2006-07-28 (Korf 2004), using the output of CEGMA v2.5 (Parra et al. 2007). We ran three consecutive iterations of MAKER2 with the draft genome sequence, the transcriptomes (Supplementary Note 9), models from Augustus, SNAP and GeneMark, the repeat library and the Swiss-Prot database (Accessed at May 23<sup>th</sup> 2016). Between the iterations, the Augustus 3.2.2 (Stanke et al. 2004) and SNAP models were retrained according to best practice of MAKER2 workflow (Supplementary Note 12). Finally, all protein sequences from MAKER2 output were assigned putative names by BLASTP searches (Camacho et al. 2009) against the Swiss-Prot database.

We created orthologous groups from protein sequences of all six annotated molluscs with OrthoFinder 0.7.1 (Emms & Kelly 2015). All proteins were functionally annotated using InterProScan 5 (Quevillon et al. 2005; Zdobnov & Apweiler 2001). The enrichment analyses were performed in TopGO (Alexa & Rahnenfuhrer 2016), a bioconductor package for R (R Development Core Team 2008). We checked for significant enrichment of GO terms in proteins private to *Radix* and

proteins found in all molluscs but *Radix*. We applied Fischer's exact test, FDR correction and filtered by q-values smaller than 0.05.

## Supplementary Material

Supplementary Note 1	Individuals and Sequencing
Supplementary Table 1	Sequenced raw data
Supplementary Table 2	Summary statistics of different assembly steps
Supplementary Note 2	Backmapping
Supplementary Figure 1	Per base coverage frequency distribution
Supplementary Figure 2	Mate-pair Insert size distribution
Supplementary Table 3	Backmapping statistics from mate pair libraries
Supplementary Note 3	Results from flow cytometric measurements
Supplementary Note 4	Genome size estimation from coverage
Supplementary Figure 3	Coverage of continuous parts of the scaffolds
Supplementary Figure 4	Annotated repeats along continuous parts of the scaffolds
Supplementary Note 5	Repeat content
Supplementary Figure 5	Classified repeat families
Supplementary Figure 6	Mollusc genome sizes
Supplementary Table 4	Proteins similar to Swiss-Prot
Supplementary Table 5	Enriched GO-terms for <i>Radix</i>
Supplementary Table 6	Enriched GO-terms in molluscs but <i>Radix</i>
Supplementary Note 6	Preprocessing and trimming
Supplementary Note 7	Contamination screening
Supplementary Figure 7	Results of contamination screening
Supplementary Note 8	Material and Methods of Flow cytometric analysis
Supplementary Note 9	Transcriptome assemblies
Supplementary Note 10	Genome assembly
Supplementary Figure 8	Transcriptome filtering
Supplementary Table 7	Genome scaffolding with transcriptomic data
Supplementary Note 11	Repeat library
Supplementary Figure 9	Read subsampling
Supplementary Note 12	Annotation

## Acknowledgments

The project was funded by SAW-network project 291.



