

**Genomic resources for the North American water vole (*Microtus richardsoni*) and the  
montane vole (*Microtus montanus*)**

Drew J. Duckett<sup>1\*</sup>, Jack Sullivan<sup>2</sup>, Stacy Pirro<sup>3</sup>, Bryan C. Carstens<sup>1</sup>

<sup>1</sup>Department of Evolution, Ecology, and Organismal Biology. The Ohio State University. 1315  
Kinnear Rd., Columbus OH, 43212

<sup>2</sup>Department of Biological Sciences, Box 443051, University of Idaho, Moscow ID, 83844-3051

1 <sup>3</sup>Iridian Genomes, Inc., 6213 Swords Way, Bethesda MD 20817

\*email: [duckettdj@gmail.com](mailto:duckettdj@gmail.com)

## 2 **Abstract**

3 *Background:* Voles of the genus *Microtus* are important research organisms, yet genomic  
4 resources in the genus are lacking. Providing such resources would benefit future studies of  
5 immunology, phylogeography, cryptic diversity, and more. *Findings:* We sequenced and  
6 assembled nuclear genomes from two subspecies of water vole (*Microtus richardsoni*) and from  
7 the montane vole (*Microtus montanus*). The water vole genomes were sequenced with Illumina  
8 and 10X Chromium plus Illumina sequencing, resulting in assemblies with ~1,600,000 and  
9 ~30,000 scaffolds respectively. The montane vole was assembled into ~13,000 scaffolds using  
10 Illumina sequencing also. In addition to the nuclear assemblies, mitochondrial genome  
11 assemblies were also performed for both species. We conducted a structural and functional  
12 annotation for the best water vole nuclear genome, which resulted in ~24,500 annotated genes,  
13 with 83% of these receiving functional annotations. Finally, we find that assembly quality  
14 statistics for our nuclear assemblies fall within the range of genomes previously published in the  
15 genus *Microtus*, making the water vole and montane vole genomes useful additions to currently  
16 available genomic resources.

17  
18 *Keywords:* genome assembly; genome annotation; mitochondrial genome; 10X Chromium;  
19 Illumina sequencing

20

21

22

23

24

## 25 **Context**

26 The genus *Microtus* consists of 62 species of voles distributed throughout North America,  
27 Europe, and Asia [1]. *Microtus* is believed to have experienced rapid speciation and  
28 diversification, with all speciation events occurring within the past four million years [2, 3], and  
29 it has been suggested that some nominal species, such as *M. pennsylvanicus*, contain cryptic  
30 diversity [4]. *Microtus* has been an important model system across multiple biological  
31 disciplines, including studies of adaptation (e.g., [5]), infectious disease (e.g., [6]), parental care  
32 (e.g., [7]), and population dynamics (reviewed in [8]). The rapid radiation of *Microtus* voles has  
33 hindered systematic classification, leading to multiple taxonomic revisions and conflicting  
34 phylogenetic analyses [1, 9, 10]. Consequently, both species boundaries and relationships among  
35 species are difficult to infer. Genomic resources within *Microtus* will help resolve these  
36 questions, and resources have steadily increased in recent years. Currently, four *Microtus* species  
37 have assembled genomes on GenBank, two European species (*M. agrestis* and *M. arvalis*) and  
38 two North American species (*M. ochrogaster*; [11], and *M. oeconomus*). The present study  
39 provides resources for two additional species: *M. richardsoni* and *M. montanus*.

40 The North American water vole (*M. richardsoni*) is adapted to a semiaquatic lifestyle, relying  
41 on alpine and sub-alpine streams for creating burrows and escaping predators [12]. Like other  
42 semiaquatic mammals (e.g., otters), it is likely that adaptations to this lifestyle have been driven  
43 by natural selection [13-15]. Water voles are among the largest species of *Microtus* and are  
44 known for making runways of stamped-down vegetation along streams through frequent  
45 movement [12, 16]. Unlike most other vole species, *M. richardsoni* does not appear to  
46 experience regular population boom and bust cycles, although population size in the species may  
47 be correlated to levels of precipitation [17]. Despite being listed as *Least Concern* by the IUCN

48 Redlist [18], the species is listed as *Critically Imperiled* by the Wyoming Natural Diversity  
49 Database due to its specific habitat requirements, which can be substantially degraded by  
50 livestock grazing [19]. *Microtus richardsoni* occupies a large, disjunct distribution in the Pacific  
51 Northwest of North America, with habitat in the Cascades Mountains and the Rocky Mountains,  
52 spanning from southern Canada into central Utah. Four subspecies are currently recognized: *M.*  
53 *r. arvicoloides* in the Cascades Mountains, *M. r. richardsoni* in the Canadian Rocky Mountains,  
54 *M. r. macropus* in the central Rocky Mountains and Wyoming, and *M. r. myllodontus* in Utah.  
55 Due to the subspecific classifications and the disjunct range of the species, *M. richardsoni* has  
56 been included in multiple studies of phylogeography in the Pacific Northwest [20-22]. These  
57 studies were based solely on mitochondrial DNA, and the results of analyses that investigated  
58 species limits and demographic history were limited to inferences that can be derived from a  
59 single gene tree. Genomic resources for *M. richardsoni* will provide a rich source of data to  
60 address these knowledge gaps.

61 The montane vole (*M. montanus*) is partially sympatric with *M. richardsoni* and can be found  
62 throughout most of the water vole's range with the exception of the Canadian Rockies. However,  
63 *M. montanus* can be found farther south and east including areas of California, Nevada,  
64 Colorado, Arizona and New Mexico [23]. The species has been divided into fifteen subspecies,  
65 including *M. m. canescens* in the Cascades Mountains., *M. m. nasus* in the central Rocky  
66 Mountains, and *M. m. amosus* in northern Utah. Notably, *M. montanus* does not exhibit a break  
67 in its range in the Columbia Basin, likely because it is not restricted to riparian areas like *M.*  
68 *richardsoni*. The species as a whole is listed as *Least Concern* by the IUCN Redlist, but *M. m.*  
69 *arizonicus* has been listed as endangered by the New Mexico State Game Commission  
70 Regulation [23], and *M. m. ricularis* has been noted as being of concern due to a small range and

71 declining population size [24]. Genomic resources in *M. montanus* will provide a wealth of data  
72 to assess subspecies boundaries, quantify gene flow among subspecies, and aid in conservation  
73 efforts of threatened subspecies.

