

1 ***De novo* genome assembly of the land snail *Candidula unifasciata* (Mollusca: Gastropoda)**

2  
3 Luis J. Chueca<sup>1,2,\*</sup>, Tilman Schell<sup>1</sup>, Markus Pfenninger<sup>1,3,4</sup>

4  
5 <sup>1</sup> LOEWE-Centre for Translational Biodiversity Genomics (LOEWE-TBG), Senckenberg Nature  
6 Research Society, Frankfurt am Main, Germany.

7 <sup>2</sup> Department of Zoology and Animal Cell Biology, University of the Basque Country (UPV-EHU),  
8 Vitoria-Gasteiz, Spain.

9 <sup>3</sup> Senckenberg Biodiversity and Climate Research Centre (SBIK-F), Frankfurt am Main, Germany.

10 <sup>4</sup> Institute of Organismic and Molecular Evolution (iOME), Faculty of Biology, Johannes  
11 Gutenberg University, Mainz, Germany.

12  
13 \*Corresponding author: Luis J. Chueca

14 Email: [luisjavier.chueca@ehu.eus](mailto:luisjavier.chueca@ehu.eus)

15  
16 **Keywords**

17  
18 Annotation, *de novo* assembly, Geomitridae, land snails, long reads, molluscs, repeats

19  
20 **Abstract**

21  
22 Among all molluscs, land snails are an economically and scientifically interesting group comprising  
23 edible species, alien species and agricultural pests. Yet, despite its high diversity, the number of  
24 whole genomes publicly available is still scarce. Here, we present the draft genome assembly of the  
25 land snail *Candidula unifasciata*, a widely distributed species along central Europe, which belongs  
26 to Geomitridae family, a group highly diversified in the Western-Palearctic region. We performed a  
27 whole genome sequencing, assembly and annotation of an adult specimen based on PacBio and  
28 Oxford Nanopore long read sequences as well as Illumina data. A genome of about 1.29 Gb was  
29 generated with a N50 length of 246 kb. More than 60% of the assembled genome was identified as  
30 repetitive elements, and 22,464 protein-coding genes were identified in the genome, where the  
31 62.27% were functionally annotated. This is the first assembled and annotated genome for a  
32 geometrid snail and will serve as reference for further evolutionary, genomic and population genetic  
33 studies of this important and interesting group.

## 35 **1. Introduction**

36

37 Gastropods are the largest group among molluscs, representing almost the 80% of the species.

38 Although most of the them are present in marine habitats, land snails diversity is estimated around

39 35.000 species (Solem 1984). Due to its low dispersal abilities, land snails have been employed in

40 many evolutionary and population genomics studies (Stankowski 2013; Schilthuizen and

41 Kellermann 2014; Chueca *et al.* 2017; Haponski *et al.* 2017). While these studies are mainly based

42 on few loci, transcriptomes or mitochondrial genomes (Kang *et al.* 2016; Romero *et al.* 2016;

43 Razkin *et al.* 2016; Korábek *et al.* 2019), only a couple of whole nuclear genomes of land snails

44 species are available so far. Geomitridae is one of the most diverse families of molluscs in Western-

45 Palearctic region. The family is composed by small to medium-size species, characterized by

46 presenting several reproductive adaptations to xeric habitats (Giusti and Manganelli 1987).

47 *Candidula unifasciata* (NCBI:txid100452) is a land snail species widely distributed along western

48 Europe, from southern France and Italy to central and northern Europe (Fig. 1). *C. unifasciata*

49 inhabits dry meadows and open lowlands with rocks, being also present in gardens and vineyards. A

50 recent molecular revision of *Candidula* (Chueca *et al.* 2018) revealed the polyphyly of the genus,

51 and split the species that composed it into six genera, questioning the traditional anatomical

52 classification. Although, there are many taxonomical, phylogeographical and evolutionary studies

53 concerning Geomitridae species (Pfenninger and Magnin 2001; Sauer and Hausdorf 2010; Brozzo

54 *et al.* 2020), the lack of reference genomes makes it difficult to investigate deeper biological and

55 evolutionary questions about geomitrids and other land snails species. Here, we present the

56 annotated draft genome of *Candidula unifasciata* that will be a valuable resource for future genomic

57 research of this important taxonomic group.

58

## 59 **2. Materials and Methods**

60

### 61 **2.1. Sample collection, library construction, sequencing**

62

63 A live population of *C. unifasciata* was collected from Winterscheid, Gilserberg, Gemany (50.93°

64 N, 9.04° E). Genomic DNA was extracted from one specimen using the phenol/chloroform method

65 and quality was checked by gel electrophoresis and NanoDrop ND-1000 spectrophotometer

66 (LabTech, USA). A total of 5.6 µg of DNA was sent to Novogene (UK) for library preparation and

67 sequencing. Then, a 300 base pair (bp) insert DNA libraries were generated using NEBNext® DNA

68 Library Prep Kit and sequenced on 3 lanes of Illumina NovaSeq 6000 platform (150 bp paired-end

69 [PE] reads). Quality of raw Illumina sequences was checked with FastQC (Andrews 2010). Low

70 quality bases and adapter sequences were subsequently trimmed by Trimmomatic v0.39 (Bolger *et*  
71 *al.* 2014). For PacBio sequencing, a DNA library was prepared from 5 µg of DNA using the  
72 SMRTbell template prep kit v.1.0. Sequencing was carried out on 10 single-molecule real-time  
73 sequencing (SMRT) cells on an RSI instrument using P6-C4 chemistry.

74

75 To obtain Oxford Nanopore Technologies (ONT) long reads, we ran two flow cells on a MinION  
76 portable sequencer. Total genomic DNA was used for library preparation with the Ligation  
77 Sequencing kit (SQK-LSK109) from ONT, using the manufacturer's protocols. Base calling of the  
78 reads from the two MinION flow cells was performed with guppy v4.0.11  
79 (<https://nanoporetech.com/nanopore-sequencing-data-analysis>), under default settings. Afterwards,  
80 ONT reads quality was checked with Nanoplot v1.28.1 (<https://github.com/wdecoster/NanoPlot>)  
81 and reads shorter than 1000 bases and mean quality below seven were discarded by running  
82 Nanofilt v2.6.0 (<https://github.com/wdecoster/nanofilt>).