## Literature cited

- Albertin CB et al. 2015. The octopus genome and the evolution of cephalopod neural and morphological novelties. *Nature*. 524:220–224. doi: 10.1038/nature14668.
- Albrecht C, Hauße T, Schreiber K, Wilke T. 2012. Mollusc biodiversity in a European ancient lake system: Lakes Prespa and Mikri Prespa in the Balkans. *Hydrobiologia*. 682:47–59. doi: 10.1007/s10750-011-0830-1.
- Alexa A, Rahnenführer J. 2016. topGO: Enrichment Analysis for Gene Ontology. R package version 2.26.0.
- Barghi N, Concepcion GP, Olivera BM, Lluisma AO. 2016. Structural features of conopeptide genes inferred from partial sequences of the *Conus tribblei* genome. *Mol. Genet. Genomics*. 291:411–422. doi: 10.1007/s00438-015-1119-2.
- Boetzer M et al. 2012. Toward almost closed genomes with GapFiller. *Genome Biol*. 13:R56. doi: 10.1186/gb-2012-13-6-r56.
- Boetzer M, Henkel C V., Jansen HJ, Butler D, Pirovano W. 2011. Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics*. 27:578–579. doi: 10.1093/bioinformatics/btq683.
- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics*. doi: 10.1093/bioinformatics/btu170.
- Camacho C et al. 2009. BLAST plus: architecture and applications. *BMC Bioinformatics*. 10:1. doi: Artn 421\nDoi 10.1186/1471-2105-10-421.
- Cantarel BL et al. 2008. MAKER: An easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Res*. 18:188–196. doi: 10.1101/gr.6743907.
- Dunn CW, Giribet G, Edgecombe GD, Hejnol A. 2014. Animal Phylogeny and Its Evolutionary Implications. *Annu. Rev. Ecol. Evol. Syst.* 45:371–395. doi: 10.1146/annurev-ecolsys-120213-091627.
- Dunn CW, Ryan JF. 2015. The evolution of animal genomes. *Curr. Opin. Genet. Dev.* 35:25–32. doi: 10.1016/j.gde.2015.08.006.
- Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. *Genome Biol*. 16:157. doi: 10.1186/s13059-015-0721-2.
- Feldmeyer B, Greshake B, Funke E, Ebersberger I, Pfenninger M. 2015. Positive selection in development and growth rate regulation genes involved in species divergence of the genus *Radix*. *BMC Evol. Biol.* 15:164. doi: 10.1186/s12862-015-0434-x.
- Feldmeyer B, Hoffmeier K, Pfenninger M. 2010. The complete mitochondrial genome of *Radix balthica* (Pulmonata, Basommatophora), obtained by low coverage shot gun next generation sequencing. *Mol. Phylogenet. Evol.* 57:1329–1333. doi: 10.1016/j.ympev.2010.09.012.
- Feldmeyer B, Wheat CW, Krezdorn N, Rotter B, Pfenninger M. 2011. Short read Illumina data for the de novo assembly of a non-model snail species transcriptome (*Radix balthica*, Basommatophora, Pulmonata), and a comparison of assembler performance. *BMC Genomics*. 12:317. doi: 10.1186/1471-2164-12-317.
- Garbar A V, Kornushin A V. 2003. Karyotypes of European species of *Radix* (Gastropoda: Pulmonata: Lymnaeidae) and their relevance to species distinction in the genus. *Malacologia*. 45:141–148. <Go to ISI>://000184914300008.
- GIGA Community of Scientists. 2014. The Global Invertebrate Genomics Alliance (GIGA): Developing Community Resources to Study Diverse Invertebrate Genomes. *J. Hered.* 105:1–18. doi: 10.1093/jhered/est084.
- González-Tizón a M, Martínez-Lage a, Rego I, Ausió J, Méndez J. 2000. DNA content, karyotypes, and chromosomal location of 18S-5.8S-28S ribosomal loci in some species of bivalve molluscs from the Pacific Canadian coast. *Genome*. 43:1065–1072. doi: 10.1139/gen-43-6-1065.
- Goubert C et al. 2015. De novo assembly and annotation of the Asian tiger mosquito (*Aedes albopictus*) repeatome with dnaPipeTE from raw genomic reads and comparative analysis with the yellow fever mosquito (*Aedes aegypti*). *Genome Biol. Evol.* 7:1192–1205. doi: 10.1093/gbe/evv050.
- Gregory TR. 2003. Genome size estimates for two important freshwater molluscs, the zebra mussel (*Dreissena polymorpha*) and the schistosomiasis vector snail (*Biomphalaria glabrata*). *Genome*. 46:841–4. doi: 10.1139/g03-069.