74 The present study provides two nuclear and one mitochondrial genome assembly for *M.*  
75 *richardsoni* along with single nuclear and mitochondrial genome assemblies for *M. montanus*.  
76 Furthermore, a structural and functional annotation are performed with one of the *M. richardsoni*  
77 genomes to aid in future studies of adaptation. Genome-level comparisons are made between the  
78 new genome assemblies and other *Microtus* genome assemblies to examine differences in  
79 assembly quality and repeat content.

80

## 81 **Sequencing and Nuclear Genome Assembly**

82 Frozen tissue from a single *M. r. arvicoloides* individual collected from the southern Cascades  
83 Mountain range (JMS\_292; 44.016667N, -121.750000E; [20]) was sent to Hudson Alpha  
84 (Huntsville, AL) for high molecular weight DNA extraction and 10X Chromium library  
85 preparation [25]. In the 10X method, each extracted DNA fragment receives a different barcode  
86 before the fragment is sheared for library preparation. After sequencing, these barcodes are used  
87 to connect sequencing reads for a more contiguous assembly. After sequencing with a single run  
88 on an Illumina HiSeqX, the resulting 150 base pair (bp) paired-end reads were input into  
89 Supernova for *de novo* genome assembly with --maxreads=all [26].

90 Additional tissue was obtained from a single *M. r. macropus* individual collected from the  
91 northern Rocky Mountains (JMG\_88; 46.333333N, -114.633333E; [20]). DNA was extracted  
92 using a Qiagen DNeasy Blood and Tissue Kit, and the DNA was sent for library preparation and  
93 sequencing by Iridian Genomes, Inc (Bethesda, MD). 150bp paired-end reads were sequenced on

94 two runs of an Illumina HiSeqX. Genome assembly was performed using two different deBruijn  
95 graph-based programs, SOAPdenovo and Discover de novo [27, 28]. For SOAPdenovo, quality  
96 trimming was performed using fastQC and Trimmomatic with settings ILLUMINACLIP:  
97 2:30:10, LEADING:3, TRAILING:3, SLIDINGWINDOW:4:15, and MINLEN:36 [29, 30].  
98 SOAPdenovo assemblies were performed with settings max\_rd\_len=150, avg\_ins=300,  
99 reverse\_seq=0, asm\_flags=3, rd\_len\_cutoff=150, rank=1, pair\_num\_cutoff=3, and map\_len=32.  
100 SOAPdenovo was run with kmer values of 63, 89, 95, and 101 based on analysis of optimal kmer  
101 values in kmerGenie [31]. Raw reads were used as input for *de novo* genome assembly with  
102 Discover as recommended in the program documentation.

103 In an attempt to provide the most contiguous assembly for *M. richardsoni*, a hybrid assembly  
104 was performed using the ARCS+LINKS pipeline [32, 33]. The ARCS+LINKS pipeline uses  
105 barcoding information from the 10X Chromium reads to scaffold the contigs from a separate  
106 genome assembly. Barcoded reads from *M. r. arvicoloides* were mapped to the *M. r. macropus*  
107 Discover assembly with bwa mem [34] before converting the mapped reads to BAM format and  
108 sorting with SAMTools [35]. ARCS and LINKS were then run with settings `-s 98 -c 5 -l 0 -z`  
109 `500 -d 0 -r 0.05 -m 50-10000 -e 30000` and `-d 4000 -k 20 -l 5 -t 2 -a 0.3 -o 0 -a 0.3 -z 500`  
110 respectively.

111 As part of a separate project, a single *M. montanus* individual from Utah  
112 (UMNH:Mamm:30891; 38.19381N, -111.5824E) was misidentified as *M. richardsoni*. DNA was  
113 extracted from the sample using a Qiagen DNeasy Blood and Tissue Kit before being sent to the  
114 University of California Davis Genome Center for library preparation and sequencing. Paired-  
115 end 150bp sequences were collected with a single shared run on an Illumina NovaSeq. Species  
116 identity was confirmed using the Barcode of Life Database (BOLD; [36]). Reads were checked

117 and trimmed for quality with fastQC and Trimmomatic as above before mapping reads to the  
118 mitochondrial cytochrome oxidase I (COI) sequence of *M. r. macropus* [37] using bwa mem.  
119 The resulting mapped reads were converted to BAM format, sorted, and indexed with  
120 SAMTools. PCR duplicates were identified and removed with Picard [38], resulting reads were  
121 piled with SAMTools mpileup using base and mapping quality scores of 30, consensus  
122 sequences were generated with bcftools [39], and consensus sequences were converted to fastq  
123 format using vcftutils with a minimum depth filter of 5 and maximum depth filter of 10000 [35].  
124 The resulting sequence was input into BOLD. Due to the low sequencing coverage, *de novo*  
125 genome assembly was not appropriate for *M. montanus*. To provide a preliminary genome  
126 sequence, a reference-guided genome assembly was performed with RaGOO [40]. Raw reads  
127 were input into Discover to generate an initial genome assembly, misassembly correction was  
128 performed with RaGOO using reads trimmed with the same settings as the *M. r. macropus* reads,  
129 and RaGOO was then used to scaffold the Discover contigs onto the *M. r. arvicoloides* assembly,  
130 which is more closely related to *M. montanus* than the other available *Microtus* genome  
131 assemblies [3]. Since *M. montanus* has less than half the chromosomes of *M. richardsoni* ( $2n =$   
132 22-24 in *montanus* versus 56 in *richardsoni* [41]), the possibility of structural errors in the *M.*  
133 *montanus* assembly was examined by calculating the percentage of reads that mapped back to the  
134 assembly using bwa mem and bamtools [42].