83

84 Two specimens, one adult and one juvenile, were ground together into small pieces using steel balls  
85 and a Retsch Mill. Then, RNA was extracted following an standard Trizol extraction. The integrity  
86 of total RNA extracted was assessed on an Agilent 4200 TapeStation (Agilent, USA), after which,  
87 approximately 1 µg of the total RNA was processed using the Universal Plus mRNA-seq library  
88 preparation kit (NuGEN, Redwood City, CA). Finally, the 300-bp insert size library was sequenced  
89 on a Illumina NovaSeq 6000 platform.

90

## 91 **2.2. Genome size estimation**

92

93 Genome size was estimated following a flow cytometry protocol with propidium iodide-stained  
94 nuclei described in (Hare and Johnston 2012). Foot tissue of one fresh adult sample of *C.*  
95 *unifasciata* and neural tissue of the internal reference standard *Acheta domesticus* (female, 1C = 2  
96 Gb) was mixed and chopped with a razor blade in a petri dish containing 2 ml of ice-cold Galbraith  
97 buffer. The suspension was filtered through a 42-µm nylon mesh and stained with the intercalating  
98 fluorochrome propidium iodide (PI, Thermo Fisher Scientific) and treated with RNase II A (Sigma-  
99 Aldrich), each with a final concentration of 25 µg/ml. The mean red PI fluorescence signal of  
100 stained nuclei was quantified using a Beckman-Coulter CytoFLEX flow cytometer with a solid-  
101 state laser emitting at 488 nm. Fluorescence intensities of 5000 nuclei per sample were recorded.  
102 We used the software CytExpert 2.3 for histogram analyses. The total quantity of DNA in the  
103 sample was calculated as the ratio of the mean red fluorescence signal of the 2C peak of the stained  
104 nuclei of the *C. unifasciata* sample divided by the mean fluorescence signal of the 2C peak of the

105 reference standard times the 1C amount of DNA in the standard reference. Four replicates were  
106 measured to minimize possible random instrumental errors. Furthermore, we estimated the genome  
107 size by coverage from mapping reads used for genome assembly back to the assembly itself using  
108 backmap v0.3 (<https://github.com/schell/backmap>; Schell *et al.* 2017). In brief, the method divides  
109 the number of mapped nucleotides by the mode of the coverage distribution. By doing so, the length  
110 of collapsed regions with many fold increased coverage is taken into account.

111

### 112 **2.3 Genome assembly workflow**

113

114 Different *de novo* genome assemblies were tested under different methods (see Table S1). The  
115 pipeline, which showed the best genome, was selected to continue further analyses. The draft  
116 genome was constructed from PacBio long reads using wtdbg2 v2.5 (Ruan and Li 2020), followed  
117 by three polishing rounds of Racon 1.4.3 (Vaser *et al.* 2017) and three polishing rounds of Pilon  
118 1.23 (Walker *et al.* 2014). After that, Illumina and PacBio reads were aligned to the assembly using  
119 backmap.pl v0.3 to evaluate coverage distribution. Then, Purge Haplotigs (Roach *et al.* 2018) was  
120 employed, under default parameters and cut off values of 15, 72 and 160 to identify and remove  
121 redundant contigs.

122

### 123 **2.4. Scaffolding and gap closing**

124

125 To further improve the assembly, we applied two rounds of scaffolding and gap closing to the  
126 selected genome assembly. The genome was first scaffolded with the SMRT and ONT reads by  
127 LINKS v1.8.7 (Warren *et al.* 2015) and then with RNA reads by Rascaf v1.0.2 (Song *et al.* 2016).  
128 Long-Read Gapcloser v1.0 (Xu *et al.* 2018) was run three times after each scaffolding step,  
129 followed by three polishing rounds of Racon v1.4.3. BlobTools v.1.0 (Kumar *et al.* 2013; Laetsch  
130 and Blaxter 2017) was employed to screen genome assembly for potential contamination by  
131 evaluating coverage, GC content and sequence similarity against the NCBI nt database of each  
132 sequence. The resulting assembly was compared in terms of contiguity using Quast v5.0.2  
133 (Gurevich *et al.* 2013), and evaluated for completeness by BUSCO v3.02 (Simão *et al.* 2015)  
134 against metazoa\_odb9 data set.

135

136

### 137 **2.5. Transcriptome assembly**

138

139 RNA reads were also checked for quality and trimmed, as was explained above, and the  
140 transcriptome was assembled using Trinity v2.9.1 (Haas *et al.* 2013). Then, the transcriptome  
141 assembly was evaluated for completeness by BUSCO v3.0.2 against the against metazoa\_odb9 data  
142 set. Moreover, the clean RNA-seq reads from different specimens were aligned against the  
143 reference genome by HISAT2 (Kim *et al.* 2015).

144

## 145 **2.6. Repeat Annotation**

146

147 RepeatModeler v2.0 (Smit and Hubley 2008) was run to construct a *de novo* repetitive library from  
148 the assembly. The resulting repetitive library created was employed by RepeatMasker v4.1.0  
149 (<http://www.repeatmasker.org/>) to annotate and masked the genome.

150

## 151 **2.7. Gene prediction and functional annotation.**

152

153 Genes were predicted by using different methods. First, genes models were predicted *ab initio*  
154 based on SNAP v. 2006-07-28 (Korf 2004) and the candidates coding regions within the assembled  
155 transcript were identified with TransDecoder v5.5.0 (<https://github.com/TransDecoder/>). Secondly,  
156 we used homology-based gene predictions by aligning protein sequences from SwissProt (2020-04)  
157 to the *Candidula unifasciata* masked genome with EXONERATE 2.2.0 (Slater and Birney 2005)  
158 and by running GeMoMa v1.7.1 (Keilwagen *et al.* 2016, 2018) taking five gastropods species as  
159 reference organisms. The selected species were *Pomacea canaliculata* (GCF\_003073045.1; (Liu *et al.*  
160 *al.* 2018), *Aplysia californica* (GCF\_000002075.1), *Elysia chlorotica* (GCA\_003991915.1; (Cai *et al.*  
161 *al.* 2019), *Radix auricularia* (GCA\_002072015.1; (Schell *et al.* 2017) and *Chrysomallon*  
162 *squamiferum* (GCA\_012295275.1; (Sun *et al.* 2020), which were downloaded from NCBI. First,  
163 from the mapped RNA-seq reads, introns were extracted and filtered by the GeMoMa modules ERE  
164 and DenoiseIntrons. Then, we ran independently the module GeMoMa pipeline for each reference  
165 species using mmseqs2 and including the RNA-seq data. The five gene annotations were then  
166 combined into a final annotation file by using the GeMoMa modules GAF and AnnotationFinalizer.  
167 Finally, we aligned *C. unifasciata* transcripts against the masked genome using PASA v2.4.1  
168 (Campbell *et al.* 2006) as implemented in autoAug.pl.