- Hallgren P, Sorita Z, Berglund O, Persson A. 2012. Effects of 17 $\alpha$ -ethinylestradiol on individual life-history parameters and estimated population growth rates of the freshwater gastropods *Radix balthica* and *Bithynia tentaculata*. *Ecotoxicology*. 21:803–810. doi: 10.1007/s10646-011-0841-8.
- Hinegardner R. 1974. Cellular DNA content of the Mollusca. *Comp. Biochem. Physiol. -- Part A Physiol.* 47:447–460. doi: 10.1016/0300-9629(74)90008-5.
- Holt C, Yandell M. 2011. MAKER2 : an annotation pipeline and genome- database management tool for second- generation genome projects. *BMC Bioinformatics*. 12:491. doi: 10.1186/1471-2105-12-491.
- Huňová K et al. 2012. *Radix* spp.: Identification of trematode intermediate hosts in the Czech Republic. *Acta Parasitol.* 57:273–84. doi: 10.2478/s11686-012-0040-7.
- Jarne P, Charlesworth D. 1993. The Evolution of the Selfing Rate in Functionally Hermaphrodite Plants and Animals. *Annu. Rev. Ecol. Syst.* 24:441–466. <http://www.jstor.org/stable/2097186>.
- Jarne P, Delay B. 1990. Inbreeding depression and self-fertilization in *Lymnaea peregra* (Gastropoda: Pulmonata). *Heredity (Edinb)*. 64:169–176. doi: 10.1038/hdy.1990.21.
- Johansson MP, Ermold F, Kristjánsson BK, Laurila A. 2016. Divergence of gastropod life history in contrasting thermal environments in a geothermal lake. *J. Evol. Biol.* 1–11. doi: 10.1111/jeb.12928.
- Kenny NJ, Namigai EKO, Marlétaz F, Hui JHL, Shimeld SM. 2015. Draft genome assemblies and predicted microRNA complements of the intertidal lophotrochozoans *Patella vulgata* (Mollusca, Patellogastropoda) and *Spirobranchus (Pomatoceros) lamarcki* (Annelida, Serpulida). *Mar. Genomics*. 24:139–146. doi: 10.1016/j.margen.2015.07.004.
- Kent WJ. 2002. BLAT — The BLAST -Like Alignment Tool. *Genome Res.* 12:656–664. doi: 10.1101/gr.229202.
- Korf I. 2004. Gene finding in novel genomes. *BMC Bioinformatics*. 5:59. doi: 10.1186/1471-2105-5-59.
- Lasek RJ, Dower WJ. 2013. *Aplysia californica* : Analysis of Nuclear DNA in Individual Nuclei of Giant Neurons *Aplysia californica* : Analysis of Nuclear DNA in Individual Nuclei of Giant Neurons. 172:278–280.
- Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv Prepr.* doi: arXiv:1303.3997 [q-bio.GN].
- Lomsadze A, Ter-Hovhannisyan V, Chernoff YO, Borodovsky M. 2005. Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res.* 33:6494–6506. doi: 10.1093/nar/gki937.
- Matty Knight, Coen M. Adema, Nithya Raghavan, Eric S. Loker FAL and HT. *Obtaining the genome sequence of the mollusc Biomphalaria glabrata: a major intermediate host for the parasite causing human schistosomiasis.*
- Moroz LL et al. Sequencing the *Aplysia* Genome: a model for single cell, real-time and comparative genomics. <http://www.genome.gov/pages/research/sequencing/seqproposals/aplysiaseq.pdf>.
- Nguyen TTT, Hayes BJ, Ingram BA. 2014. Genetic parameters and response to selection in blue mussel (*Mytilus galloprovincialis*) using a SNP-based pedigree. *Aquaculture*. 420–421:295–301. doi: 10.1016/j.aquaculture.2013.11.021.
- Otto F. 1990. DAPI Staining of Fixed Cells for High-Resolution Flow Cytometry of Nuclear DNA. *Methods Cell Biol.* 33:105–110. doi: 10.1016/S0091-679X(08)60516-6.
- Parra G, Bradnam K, Korf I. 2007. CEGMA: A pipeline to accurately annotate core genes in eukaryotic genomes. *Bioinformatics*. 23:1061–1067. doi: 10.1093/bioinformatics/btm071.
- Patel S, Schell T, Eifert C, Feldmeyer B, Pfenninger M. 2015. Characterizing a hybrid zone between a cryptic species pair of freshwater snails. *Mol. Ecol.* 24:643–655. doi: 10.1111/mec.13049.
- Peñarrubia L, Araguas RM, et al. 2015. Identification of 246 microsatellites in the Asiatic clam (*Corbicula fluminea*). *Conserv. Genet. Resour.* 7:393–395. doi: 10.1007/s12686-014-0378-2.
- Peñarrubia L, Sanz N, Pla C, Vidal O, Viñas J. 2015. Using massive parallel sequencing for the development, validation, and application of population genetics markers in the invasive bivalve zebra mussel (*Dreissena polymorpha*). *PLoS One*. 10:1–14. doi: 10.1371/journal.pone.0120732.
- Pfenninger M, Cordellier M, Streit B. 2006. Comparing the efficacy of morphologic and DNA-based taxonomy in the freshwater gastropod genus *Radix* (Basommatophora, Pulmonata). *BMC Evol. Biol.* 6:100. doi: 10.1186/1471-2148-6-100.