135 The final assemblies were submitted to GenBank [43], where screening was performed to  
136 identify any contamination, and contaminated scaffolds were removed. All assemblies were  
137 evaluated with QUAST [44], bbmap [45], custom Python scripts  
138 ([https://github.com/djlduckett/Genome\\_Resources/](https://github.com/djlduckett/Genome_Resources/)), and BUSCO using the Euarchontoglires  
139 reference set [46]. After comparing assembly statistics from the different assemblies of *M. r.*

140 *macropus*, the Discover assembly was selected as best because it had less fragmentation, higher  
141 N50 and L50, and a higher BUSCO score than the SOAPdenovo assemblies (Table 1). Genome  
142 sequencing of *M. r. arvicoloides* produced over 800 million (M) reads and 47x genome  
143 sequencing coverage. The final genome assembly consisted of ~32 thousand (K) scaffolds with  
144 an N50 of 2.3 megabase pairs (Mb), 1.3% missing data (N), and a BUSCO score of 85.8%.  
145 Supernova estimated the length of the genome assembled to be ~2.4Gb and the total genome size  
146 to be ~2.6 gigabase pairs (Gb). *Microtus richardsoni macropus* sequencing produced over 600M  
147 reads and 35x coverage. Genome assembly with Discover resulted in ~1.6M scaffolds with an  
148 N50 of 16 kilobase pairs (Kb), 0.06% Ns, and a BUSCO score of 54.5%. Given that there are  
149 many programs to perform *de novo* genome assembly from short reads, it is possible that another  
150 program would have produced a more contiguous *M. r. macropus* assembly, but previous studies  
151 have shown Discover performs well compared to other programs [47, 48]. The hybrid assembly  
152 produced with the ARCS+LINKS pipeline had ~1.6M scaffolds, an N50 of 38Kb, 0.09% Ns, and  
153 a BUSCO score of 59.8%. Because of the poor quality of the hybrid assembly, it was not used  
154 for further analyses, and the *M. richardsoni* subspecies assemblies were kept separate. It seems  
155 likely that the high fragmentation of the hybrid assembly is due to the fragmentation of the  
156 Discover input assembly. Published results with this hybrid pipeline often include a much higher  
157 sequencing coverage of the input contigs to produce a better starting point for the pipeline.  
158 Therefore, additional Illumina sequencing with *M. r. macropus* in the future could substantially  
159 improve the hybrid assembly. 108M reads (13x coverage) were used to produce the preliminary  
160 *M. montanus* genome, resulting in ~13K scaffolds, an N50 of ~3.1Mb, 8.8% Ns, and a BUSCO  
161 score of 82.6%. Additionally, 89.3% of reads mapped back to the *M. montanus* assembly.

162

163



164

**Table 1**

	Discover	SOAPdenovo	SOAPdenovo	SOAPdenovo	SOAPdenovo
Kmer	NA	63	89	95	101
Length	2.54Gb	2.72Gb	2.88Gb	2.89Gb	3.21Gb
Scaffolds	1.6M	4.1M	4.0M	4.1M	6.7M
Max Scaffold	264Kb	186Kb	146Kb	174Kb	139Kb
N50	16.1Kb	4.5Kb	3.4Kb	3.4Kb	1.5Kb
L50	35.7K	117K	156K	163K	371K
BUSCO	54.5%	38.1%	37.1%	35.9%	25.9%
% N	0.06	1.45	0.99	0.94	0.90
% GC	42.13	41.92	41.91	41.92	41.98

165 Comparison of genome assembly strategies for *M. r. macropus*. NA: not applicable.

166

### 167 **Mitochondrial Genomes**

168 The complete mitochondrial genomes of *M. r. arvicoloides* and *M. montanus* were assembled  
169 using the genomic sequencing reads. The mitochondrial genomes were assembled by both  
170 mapping reads to a reference mitochondrial genome and using the reference-guided assembly  
171 program Novoplasty [49]. For the mapping assembly, reads were mapped to the *M. r. macropus*  
172 mitochondrial genome, using the same steps as the *M. montanus* BOLD analysis. The  
173 mitochondrial assemblies were 16,285bp and 16,268bp in length with an average depth of  
174 coverage of 7886x and 6805x for *M. r. arvicoloides* and *M. montanus* respectively. Reference  
175 guided mitochondrial assemblies with Novoplasty used the *M. r. macropus* mitochondrial  
176 genome as the reference along with settings *Genome Range=12000-22000, K-mer=33, Read*

177 *Length=150, and Insert size=400*. Because the *M. r. arvicoloides* dataset contained many reads,  
178 25% of reads were subsampled to use for assembly, as suggested in the program documentation.  
179 The assemblies for *M. r. arvicoloides* and *M. montanus* were 16,298bp and 16,319bp in length  
180 with average depths of coverage of 5131x and 14,713x respectively. To compare mitochondrial  
181 assemblies between methods, the assemblies were aligned using the MUSCLE plugin in  
182 Geneious v. R9 with eight iterations and an open gap score of -1 [50, 51]. This comparison  
183 showed the Novoplasty assemblies contained multiple insertions compared to the mapped  
184 assemblies and the reference mitochondrial genome. These insertions were up to 13bp long in  
185 multiple genes, including *trnT*, *trnK*, and *ATP8*. Comparison to other *Microtus* mitochondrial  
186 genomes (*M. ochrogaster*; NC\_027945.1 and *M. fortis*; NC\_015243.1) showed that the  
187 Novoplasty assemblies were the only mitochondrial assemblies to exhibit these insertions.  
188 Therefore, the mapping assemblies were used for further analyses. The mapping assemblies for  
189 both species included ambiguous bases, which were much more frequent for *M. montanus* than  
190 *M. r. arvicoloides*. These may be the result of using the mitochondrial genome of a different  
191 subspecies (for *M. r. arvicoloides*) or species (for *M. montanus*) for mapping the reads.  
192 Additionally, the presence of nuclear DNA of mitochondrial origin (NUMTs; [52, 53]) may have  
193 influenced these results. If mitochondrial segments have been incorporated into the nuclear  
194 genomes and subsequent mutations have occurred, both nuclear and mitochondrial sequences  
195 could be mapped to the same mitochondrial region during assembly and result in the ambiguous  
196 bases observed here. It is likely that NUMTs are present, as they have been documented in other  
197 species of *Microtus* [54-56]. Both mitochondrial genomes were annotated using MITOS [57].  
198 The annotations each consisted of 22 tRNA genes, 2 rRNA genes, and 13 protein coding genes.  
199