169

170 Gene prediction data from each method were combined using EVidenceMolder v1.1.1 (Haas *et al.*  
171 2008) to obtain a consensus gene set for the raccoon-dog genome. Gene models from GeMoMa and  
172 SNAP were converted to EVM compatible gff3 files and combined with CDS identified by  
173 TransDecoder into a gene predictions file. After that, EVM was run including gene model

174 predictions, protein and transcript alignments and repeat regions to produce a reliable consensus  
175 gene set.

176

177 Predicted genes were annotated by BLAST search against the Swiss-Prot database with an e-value  
178 cutoff of  $10^{-6}$ . InterProScan v5.39.77 (Quevillon *et al.* 2005) was used to predict motifs and  
179 domains, as well as Gene ontology (GO) terms.

180

### 181 **3. Results and Discussion**

182

#### 183 **3.1 Genome assembly**

184

185 The calculated DNA content through flow cytometry experiments was 1.54 Gb. The genome size  
186 estimation by Illumina read coverage resulted in 1.42 Gb. The estimated heterozygosity by  
187 GenomeScope of the specimen employed for genome assembly was around 1.09% (Fig. 2.a), being  
188 in the range of other land snail genomes (Guo *et al.* 2019; Saenko *et al.* 2021). We generated  
189 sequence data for a total coverage of approximately 120.6X and 25.6X of Illumina and PacBio  
190 reads respectively. After scaffolding with long reads (PacBio and ONT) and RNA data, we  
191 produced a draft genome assembly of 1.29 Gb with 8,586 scaffolds and a scaffold N50 of 246 kb  
192 (Table 1). Completeness evaluation by BUSCO against the metazoan\_odb9 data set showed high  
193 values, recovering more than the 92% as complete and less than the 6% as missing genes for both,  
194 assembly and annotation, analyses (Table 1). This results were in the range of other gastropods  
195 genome assemblies (Schell *et al.* 2017; Liu *et al.* 2018; Guo *et al.* 2019; Sun *et al.* 2020), being  
196 slightly better than closest relative assembly of *Cepaea nemoralis* (Saenko *et al.* 2021). For genome  
197 quality evaluation, we compared the *C. unifasciata* draft genome generated with other mollusc  
198 genomes publicly available. This comparison showed high quality in terms of contig number and  
199 scaffold N50 among land snail genomes. The mapping of the Illumina reads against the final  
200 genome assembly showed that the 98.56% of them were aligned to it, as well as a good removal of  
201 redundant contigs (Fig. 2b). Finally, BlobTools analysis didn't reflect substantial contamination  
202 (Fig. 3), indicating the reliability of the data.

203

#### 204 **3.2 Genome annotation**

205

206 We estimated the total repeat content of the *C. unifasciata* genome assembly around 61.10% (Table  
207 2), values slightly smaller than other land snails genomes (Guo *et al.* 2019; Saenko *et al.* 2021).  
208 Approximately one third of the assembled genome (33.96%) was identified as Transposable

209 elements (TEs) such as long interspersed nuclear elements (LINEs; 25.03%), short interspersed  
210 nuclear elements (SINEs; 4.23%), long tandem repeats (LTR; 0.60%) and DNA transposons  
211 (4.10%).

212

213 We predicted 22,464 genes in the *C. unifasciata* genome (Table 3) by using a homology-based gene  
214 prediction and EVM. Among the identified proteins, 13,221 (62.27%) were annotated to have at  
215 least one GO term. Finally, 21,231 proteins (94.51%) were assigned to at least one of the database  
216 from InterProScan (Table 3). BUSCO and functional annotations results indicated high quality.  
217 Total protein-coding genes was in the range of other gastropods annotations (Schell *et al.* 2017; Liu  
218 *et al.* 2018; Guo *et al.* 2019), however this number represented only the half of its closest relative  
219 *Cepaea nemoralis* (Saenko *et al.* 2021).

220

#### 221 **4. Conclusions**

222

223 Here, we present a draft assembled and annotated genome of the land snail *Candidula unifasciata*.  
224 The obtained genome is comparable with other land snail and Gastropoda genomes publicly  
225 available. The new genome resource will be reference for further population genetics, evolutionary  
226 and genomic studies of this highly world-wide diverse group.

227

#### 228 **Data Availability Statement**

229

230 All raw data generated for this study (Illumina, PacBio, MinION, and RNA-seq reads) are available  
231 at the European Nucleotide Archive database (ENA) under the Project number: PRJEB41346. The  
232 final genome assembly and annotation can be found under the accession number GCA\_905116865.

233

#### 234 **Competing interests**

235

236 The authors declare that they have no competing interests.

237

#### 238 **Acknowledgments**

239

240 This work was funded by LOEWE-Centre for Translational Biodiversity Genomics (LOEWE-  
241 TBG). We thank Damian Baranski for help with the DNA isolation and library preparations. Luis J.  
242 Chueca was supported by a Post-doctoral Fellowship awarded by the Department of Education,  
243 Universities and Research of the Basque Government (Ref.: POS\_2018\_1\_0012).

244

245 **Author contributions**

246

247 M.P. and L.J.C. conceived the idea. M.P. collected the specimens. L.J.C. designed and performed  
248 the bioinformatic analyses with support of T.S. L.J.C. prepared the manuscript, and all authors  
249 edited and approved the final version.