- Pfenninger M, Salinger M, Haun T, Feldmeyer B. 2011. Factors and processes shaping the population structure and distribution of genetic variation across the species range of the freshwater snail *radix balthica* (Pulmonata, Basommatophora). *BMC Evol. Biol.* 11:135. doi: 10.1186/1471-2148-11-135.
- Quevillon E et al. 2005. InterProScan: Protein domains identifier. *Nucleic Acids Res.* 33:116–120. doi: 10.1093/nar/gki442.
- Quintela M, Johansson MP, Kristjánsson BK, Barreiro R, Laurila A. 2014. AFLPs and Mitochondrial haplotypes reveal local adaptation to extreme thermal environments in a freshwater gastropod. *PLoS One.* 9. doi: 10.1371/journal.pone.0101821.
- R Development Core Team. 2008. R: A language and environment for statistical computing.
- Rodríguez-Juíz AM, Torrado M, Méndez J. 1996. Genome-size variation in bivalve molluscs determined by flow cytometry. *Mar. Biol.* 126:489–497. doi: 10.1007/BF00354631.
- Romero PE, Pfenninger M, Kano Y, Klussmann-Kolb A. 2015. Molecular phylogeny of the Ellobiidae (Gastropoda: Panpulmonata) supports independent terrestrial invasions. *Mol. Phylogenet. Evol.* doi: 10.1016/j.ympev.2015.12.014.
- Romero PE, Weigand AM, Pfenninger M. 2016. Positive selection on panpulmonate mitogenomes provide new clues on adaptations to terrestrial life. *BMC Evol. Biol.* 16. doi: 10.1186/s12862-016-0735-8.
- Rosenbaum DM, Rasmussen SGF, Kobilka BK. 2009. The structure and function of G-protein-coupled receptors. *Nature.* 459:356–363. doi: 10.1038/nature08144.
- Rundle SD, Smirthwaite JJ, Colbert MW, Spicer JJ. 2011. Predator cues alter the timing of developmental events in gastropod embryos. *Biol. Lett.* 7:285–287. doi: 10.1098/rsbl.2010.0658.
- Simakov O et al. 2013. Insights into bilaterian evolution from three spiralian genomes. *Nature.* 493:526–31. doi: 10.1038/nature11696.
- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM. 2015. BUSCO: Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.* 31:3210–3212. doi: 10.1093/bioinformatics/btv351.
- Simit A, Hubley R. 2015. RepeatModeler Open-1.0.
- Sojo V, Dessimoz C, Pomiankowski A, Lane N. 2016. Membrane proteins are dramatically less conserved than water-soluble proteins across the tree of life. *Mol. Biol. Evol.* 33:msw164. doi: 10.1093/molbev/msw164.
- Sommer U, Adrian R, Bauer B, Winder M. 2012. The response of temperate aquatic ecosystems to global warming: Novel insights from a multidisciplinary project. *Mar. Biol.* 159:2367–2377. doi: 10.1007/s00227-012-2085-4.
- Stanke M, Steinkamp R, Waack S, Morgenstern B. 2004. AUGUSTUS: A web server for gene finding in eukaryotes. *Nucleic Acids Res.* 32:309–312. doi: 10.1093/nar/gkh379.
- Tills O et al. 2011. A genetic basis for intraspecific differences in developmental timing? *Evol. Dev.* 13:542–548. doi: 10.1111/j.1525-142X.2011.00510.x.
- Tills O, Truebano M, Rundle S. 2015. An embryonic transcriptome of the pulmonate snail *Radix balthica*. *Mar. Genomics.* 24:259–260. doi: 10.1016/j.margen.2015.07.014.
- Vinogradov AE. 1998. Variation in ligand-accessible genome size and its ecomorphological correlates in a pond snail. *Cytometry.* 65:59–65. doi: 10.1111/j.1601-5223.1998.00059.x.
- Wiehn JR, Kopp K, Rezzonico S, Karttunen S, Jokela J. 2002. Family-Level Covariation Between Parasite Resistance and Mating System in a Hermaphroditic Freshwater Snail. *Evolution (N. Y.)*. 56:1454–1461. doi: 10.1111/j.0014-3820.2002.tb01457.x.
- Xue W et al. 2013. L\_RNA\_scaffolder: scaffolding genomes with transcripts. *BMC Genomics.* 14. doi: 10.1186/1471-2164-14-604.
- Yu TL, Deng YH, Zhang J, Duan LP. 2016. Size-assortative copulation in the simultaneously hermaphroditic pond snail *Radix auricularia* (Gastropoda: Pulmonata). *Anim. Biol.* doi: 10.1163/15707563-00002501.
- Zdobnov EM, Apweiler R. 2001. InterProScan - an integration platform for the signature-recognition methods in InterPro. *Bioinformatics.* 17:847–848. doi: 10.1093/bioinformatics/17.9.847.
- Zhang G et al. 2012. The oyster genome reveals stress adaptation and complexity of shell formation. *Nature.* 490:49–54. doi: 10.1038/nature11413.