## 200 ***Microtus* Genome Assembly Comparison**

201 The available *Microtus* genome assemblies, *M. agrestis* (GCA\_902806755.1), *M. arvalis*  
202 (GCA\_007455615.1), *M. ochrogaster* (GCA\_000317375.1), and *M. oeconomus*  
203 (GCA\_007455595.1), were downloaded from GenBank. Assembly summary statistics were  
204 calculated using QUAST, bbmap, and custom Python scripts  
205 ([https://github.com/djlduckett/Genome\\_Resources/](https://github.com/djlduckett/Genome_Resources/)). To compare repeat content among all  
206 genomes, including the three produced by the current study, repeats were first identified *de novo*  
207 using RepeatModeler [58]. RepeatMasker was then used to further identify repeats using a  
208 combined repeat library that included the repeats identified from RepeatModeler and those from  
209 the RepeatMasker *Rodentia* database [59]. The percentage of the genome consisting of each type  
210 of repeat element was extracted from the RepeatMasker log file for each genome assembly.

211 All genome assemblies used some form of Illumina sequencing (Table 2), although assembly  
212 continuity varied greatly among assemblies from 1366 scaffolds in *M. agrestis* to 1.6 M scaffolds  
213 in *M. r. macropus*. Genome coverage was similarly varied, from 13x in *M. montanus* to 35x in  
214 *M. r. macropus* to 77x in *M. arvalis* and *M. oeconomus*. The percent of repetitive regions ranged  
215 from 31.7% in *M. montanus* to 44.1% in *M. arvalis* (Figure 1), and repeat content did not appear  
216 to be associated with phylogenetic relatedness as repeats between the two subspecies of *M.*  
217 *richardsoni* were not more similar to each other than to other *Microtus* species. However, it is  
218 possible that the repeat content is affected by the continuity of the genome assemblies, and  
219 further research is needed to confirm this relationship.

Table 2

Species	<i>M. agrestis</i>	<i>M. arvalis</i>	<i>M. montanus</i> *	<i>M. ochrogaster</i>	<i>M. oeconomus</i>	<i>M. r. arvicoloides</i> *	<i>M. r. macropus</i> *
<b>Distribution</b>	Europe	Europe	North America	North America	North America	North America	North America
<b>Year</b>	2020	2019	2020	2012	2019	2020	2020
<b>Accession (GCA_)</b>	902806775.1	7455615.1	xxxxxxxxxxx	317375.1	7455595.1	xxxxxxxxxxx	xxxxxxxxxxx
<b>Sequencing</b>	10X Chromium + Illumina	Illumina	Illumina	Illumina	Illumina	10X Chromium + Illumina	Illumina
<b>Assembler</b>	Supernova	Discover	RaGOO	ALLPATHS	Discover	Supernova	Discover
<b>Length</b>	2.03Gb	2.62Gb	2.34Gb	2.29Gb	2.31Gb	2.36Gb	2.54Gb
<b>Coverage</b>	50	77	13	94	77	47	35
<b># Scaffolds</b>	1,366	1,081,432	12,962	6,341	562,436	31,632	1,648,927
<b>Longest Scaffold</b>	56.96Mb	0.80Mb	748.72Mb	126.73Mb	0.93Mb	16.00Mb	0.26Mb
<b>N50</b>	13.35Mb	0.53Mb	3.08Mb	61.81Mb	0.11Mb	2.30Mb	0.02Mb
<b>L50</b>	45	11,870	91	14	5,556	278	35,660
<b>%N</b>	2.87	0.07	8.81	8	0.12	1.29	0.06
<b>%GC</b>	42.33	41.71	42.38	42.25	42.18	42.21	42.13

221 Genome assembly comparison among *Microtus* species. Assemblies with a \* were produced by the present study. Note: in-depth methods for *M. agrestis* are not  
 222 available, and it is possible that the assembly includes additional sequencing and/or methods.