250

251 **References**

252

253 Andrews, S., 2010 FastQC: a quality control tool for high throughput sequence data.

254 Bolger, A. M., M. Lohse, and B. Usadel, 2014 Trimmomatic: A flexible trimmer for Illumina  
255 sequence data. *Bioinformatics* 30: 2114–2120.

256 Brozzo, A., J. Harl, W. De Mattia, D. Teixeira, F. Walther *et al.*, 2020 Molecular phylogeny and trait  
257 evolution of Madeiran land snails: radiation of the Geomitridae (Stylommatophora: Helicoidea:  
258 Geomitridae). *Cladistics* 36: 594–616.

259 Cai, H., Q. Li, X. Fang, J. Li, N. E. Curtis *et al.*, 2019 A draft genome assembly of the solar-  
260 powered sea slug *Elysia chlorotica*. *Sci. Data* 6: 190022.

261 Campbell, M. A., B. J. Haas, J. P. Hamilton, S. M. Mount, and C. R. Robin, 2006 Comprehensive  
262 analysis of alternative splicing in rice and comparative analyses with *Arabidopsis*. *BMC*  
263 *Genomics* 7: 1–17.

264 Chueca, L. J., B. J. Gómez-Moliner, M. Forés, and M. J. Madeira, 2017 Biogeography and radiation  
265 of the land snail genus *Xerocrassa* (Geomitridae) in the Balearic Islands. *J. Biogeogr.* 44: 760–  
266 772.

267 Chueca, L. J., B. J. Gómez-Moliner, M. J. Madeira, and M. Pfenninger, 2018 Molecular phylogeny  
268 of *Candidula* (Geomitridae) land snails inferred from mitochondrial and nuclear markers  
269 reveals the polyphyly of the genus. *Mol. Phylogenet. Evol.* 118:.

270 Giusti, F., and G. Manganelli, 1987 Notulae malacologicae, XXXVI. On some Hygromiidae  
271 (Gastropoda: Helicoidea) living in Sardinia and in Corsica.(Studies on the Sardinian and  
272 Corsican malacofauna VI). *Boll. Malacol.* 23: 123–206.

273 Guo, Y., Y. Zhang, Q. Liu, Y. Huang, G. Mao *et al.*, 2019 A chromosomal-level genome assembly  
274 for the giant African snail *Achatina fulica*. *Gigascience* 8: 1–8.

275 Gurevich, A., V. Saveliev, N. Vyahhi, and G. Tesler, 2013 QUAST: Quality assessment tool for  
276 genome assemblies. *Bioinformatics* 29: 1072–1075.

277 Haas, B. J., A. Papanicolaou, M. Yassour, M. Grabherr, P. D. Blood *et al.*, 2013 De novo transcript  
278 sequence reconstruction from RNA-seq using the Trinity platform for reference generation and



- 279 analysis. Nat. Protoc. 8: 1494–1512.
- 280 Haas, B. J., S. L. Salzberg, W. Zhu, M. Pertea, J. E. Allen *et al.*, 2008 Automated eukaryotic gene  
281 structure annotation using EVIDENCEModeler and the Program to Assemble Spliced  
282 Alignments. Genome Biol. 9: 1–22.
- 283 Haponski, A. E., T. Lee, and D. Ó Foighil, 2017 Moorean and Tahitian *Partula* tree snail survival  
284 after a mass extinction: New genomic insights using museum specimens. Mol. Phylogenet.  
285 Evol. 106: 151–157.
- 286 Hare, E. E., and J. S. Johnston, 2012 Chapter 1 of Propidium Iodide-Stained Nuclei. Methods 772:  
287 3–12.
- 288 Kang, S. W., B. B. Patnaik, H. J. Hwang, S. Y. Park, J. M. Chung *et al.*, 2016 Transcriptome  
289 sequencing and de novo characterization of Korean endemic land snail, *Koreanohadra*  
290 *kurodana* for functional transcripts and SSR markers. Mol. Genet. Genomics 291: 1999–2014.
- 291 Keilwagen, J., F. Hartung, M. Paulini, S. O. Twardziok, and J. Grau, 2018 Combining RNA-seq  
292 data and homology-based gene prediction for plants, animals and fungi. BMC Bioinformatics  
293 19:.
- 294 Keilwagen, J., M. Wenk, J. L. Erickson, M. H. Schattat, J. Grau *et al.*, 2016 Using intron position  
295 conservation for homology-based gene prediction. Nucleic Acids Res. 44:.
- 296 Kim, D., B. Langmead, and S. L. Salzberg, 2015 HISAT: A fast spliced aligner with low memory  
297 requirements. Nat. Methods 12: 357–360.
- 298 Korábek, O., A. Petrusek, and M. Rovatsos, 2019 The complete mitogenome of *Helix pomatia* and  
299 the basal phylogeny of Helicinae (Gastropoda, Stylommatophora, Helicidae). Zookeys 2019:  
300 19–30.
- 301 Korf, I., 2004 Gene finding in novel genomes. BMC Bioinformatics 5: 1–9.
- 302 Kumar, S., M. Jones, G. Koutsovoulos, M. Clarke, and M. Blaxter, 2013 Blobology: exploring raw  
303 genome data for contaminants, symbionts and parasites using taxon-annotated GC-coverage  
304 plots. Front. Genet. 4: 1–12.
- 305 Laetsch, D. R., and M. L. Blaxter, 2017 BlobTools : Interrogation of genome assemblies [ version  
306 1 ; peer review : 2 approved with reservations ]. F1000Research 6: 1287.
- 307 Liu, C., Y. Zhang, Y. Ren, H. Wang, S. Li *et al.*, 2018 The genome of the golden apple snail  
308 *Pomacea canaliculata* provides insight into stress tolerance and invasive adaptation.  
309 Gigascience 7: 1–13.
- 310 Pfenninger, M., and F. Magnin, 2001 Phenotypic evolution and hidden speciation in *Candidula*  
311 *unifasciata* ssp. (Helicellinae, Gastropoda) inferred by 16S variation and quantitative shell  
312 traits. Mol. Ecol. 10: 2541–2554.
- 313 Quevillon, E., V. Silventoinen, S. Pillai, N. Harte, N. Mulder *et al.*, 2005 InterProScan: Protein