## Tables

Table 1. Available mollusc genomes. An overview from column 2 can be found in Supplementary Figure 6. Column 5: Fraction of N's in the assembly. Column 6: BUSCOs: (Benchmarking Universal Single-Copy Orthologs) N<sub>Metazoa</sub>=843; Present = complete + fragmented.

Species	Assembly length / Estimated genome size = % assembled	#sequences / N50 (*contigs)	Coverage / Technology	Gap [%]	BUSCOs present	Number of annotated proteins
<i>Octopus bimaculoides</i> <sup>a</sup>	2.4Gb/2.7Gb = 89%	151674 / 475kb	92 Illumina	15.1	73,8	23,994
<i>Dreissena polymorpha</i> <sup>b</sup>	906Mb / 1.7 Gb <sup>c</sup> = 0.06%	* 1057 / 855 bp	3 454	0	0	-
<i>Corbicula fluminea</i> <sup>d</sup>	663 Mb / ?	* 778 / 849 bp	3 454	0	0	-
<i>Crassostrea gigas</i> <sup>e</sup>	558 Mb / 890 Mb <sup>f</sup> = 62.7%	7659 / 402 kb	100 Illumina	11.8	82	45,406
<i>Mytilus galloprovincialis</i> <sup>g</sup>	1.6Gb / 1.9Gb <sup>h</sup> = 86%	* 2315965 / 1067 bp	17 Illumina	0	1.6	-
<i>Lottia gigantea</i> <sup>i</sup>	360 Mb / 421 Mb <sup>j</sup> = 85%	4469 / 1870 kb	8.87 Sanger	16.9	97.0	23,822
<i>Patella vulgata</i> <sup>k</sup>	579 Mb / 1,460 Mb = 39.7%	295348 / 3160	25.6 Illumina	0.00062	16.6	-

<i>Conus tribblei</i> <sup>1</sup>	2,160 Mb / 2,757 Mb = 78%	1126156 / 2681 bp	28.5 Illumina	0	44	-
<i>Aplysia californica</i> <sup>m</sup>	927 Mb / 1,760 Mb <sup>j</sup> resp. 1,956 Mb <sup>n</sup> = 53 resp. 47%	4,332 / 918 kb	66 Illumina	20.4	94.1	27,591
<i>Biomphalaria glabrata</i> <sup>o</sup>	916 Mb / 929 Mb <sup>c</sup> = 99%	331,401 / 48 kb	27.5 454	1.9	89.1	36,675
<i>Lymnaea stagnalis</i> <sup>p</sup>	833 Mb / 1,193 Mb <sup>q</sup> = 70%	* 328,378 / 5.8 kb		0	88	-

---

References: Genome sizes are from the genome publications, if not cited separately. <sup>a</sup>(Albertin et al. 2015); <sup>b</sup>(Peñarrubia, Sanz, et al. 2015); <sup>c</sup>(Gregory 2003); <sup>d</sup>(Peñarrubia, Araguas, et al. 2015); <sup>e</sup>(Zhang et al. 2012); <sup>f</sup>(González-Tizón et al. 2000); <sup>g</sup>(Nguyen et al. 2014); <sup>h</sup>(Rodríguez-Juiz et al. 1996); <sup>i</sup>(Simakov et al. 2013); <sup>j</sup>(Hinegardner 1974); <sup>k</sup>(Kenny et al. 2015); <sup>l</sup>(Barghi et al. 2016); <sup>m</sup>(Moroz et al.) GCF\_000002075.1; <sup>n</sup>(Lasek & Dower 2013); <sup>o</sup>(Matty Knight, Coen M. Adema, Nithya Raghavan, Eric S. Loker) GCF\_000457365.1; <sup>p</sup>(unpublished - Ashworth Laboratories 2016) GCA\_900036025.1; <sup>q</sup>(Vinogradov 1998)