## 223 **Genome Annotation**

224 The *M. r. arvicoloides* genome assembly was annotated with the MAKER pipeline [60],  
225 loosely following the tutorial provided by Daren Card (<https://gist.github.com/darencard>  
226 [/bb1001ac1532dd4225b030cf0cd61ce2](https://gist.github.com/darencard/bb1001ac1532dd4225b030cf0cd61ce2)). Briefly, the pipeline consists of masking repeats  
227 followed by multiple rounds of annotation with both evidence-based and ab-initio gene models.  
228 Repeats were identified as described above. Complex repeats were then extracted from  
229 RepeatMasker results using grep with keywords “Satellite” and “rich”. Within Maker, the  
230 model\_org argument was set to “simple” so Maker would soft mask simple repeats, and the  
231 RepeatMasker results were provided to hard mask complex repeats. Evidence-based gene  
232 discovery used protein and mRNA sequences from the previous genome annotation of *M.*  
233 *ochrogaster* (GCF\_000317375.1) as well as an additional RNASeq assembly from *M.*  
234 *pennsylvanicus* (GSM3499528; [61]). Hidden Markov models (HMMs) for ab-initio gene  
235 prediction were trained using both SNAP and Augustus [62, 63]. With SNAP, gene models  
236 identified by MAKER were filtered using an Annotation Edit Distance (AED) of 0.5 and an  
237 amino acid length of 50. After validating these models with SNAP’s Fathom utility, removing  
238 likely errors, and including 1000bp surrounding each training sequence, the training sequences  
239 were passed to the hmm-assembler script. For Augustus, training sequences plus 1000bp on each  
240 side were obtained from the first round of MAKER mRNA annotations. Augustus was used to  
241 train the HMM using the --long option in BUSCO and the Euarchontoglires reference set.  
242 MAKER was then run again with the previously annotated gene models and the HMM models  
243 from SNAP and Augustus. After the initial MAKER run, two cycles of ab-initio gene prediction  
244 and annotation with MAKER were performed. To prevent overfitting, results were compared  
245 after each round of MAKER. Because the increase in AED score was minimal between the first

246 and second rounds of ab-initio gene prediction, further analysis was conducted on the results  
247 after the first round only. This round annotated ~24K genes with a mean gene length of 7445bp  
248 (Table 3), which is within the range found in previous studies of *M. ochrogaster* (22,427 genes;  
249 GCF\_000317375.1) and *Arvicola amphibious* (25,136 genes; GCF\_903992535.1). Of these  
250 annotations all occurred on scaffolds greater than 1Kb in length and 97% occurred on scaffolds  
251 greater than 10Kb in length.

252 Functional annotation of the *M. r. arvicoloides* genome was performed using GOfeat, an  
253 online functional annotation tool that uses multiple protein databases including UniProt,  
254 InterPro, and Pfam [64-67]. An input file for GOfeat was generated by supplying the genome  
255 assembly FASTA file and the MAKER General Feature Format (GFF3) file to the Python  
256 package gffread [68]. GOfeat annotated 83.49% of genes. Biological Processes accounted for  
257 42.46% of annotations, Cellular Components accounted for 30.29%, and Molecular Functions  
258 comprised 27.25%. The most frequent gene ontology (GO) terms were *positive regulation of*  
259 *transcription by RNA polymerase II*, *negative regulation of transcription by RNA polymerase II*,  
260 and *DNA-templated regulation of transcription* for Biological Processes, *cytoplasm* and *plasma*  
261 *membrane* for Cellular Components, and *metal ion binding* and *calcium ion binding* for  
262 Molecular Functions.

263

264

265

266

267

268

269

270

**Table 3**

	<b>Before Gene Modeling</b>	<b>Gene Modeling Round 1</b>	<b>Gene Modeling Round 2</b>
<b>Genes</b>	20,945	24,548	23,811
<b>Exons</b>	139,845	192,974	179,225
<b>mRNA</b>	20,945	24,548	23,811
<b>tRNA</b>	-	24,504	24,539
<b>5' UTR</b>	-	1,229	1,180
<b>3' UTR</b>	-	503	642
<b>Mean Gene Length</b>	-	7,445	7,132
<b>AED &lt; 0.50</b>	0.993	0.881	0.888
<b>AED &lt; 0.25</b>	0.672	0.543	0.520
<b>BUSCO (Complete)</b>	-	67.7%	70.5%

271 Structural annotation summary after each round of MAKER. UTR: untranslated region; AED: annotation edit

272 distance. Values with dashes were not analyzed prior to gene modeling with SNAP and Augustus.

273

## 274 **Conclusion**

275 The current study details the assembly and annotation of three nuclear and two mitochondrial  
 276 genomes. Compared to previously published nuclear genomes, the *M. r. arvicoloides* and *M.*  
 277 *montanus* genomes are of high quality as evidenced by the low number of scaffolds, high  
 278 N50/L50 values, and high BUSCO scores. While not as complete as the other *Microtus* genomes,  
 279 the nuclear genome of *M. r. macropus* will still be useful for mapping low coverage reads or  
 280 reduced representation sequencing data. Furthermore, the mitochondrial genomes contributed  
 281 here add to a growing number for the genus *Microtus* and reinforce earlier suggestions that high-  
 282 quality mitochondrial genomes can be obtained as byproducts of nuclear sequencing (e.g., [69,  
 283 70]). Overall, the data presented serve as an example that even though they do not include

284 chromosomal information, high-quality draft genomes can be produced from widely available  
285 and very cost-effective methods like the 10X Chromium protocol. These references can aid a  
286 variety of studies including those examining genus and species adaptation [71, 72],  
287 phylogenetics [10], phylogeography [22, 73], and disease dynamics [6, 74]. However, some  
288 activities, like exploring changes to chromosome structure, will not be possible due to the  
289 fragmentation and lack of chromosomal mapping for these assemblies. Finally, the *M. r.*  
290 *macropus* and *M. montanus* sequencing data and preliminary assemblies will serve as the  
291 building blocks of more accurate reference genomes in the future.

292

### 293 **Availability of Supporting Data and Materials**

294 Raw sequences, nuclear assemblies, and mitochondrial assemblies are available from GenBank  
295 under BioProjects PRJNA673719, PRJNA509068, and PRJNA673873 for *M. r. arvicoloides*, *M.*  
296 *r. macropus*, and *M. montanus* respectively. The custom python script used to calculate genome  
297 assembly summary information is available on GitHub  
298 ([https://github.com/djlduckett/Genome\\_Resources/](https://github.com/djlduckett/Genome_Resources/)). Full BUSCO tables, structural annotation  
299 gff files, functional annotation tables, and repeat libraries are available in the GigaScience data  
300 repository (<http://gigadb.org/>).