- 314 domains identifier. *Nucleic Acids Res.* 33: 116–120.
- 315 Razkin, O., G. Sonet, K. Breugelmans, M. J. Madeira, B. J. Gómez-Moliner *et al.*, 2016 Species  
316 limits, interspecific hybridization and phylogeny in the cryptic land snail complex *Pyramidula*:  
317 The power of RADseq data. *Mol. Phylogenet. Evol.* 101: 267–278.
- 318 Roach, M. J., S. A. Schmidt, and A. R. Borneman, 2018 Purge Haplotigs: allelic contig  
319 reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* 19: 460.
- 320 Romero, P. E., A. M. Weigand, and M. Pfenninger, 2016 Positive selection on panpulmonate  
321 mitogenomes provide new clues on adaptations to terrestrial life. *BMC Evol. Biol.* 16: 1–13.
- 322 Ruan, J., and H. Li, 2020 Fast and accurate long-read assembly with wtdbg2. *Nat. Methods* 17:  
323 155–158.
- 324 Saenko, S. V, D. S. J. Groenenberg, A. Davison, and M. Schilthuizen, 2021 The draft genome  
325 sequence of the grove snail *Cepaea nemoralis*. *G3 Genes, Genomes, Genet.* jkaa071:.
- 326 Sauer, J., and B. Hausdorf, 2010 Reconstructing the evolutionary history of the radiation of the land  
327 snail genus *Xerocrassa* on Crete based on mitochondrial sequences and AFLP markers. *BMC*  
328 *Evol. Biol.* 10: 299.
- 329 Schell, T., B. Feldmeyer, H. Schmidt, B. Greshake, O. Tills *et al.*, 2017 An Annotated Draft  
330 Genome for *Radix auricularia* (Gastropoda, Mollusca). *Genome Biol. Evol.* 9: 585–592.
- 331 Schilthuizen, M., and V. Kellermann, 2014 Contemporary climate change and terrestrial  
332 invertebrates: evolutionary versus plastic changes. *Evol. Appl.* 7: 56–67.
- 333 Simão, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, and E. M. Zdobnov, 2015 BUSCO:  
334 Assessing genome assembly and annotation completeness with single-copy orthologs.  
335 *Bioinformatics* 31: 3210–3212.
- 336 Slater, G. S. C., and E. Birney, 2005 Automated generation of heuristics for biological sequence  
337 comparison. *BMC Bioinformatics* 6: 1–11.
- 338 Smit, A., and R. Hubley, 2008 RepeatModeler Open-1.0. Available fom [http://www. repeatmasker.](http://www.repeatmasker.org)  
339 [org](http://www.repeatmasker.org).
- 340 Solem, A., 1984 A world model of land snail diversity and abundance, pp. 6–22 in *World-wide*  
341 *Snails, Biogeographical studies on non-marine mollusca*, Brill and Backhuys, Leiden.
- 342 Song, L., D. S. Shankar, and L. Florea, 2016 Rascaf: Improving Genome Assembly with RNA  
343 Sequencing Data. *Plant Genome* 9: 1–12.
- 344 Stankowski, S., 2013 Ecological speciation in an island snail: Evidence for the parallel evolution of  
345 a novel ecotype and maintenance by ecologically dependent postzygotic isolation. *Mol. Ecol.*  
346 *22: 2726–2741.*
- 347 Sun, J., C. Chen, N. Miyamoto, R. Li, J. D. Sigwart *et al.*, 2020 The Scaly-foot Snail genome and  
348 implications for the origins of biomineralised armour. *Nat. Commun.* 11:1-12.

- 349 Vaser, R., I. Sović, N. Nagarajan, and M. Šikić, 2017 Fast and accurate de novo genome assembly  
 350 from long uncorrected reads. *Genome Res.* 27: 737–746.
- 351 Walker, B. J., T. Abeel, T. Shea, M. Priest, A. Abouelliel *et al.*, 2014 Pilon: An Integrated Tool for  
 352 Comprehensive Microbial Variant Detection and Genome Assembly Improvement (J. Wang,  
 353 Ed.). *PLoS One* 9: e112963.
- 354 Warren, R. L., C. Yang, B. P. Vandervalk, B. Behsaz, A. Lagman *et al.*, 2015 LINKS: Scalable,  
 355 alignment-free scaffolding of draft genomes with long reads. *Gigascience* 4:.
- 356 Xu, G. C., T. J. Xu, R. Zhu, Y. Zhang, S. Q. Li *et al.*, 2018 LR-Gapcloser: A tiling path-based gap  
 357 closer that uses long reads to complete genome assembly. *Gigascience* 8: 1–14.

358

### 359 **Figures and Tables**

360

361 **Table 1.** Genome assembly and annotation statistics for *C. unifasciata* and comparison with other  
 362 land snails genomes.

Statistic	<i>Candidula unifasciata</i>	<i>Cepaea nemoralis</i>	<i>Achatina fulica</i>
<b>Total sequence length</b>	1,286,461,591	3,490,924,950	1,850,322,141
<b>No. of contigs</b>	11,756	28,698	8,122
<b>Contig N50</b>	246,413	330,079	721,038
<b>Contig L50</b>	1,602	3,071	697
<b>No. of scaffolds</b>	8,586	28,537	921
<b>scaffolds &gt; 10000 bp</b>	7,180	26,580	189
<b>Scaffold N50</b>	246,413	333,110	59,589,303
<b>Scaffold L50</b>	940	3,035	13
<b>GC content (%)</b>	40.69	41.25	38.77
<b>BUSCO complete (genome)</b>	92.4% (S:85.3%; D:7.1%)	87.2% (S:74.3%; D:12.9%)*	91.5% (S:84.6%; D:6.9%)
<b>BUSCO fragmented (genome)</b>	1.6%	3.8%*	2.5%
<b>BUSCO missing (genome)</b>	6.0%	9.0%*	6.0%
<b>BUSCO complete (annotation)</b>	94.5% (S:86.0%; D:8.5%)	na	95.6% (S:86.8%; D:8.8%)

<b>BUSCO fragmented (annotation)</b>	2.6%	na	1.9%
<b>BUSCO missing (annotation)</b>	2.9%	na	2.5%
<b>BUSCO complete (transcriptome)</b>	94.7% (S:52.6%; D:42.1%)		
<b>BUSCO fragmented (transcriptome)</b>	3.8%		
<b>BUSCO missing (transcriptome)</b>	1.5%		

363 \*against metazoa\_odb10 dataset (n=954)

364

365 **Table 2.** Repeat statistics. *De novo* and homology based repeat annotations as reported by

366 RepeatMasker and RepeatModeler for *C. unifasciata* and comparison with *Cepaea nemoralis*.

367 Families of repeats included here are long interspersed nuclear elements (LINEs), short interspersed

368 nuclear elements (SINEs), long tandem repeats (LTR), DNA transposons (DNA), unclassified

369 (unknown) repeat families, small RNA repeats (SmRNA), and others (consisting of small, but

370 classified repeat groups). The total is the total percentage of base pairs made up of repeats in each

371 genome assembly, respectively.