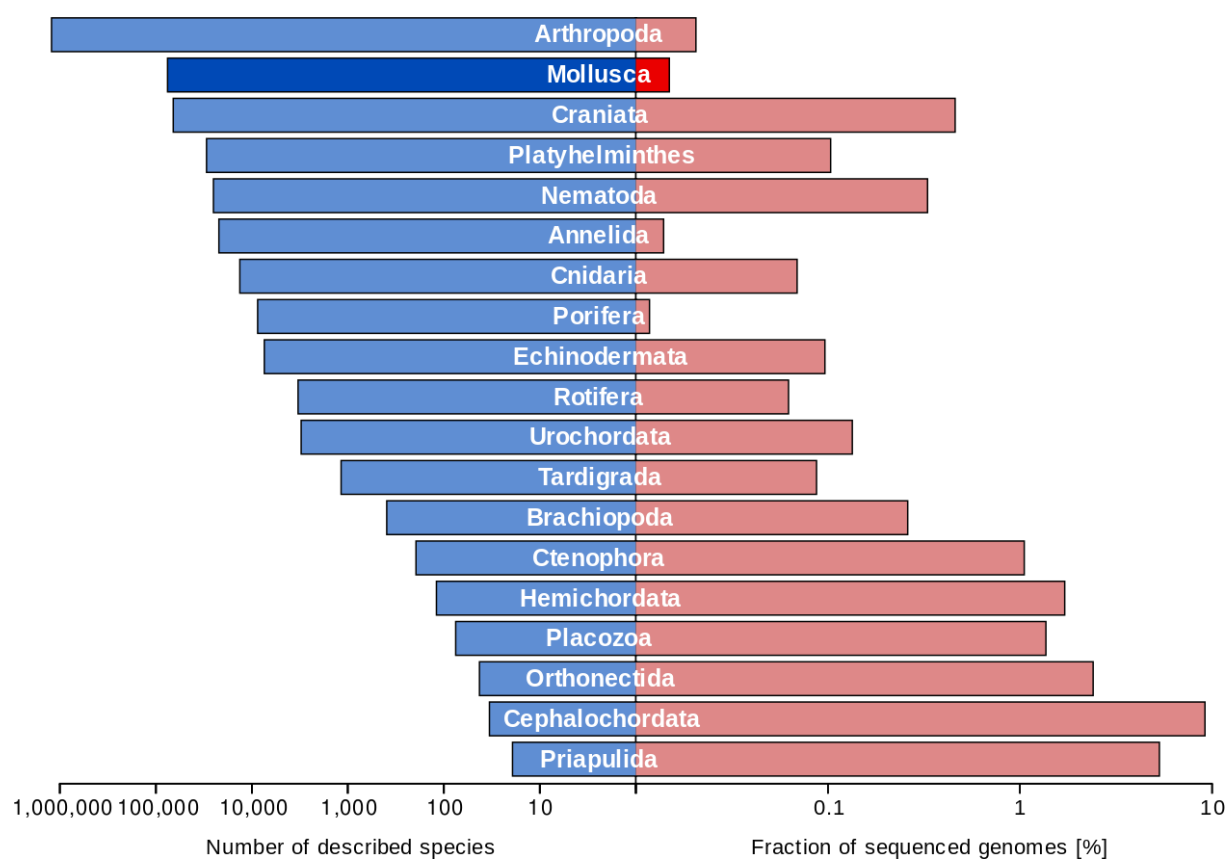
Table 2. Characteristics of the *R. auricularia* genome and draft assembly. BUSCOs: (Benchmarking Universal Single-Copy Orthologs) N<sub>Metazoa</sub>=843; Present = complete + fragmented

Parameter	Value
Haploid chromosome number	17 (Garbar & Korniushev 2003)
Estimated genome length	
Flow cytometry	1.51 Gb (Vinogradov 1998)
	1.58 Gb ± 21.5 Mb (this study)
Sequencing coverage	1.60 Gb
Total assembly length	0.91 Gb single copy or high complexity regions
#scaffolds	4,823
N50	578,730 bp
Gaps	6.4% N
Coverage	72x
Estimation of gene completeness	94.6 % of BUSCO genes present
Gene prediction	17,338 genes
Gene space (UTR, Exons, Introns etc.)	200.6 Mb = 21.9% of assembly
Gene length (median)	8.0 kb
Gene fragmentation	147,195 exons
Exon space	25.3 Mb = 2.8% of assembly (1.6% of total genome)
Exon length (median)	125 bp
Protein length (median)	332 AA

## Figure Legends

**Fig. 1.** The number of described species (Dunn & Ryan 2015; GIGA Community of Scientists 2014) and the fraction of sequenced genomes (<http://www.ncbi.nlm.nih.gov/genome/browse/> on September 1<sup>st</sup>, 2016). Animal phyla were obtained from (Dunn et al. 2014). Phyla with genomic record are displayed only.

**Fig. 2.** Photograph of *Radix auricularia*. Photo by Markus Pfenninger.



**Fig. 1.**





**Fig. 2.**