301

302

303

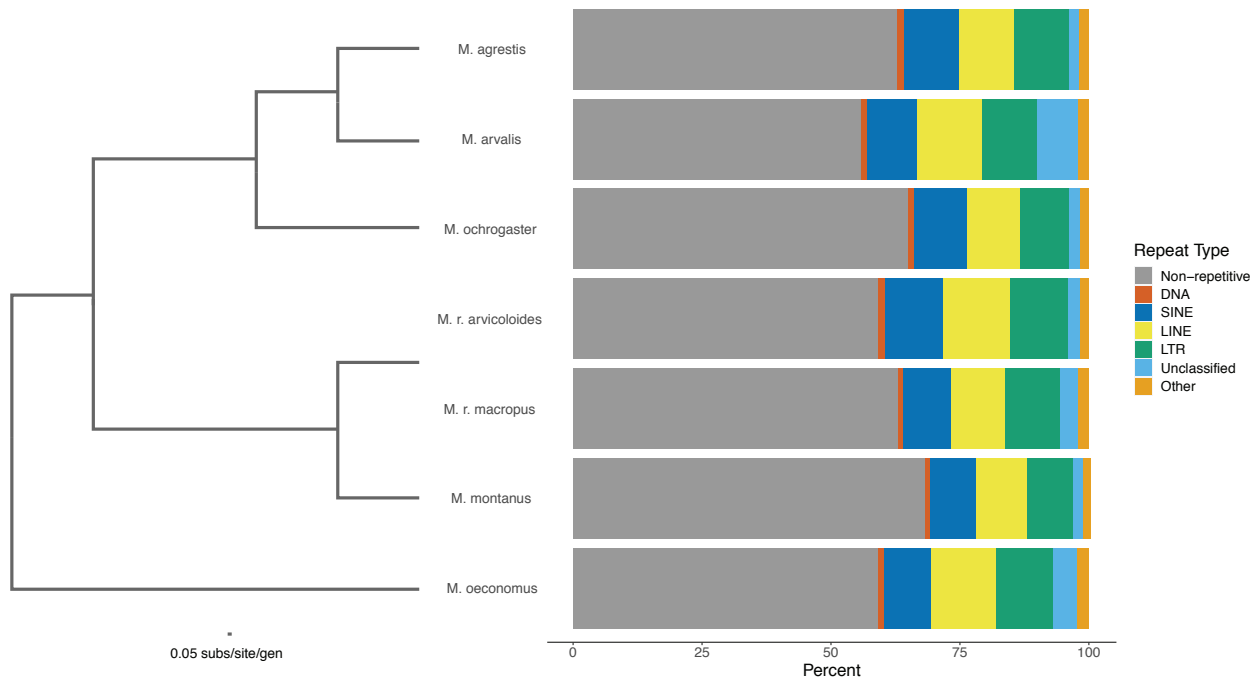
304

305

306



307 **Figures:**



322 **Competing Interests**

323 SP is the director of Iridian Genomes, Inc.

324

325 **Funding**

326 Sequencing was funded by Iridian Genomes, Inc., as well as the National Science Foundation

327 (DEB-1457519). Salary support for DD was provided by The Ohio State University and the

328 National Science Foundation (DBI-1945347).

329

330 **Author Contributions**

331 DD, JS, and BC conceived the study. JS, SP, and BC provided funding for sequencing. DD

332 performed DNA extractions, assembled genomes, and annotated genomes with input from SP.

333 DD and BC wrote the manuscript with input from JS and SP. DD and SP submitted the resources

334 to GenBank.

335

336 **Acknowledgements**

337 We thank Jeffrey Good and Eric Rickart/Utah Museum of Natural History for tissue samples,

338 Michael Broe for advice with genome assembly and annotation, and the Ohio Supercomputer

339 Center (OSC) for computational resources.

340

341

342

343

344

345 **References**

- 346 1. Wilson DE, Reeder DM, editors. Mammal species of the world: a taxonomic and geographic  
347 reference. JHU Press; 2005.
- 348 2. Reig OA. Karyotypic repatterning as one triggering factor in cases of explosive speciation. In:  
349 Fontdevila A, editor. Evolutionary biology of transient unstable populations. Springer, Berlin,  
350 Heidelberg; 1989 p. 246-289.
- 351 3. Stepan SJ, Schenk JJ. Murid rodent phylogenetics: 900-species tree reveals increasing  
352 diversification rates. PloS one. 2017;12:8.
- 353 4. Jackson DJ, Cook JA. A precarious future for distinctive peripheral populations of meadow voles  
354 (*Microtus pennsylvanicus*). Journal of Mammalogy. 2020;101:1:36-51.
- 355 5. Monarca RI, Speakman JR, da Luz Mathias M. Energetics and thermal adaptation in semifossorial  
356 pine-voles *Microtus lusitanicus* and *Microtus duodecimcostatus*. Journal of Comparative  
357 Physiology B. 2019;189:2:309-18.
- 358 6. Wanelik KM, Begon M, Birtles RJ, Bradley JE, Friberg IM, Jackson JA, Taylor CH, Thomason AG,  
359 Turner AK, Paterson S. A candidate tolerance gene identified in a natural population of field  
360 voles (*Microtus agrestis*). Molecular ecology. 2018;27:4:1044-52.
- 361 7. Seelke AM, Perkeybile AM, Grunewald R, Bales KL, Krubitzer LA. Individual differences in cortical  
362 connections of somatosensory cortex are associated with parental rearing style in prairie voles  
363 (*Microtus ochrogaster*). Journal of Comparative Neurology. 2016;524:3:564-77.
- 364 8. Oli MK. Population cycles in voles and lemmings: state of the science and future directions.  
365 Mammal Review. 2019;49:3:226-39.
- 366 9. Bailey V. Revision of American voles of the genus *Microtus*. North American Fauna. 1900;17:1-  
367 88.