Assembly	LINE	SINE	LTR	DNA	Unclassified	SmRNA	Others	Total (%)
<i>Candidula unifasciata</i>	1,253,318	427,509	11,975	298,828	1,334,718	413,197	708,740	61.1
<i>Cepaea nemoralis</i>	2,820,864	342,120	209,476	443,363	4,400,828	444,489	1,267,814	77.0

372

373 **Table 3.** Functional annotation of the predicted protein-coding genes for *C. unifasciata* genome.

		<i>C. unifasciata</i>
<b>Number</b>		
	<b>Gene</b>	22,464
	<b>mRNA</b>	22,464
	<b>Exon</b>	147,783
	<b>CDS</b>	147,783
<b>Mean</b>		
	<b>mRNAs/gene</b>	1
	<b>CDSs/mRNA</b>	6.58
<b>Median length</b>		

	<b>Gene</b>	11,931
	<b>mRNA</b>	11,931
	<b>Exon</b>	129
	<b>Intron</b>	2,025
	<b>CDS</b>	129
<b>Total space (Mb)</b>		
	<b>Gene</b>	379,573,459
	<b>CDS</b>	26,582,739
<b>Single</b>		
	<b>CDS mRNA</b>	3,562
<b>InterproScan</b>		21,231 (94.51%)
<b>GO</b>		13,221 (62.27%)
<b>Reactome</b>		5,069 (22.56%)
<b>SwissProt</b>		16,809 (74.83%)

374

375 **Table 4.** Software employed in this work, their package version and source availability.

<b>Name</b>	<b>Version</b>	<b>Url</b>
Flye	2.6	<a href="https://github.com/fenderglass/Flye">https://github.com/fenderglass/Flye</a>
wtdbg2	2.5	<a href="https://github.com/ruanjue/wtdbg2">https://github.com/ruanjue/wtdbg2</a>
Canu	1.9	<a href="https://github.com/marbl/canu">https://github.com/marbl/canu</a>
Racon	1.4.3	<a href="https://github.com/isovic/racon">https://github.com/isovic/racon</a>
Pilon	1.23	<a href="https://github.com/broadinstitute/pilon">https://github.com/broadinstitute/pilon</a>
Quast	5.0.2	<a href="https://github.com/ablab/quast">https://github.com/ablab/quast</a>
BUSCO	3.0.2	<a href="https://busco.ezlab.org/">https://busco.ezlab.org/</a>
BlobTools	1.1.1	<a href="https://github.com/DRL/blobtools">https://github.com/DRL/blobtools</a>
LINKS	1.8.7	<a href="https://github.com/bcgsc/LINKS">https://github.com/bcgsc/LINKS</a>
Rascaf	1.0.2	<a href="https://github.com/mourisl/Rascaf">https://github.com/mourisl/Rascaf</a>
Long-Read Gapcloser	1.0	<a href="https://github.com/CAFS-bioinformatics/LR_Gapcloser">https://github.com/CAFS-bioinformatics/LR_Gapcloser</a>
FastQC	0.11.9	<a href="https://www.bioinformatics.babraham.ac.uk/projects/fastqc/">https://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>
Trimmomatic	0.39	<a href="http://www.usadellab.org/cms/?page=trimmomatic">http://www.usadellab.org/cms/?page=trimmomatic</a>
MultiQC	1.9	<a href="https://multiqc.info/">https://multiqc.info/</a>
GenomeScope	1.0	<a href="http://qb.cshl.edu/genomescope/">http://qb.cshl.edu/genomescope/</a>
Trinity	2.9.1	<a href="https://github.com/trinityrnaseq/trinityrnaseq/wiki">https://github.com/trinityrnaseq/trinityrnaseq/wiki</a>
GeMoMa	1.6.4	<a href="http://www.jstacs.de/index.php/GeMoMa">http://www.jstacs.de/index.php/GeMoMa</a>
MMseqs2		<a href="https://github.com/soedinglab/MMseqs2">https://github.com/soedinglab/MMseqs2</a>
TransDecoder	5.5.0	<a href="https://github.com/TransDecoder">https://github.com/TransDecoder</a>
SNAP	2006-07-28	
EXONERATE	2.2.0	<a href="https://www.ebi.ac.uk/about/vertebrate-genomics/software/exonerate-manual">https://www.ebi.ac.uk/about/vertebrate-genomics/software/exonerate-manual</a>
PASA	2.4.1	<a href="https://github.com/PASApipeline/PASApipeline">https://github.com/PASApipeline/PASApipeline</a>