- 368 10. Barbosa S, Paupério J, Pavlova SV, Alves PC, Searle JB. The *Microtus* voles: Resolving the  
369 phylogeny of one of the most speciose mammalian genera using genomics. *Molecular*  
370 *phylogenetics and evolution*. 2018;125:85-92.
- 371 11. McGraw LA, Davis JK, Lowman JJ, ten Hallers BF, Koriabine M, Young LJ, De Jong PJ, Rudd MK,  
372 Thomas JW. Development of genomic resources for the prairie vole (*Microtus ochrogaster*):  
373 construction of a BAC library and vole-mouse comparative cytogenetic map. *BMC genomics*.  
374 2010;11:1:1-8.
- 375 12. Klaus M, Beauvais GP. Water Vole (*Microtus richardsoni*): A Technical Conservation Assessment.  
376 Prepared for USDA Forest Service, Rocky Mountain Region, Species Conservation Project. 2004.
- 377 13. Stein BR. Bone density and adaptation in semiaquatic mammals. *Journal of mammalogy*.  
378 1989;70:3:467-76.
- 379 14. Dunstone N, Gorman ML, editors. Behaviour and ecology of riparian mammals. Cambridge  
380 University Press; 2007.
- 381 15. Beichman AC, Koepfli KP, Li G, Murphy W, Dobrynin P, et al. Aquatic adaptation and depleted  
382 diversity: a deep dive into the genomes of the sea otter and giant otter. *Molecular biology and*  
383 *evolution*. 2019;36:12:2631-55.
- 384 16. Ludwig DR. *Microtus richardsoni*. *Mammalian Species*. 1984;223:1-6.
- 385 17. Klaus M, Moore RE, Vyse E. Impact of precipitation and grazing on the water vole in the  
386 Beartooth Mountains of Montana and Wyoming, USA. *Arctic, Antarctic, and Alpine Research*.  
387 1999;31:3:278-82.
- 388 18. Cassola, F. *Microtus richardsoni*. The IUCN Red List of Threatened Species 2016.
- 389 19. Wyoming Natural Diversity Database. University of Wyoming. 2020.  
390 [https://wyndd.org/species\\_list/](https://wyndd.org/species_list/).

- 391 20. Cartens BC, Brunsfeld SJ, Demboski JR, Good JM, Sullivan J. Investigating the evolutionary history  
392 of the Pacific Northwest mesic forest ecosystem: hypothesis testing within a comparative  
393 phylogeographic framework. *Evolution*. 2005;59:8:1639-52.
- 394 21. Carstens BC, Richards CL. Integrating coalescent and ecological niche modeling in comparative  
395 phylogeography. *Evolution: International Journal of Organic Evolution*. 2007;61:6:1439-54.
- 396 22. Espíndola A, Ruffley M, Smith ML, Carstens BC, Tank DC, Sullivan J. Identifying cryptic diversity  
397 with predictive phylogeography. *Proceedings of the Royal Society B: Biological Sciences*.  
398 2016;283:1841:20161529.
- 399 23. Sera WE, Early CN. *Microtus montanus*. *Mammalian Species*. 2003;716:1-10.
- 400 24. Cassola, F. *Microtus montanus*. The IUCN Red List of Threatened Species 2016.
- 401 25. Zheng GX, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, Kyriazopoulou-  
402 Panagiotopoulou S, Masquelier DA, Merrill L, Terry JM, Mudivarti PA. Haplotyping germline and  
403 cancer genomes with high-throughput linked-read sequencing. *Nature biotechnology*.  
404 2016;34:3:303-11.
- 405 26. Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. Direct determination of diploid genome  
406 sequences. *Genome research*. 2017;27:5:757-67.
- 407 27. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, He G, Chen Y, Pan Q, Liu Y, Tang J. SOAPdenovo2: an  
408 empirically improved memory-efficient short-read de novo assembler. *Gigascience*.  
409 2012;1:1:2047-17.
- 410 28. Weisenfeld NI, Yin S, Sharpe T, Lau B, Hegarty R, Holmes L, Sogoloff B, Tabbaa D, Williams L, Russ  
411 C, Nusbaum C. Comprehensive variation discovery in single human genomes. *Nature genetics*.  
412 2014;46:12:1350.
- 413 29. Andrews S. FastQC: a quality control tool for high throughput sequence data. 2010.  
414 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/> .

- 415 30. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.  
416 Bioinformatics. 2014;30:15:2114-20.
- 417 31. Chikhi R, Medvedev P. Informed and automated k-mer size selection for genome assembly.  
418 Bioinformatics. 2014;30:1:31-7.
- 419 32. Yeo S, Coombe L, Warren RL, Chu J, Birol I. ARCS: scaffolding genome drafts with linked reads.  
420 Bioinformatics. 2018;34:5:725-31.
- 421 33. Warren RL, Yang C, Vandervalk BP, Behsaz B, Lagman A, Jones SJ, Birol I. LINKS: Scalable,  
422 alignment-free scaffolding of draft genomes with long reads. GigaScience. 2015;4:1:s13742-015.
- 423 34. Li H, Durbin R. Fast and accurate long-read alignment with Burrows–Wheeler transform.  
424 Bioinformatics. 2010;26:5:589-95.
- 425 35. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. The  
426 sequence alignment/map format and SAMtools. Bioinformatics. 2009;25:16:2078-9.
- 427 36. Ratnasingham S, Hebert PD. BOLD: The Barcode of Life Data System ([http://www. barcodinglife.](http://www.barcodinglife.org)  
428 [org](http://www.barcodinglife.org)). Molecular ecology notes. 2007;7:3:355-64.
- 429 37. Alqahtani F, Duckett D, Pirro S, Mandoiu II. Complete mitochondrial genome of the water vole,  
430 *Microtus richardsoni* (Cricetidae, Rodentia). Mitochondrial DNA Part B. 2020;5:3:2498-9. Broad  
431 Institute. Picard tools. 2016. <http://github.com/broadinstitute/picard>.
- 432 38. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and  
433 population genetical parameter estimation from sequencing data. Bioinformatics.  
434 2011;27:21:2987-93.
- 435 39. Alonge M, Soyk S, Ramakrishnan S, Wang X, Goodwin S, Sedlazeck FJ, Lippman ZB, Schatz MC.  
436 RaGOO: fast and accurate reference-guided scaffolding of draft genomes. Genome biology.  
437 2019;20:1:1-7.