EVidenceMolder	1.1.1	<a href="https://evidencemodeler.github.io">https://evidencemodeler.github.io</a>
guppy	4.0.11	<a href="https://nanoporetech.com/nanopore-sequencing-data-analysis">https://nanoporetech.com/nanopore-sequencing-data-analysis</a>
Nanoplots	1.28.1	<a href="https://github.com/wdecoster/NanoPlot">https://github.com/wdecoster/NanoPlot</a>
Nanofilt	2.6.0	<a href="https://github.com/wdecoster/nanofilt">https://github.com/wdecoster/nanofilt</a>
backmap.pl	0.3	<a href="https://github.com/schelllt/backmap">https://github.com/schelllt/backmap</a>
SAMtools	1.10	<a href="https://github.com/samtools/samtools">https://github.com/samtools/samtools</a>
BWA	0.7.17	<a href="https://github.com/lh3/bwa">https://github.com/lh3/bwa</a>
minimap2	2.17	<a href="https://github.com/lh3/minimap2">https://github.com/lh3/minimap2</a>
Qualimap	2.2.1	<a href="http://qualimap.conesalab.org/">http://qualimap.conesalab.org/</a>
bedtools	2.28.0	<a href="https://bedtools.readthedocs.io/en/latest/">https://bedtools.readthedocs.io/en/latest/</a>
Rscript	3.6.3	<a href="https://www.r-project.org/">https://www.r-project.org/</a>
RepeatModeler	2.0	<a href="http://www.repeatmasker.org/RepeatModeler/">http://www.repeatmasker.org/RepeatModeler/</a>
RepeatMasker	4.1.0	<a href="http://www.repeatmasker.org/">http://www.repeatmasker.org/</a>
HISAT2	2.1.0	<a href="http://daehwankimlab.github.io/hisat2/">http://daehwankimlab.github.io/hisat2/</a>

376

377 **Table S1.** Comparison between draft genomes assemblies obtained by the different tools.

	<b>Platanus 1.2.4</b>	<b>SOAPdenovo2</b>	<b>MaSuRCA 3.3.3</b>	<b>wtdbg 2.5</b>	<b>Flye 2.6</b>
<b>Total sequence length</b>	1,272,133,741	981,942,849	1,609,244,920	1,390,813,883	1,505,080,485
<b>No. of contigs</b>	879,520	848,801	23,717	16,291	23,552
<b>contigs &gt; 10000 bp</b>	1,947	292	18,340	12,288	18,725
<b>Largest contig</b>	29,328	20,266	1,951,786	1,581,874	1,244,054
<b>Contig N50</b>	1,818	1,308	172,678	222,260	117,519
<b>Contig L50</b>	194,857	215,144	2,483	1,789	3,522
<b>GC content (%)</b>	40.86	40.70	40.87	40.66	40.69
<b>BUSCO complete</b>			91.6% (S:78.0%; D:13.6%)	91.7% (S:83.9%; D:7.8%)	91.0% (S:79.4%; D:11.6%)
<b>BUSCO fragmented</b>			2.0%	2.4%	2.8%
<b>BUSCO missing</b>			6.4%	5.9%	6.2%

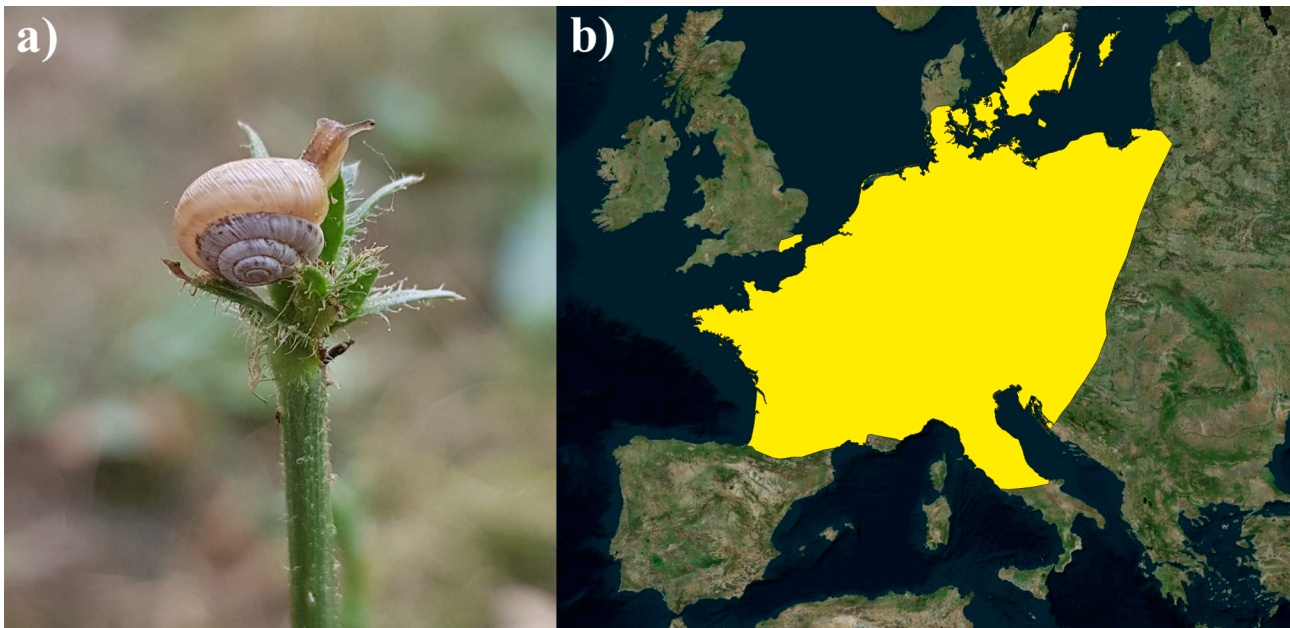
378

379 **Figures:**

380

381 **Figure 1. a)** Picture of an adult specimen of *Candidula unifasciata*, copyright © Luis J. Chueca. **b)**

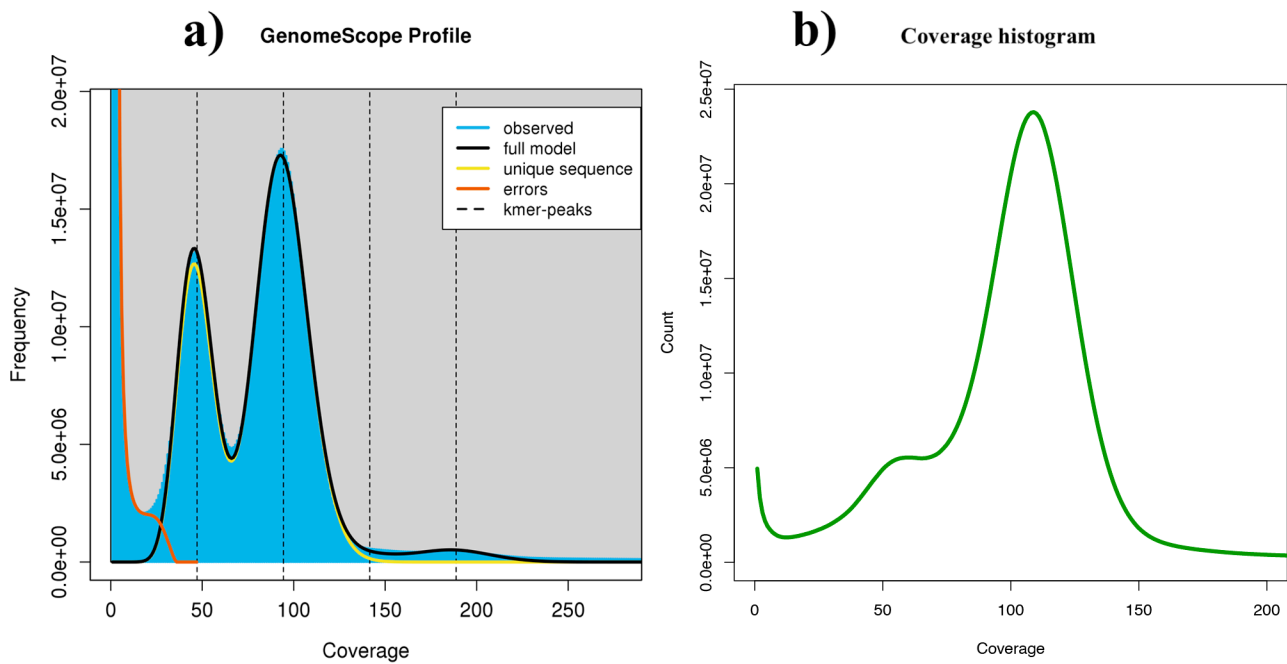
382 Distribution range of *C. unifasciata* in Europe.



383

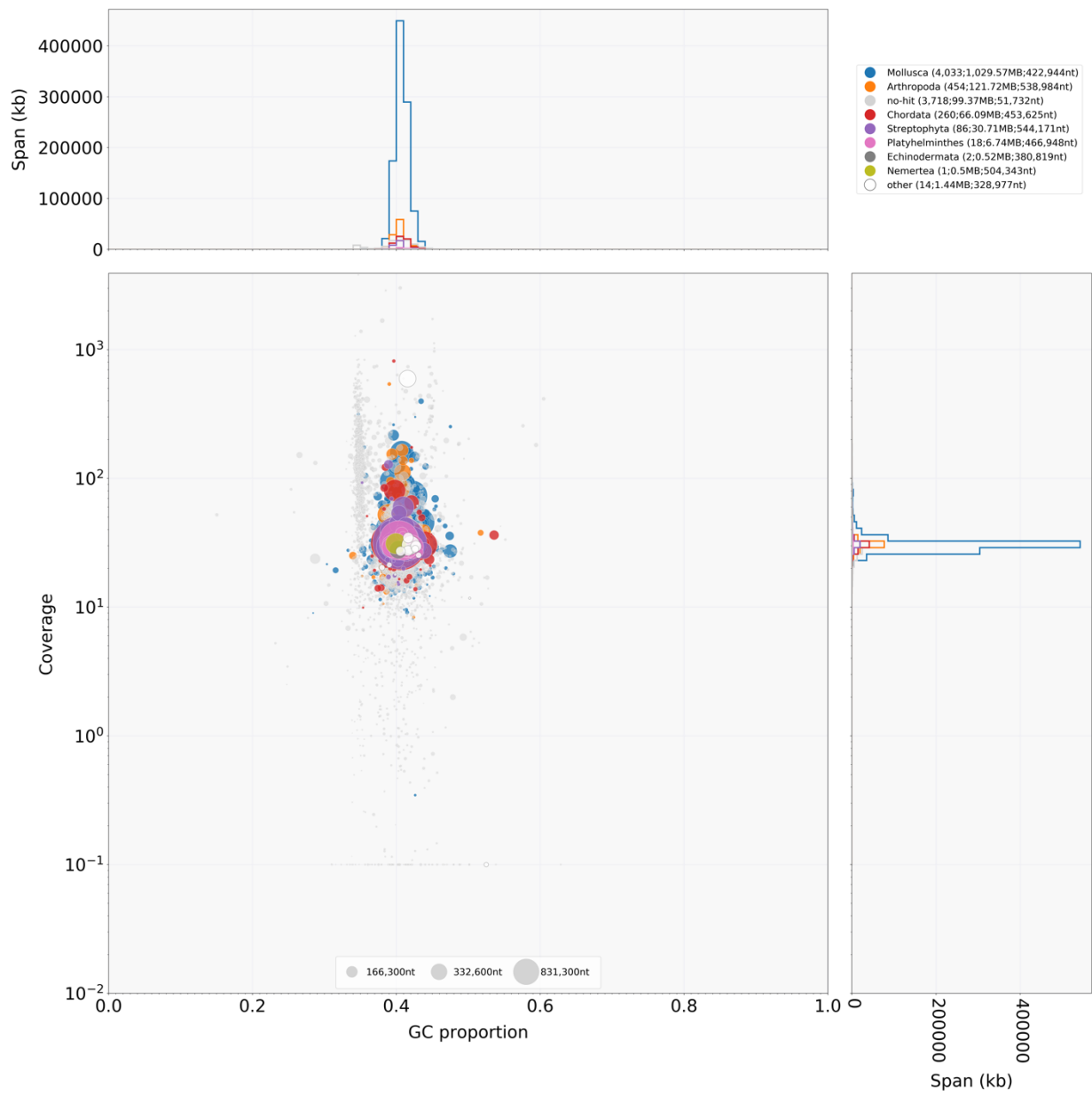
384

385 **Figure 2. a)** GenomeScope k-mer profile plot for *Candidula unifasciata* genome based on 21-mers  
386 in Illumina reads. **b)** Coverage histogram for the final assembly based on the Illumina reads.



387

388 **Figure 3.** Blob plot showing read depth of coverage, GC content and size of each scaffold. Size of  
389 the blobs correspond to size of the scaffold and color corresponds to taxonomic assignment of  
390 BLAST.



391