- 438 40. Modi WS. Phylogenetic analyses of chromosomal banding patterns among the Nearctic  
439 Arvicolidae (Mammalia: Rodentia). *Systematic Zoology*. 1987;36:2:109-36.
- 440 41. Barnett DW, Garrison EK, Quinlan AR, Strömberg MP, Marth GT. BamTools: a C++ API and toolkit  
441 for analyzing and managing BAM files. *Bioinformatics*. 2011;27:12:1691-2.
- 442 42. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank.  
443 *Nucleic acids research*. 2012;41:D1:D36-42.
- 444 43. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUASt: quality assessment tool for genome  
445 assemblies. *Bioinformatics*. 2013;29:8:1072-5.
- 446 44. Bushnell B. BBMap: a fast, accurate, splice-aware aligner. Lawrence Berkeley National  
447 Lab.(LBNL), Berkeley, CA (United States); 2014. <https://jgi.doe.gov/data-and-tools/bbtools/>.
- 448 45. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome  
449 assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.  
450 2015;31:19:3210-2.
- 451 46. Mikheenko A, Prjibelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome assembly  
452 evaluation with QUASt-LG. *Bioinformatics*. 2018;34:13:i142-50.
- 453 47. Fernandez-Silva I, Henderson JB, Rocha LA, Simison WB. Whole-genome assembly of the coral  
454 reef Pearlscale Pygmy Angelfish (*Centropyge vrolikii*). *Scientific reports*. 2018;8:1:1-1.
- 455 48. Dierckxsens N, Mardulyn P, Smits G. NOVOPlasty: de novo assembly of organelle genomes from  
456 whole genome data. *Nucleic acids research*. 2017;45:4:e18-.
- 457 49. Kears e M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A,  
458 Markowitz S, Duran C, Thierer T. Geneious Basic: an integrated and extendable desktop  
459 software platform for the organization and analysis of sequence data. *Bioinformatics*.  
460 2012;28:12:1647-9.

- 461 50. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
462 Nucleic acids research. 2004;32:5:1792-7.
- 463 51. Lopez JV, Yuhki N, Masuda R, Modi W, O'Brien SJ. Numt, a recent transfer and tandem  
464 amplification of mitochondrial DNA to the nuclear genome of the domestic cat. Journal of  
465 molecular evolution. 1994;39:2:174-90.
- 466 52. Hazkani-Covo E, Zeller RM, Martin W. Molecular poltergeists: mitochondrial DNA copies  
467 (nuMountains) in sequenced nuclear genomes. PLoS Genet. 2010;6:2:e1000834.
- 468 53. Triant DA, DeWoody JA. Extensive mitochondrial DNA transfer in a rapidly evolving rodent has  
469 been mediated by independent insertion events and by duplications. Gene. 2007;401:1-2:61-70.
- 470 54. Triant DA, DeWoody JA. Molecular analyses of mitochondrial pseudogenes within the nuclear  
471 genome of arvicoline rodents. Genetica. 2008;132:1:21-33.
- 472 55. Triant DA, DeWoody JA. Demography and phylogenetic utility of numt pseudogenes in the  
473 Southern Red-Backed Vole (*Myodes gapperi*). Journal of mammalogy. 2009;90:3:561-70.
- 474 56. Bernt M, Donath A, Jühling F, Externbrink F, Florentz C, Fritzsich G, Pütz J, Middendorf M, Stadler  
475 PF. MITOS: improved de novo metazoan mitochondrial genome annotation. Molecular  
476 phylogenetics and evolution. 2013;69:2:313-9.
- 477 57. Smit AF, Hubley R. RepeatModeler Open-1.0. 2008. <http://www.repeatmasker.org> .
- 478 58. Smit AF, Hubley R, Green P. RepeatMasker Open-4.0. 2013. <http://www.repeatmasker.org> .
- 479 59. Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B, Holt C, Alvarado AS, Yandell M. MAKER:  
480 an easy-to-use annotation pipeline designed for emerging model organism genomes. Genome  
481 research. 2008;18:1:188-96.
- 482 60. Young RL, Ferkin MH, Ockendon-Powell NF, Orr VN, Phelps SM, Pogány Á, Richards-Zawacki CL,  
483 Summers K, Székely T, Trainor BC, Urrutia AO. Conserved transcriptomic profiles underpin



- 484 monogamy across vertebrates. *Proceedings of the National Academy of Sciences*.  
485 2019;116:4:1331-6.
- 486 61. Korf I. Gene finding in novel genomes. *BMC bioinformatics*. 2004;5:1:59.  
487 62. Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron submodel.  
488 *Bioinformatics*. 2003;19:ii215-25.
- 489 63. Araujo FA, Barh D, Silva A, Guimarães L, Ramos RT. GO FEAT: a rapid web-based functional  
490 annotation tool for genomic and transcriptomic data. *Scientific reports*. 2018;8:1:1-4.  
491 64. UniProt Consortium. UniProt: a hub for protein information. *Nucleic acids research*.  
492 2015;43:D1:D204-12.
- 493 65. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P, Das U, Daugherty L,  
494 Duquenne L, Finn RD. InterPro: the integrative protein signature database. *Nucleic acids*  
495 *research*. 2009;37:D211-5.
- 496 66. Bateman A, Coin L, Durbin R, Finn RD, Hollich V, Griffiths-Jones S, Khanna A, Marshall M, Moxon  
497 S, Sonnhammer EL, Studholme DJ. The Pfam protein families database. *Nucleic acids research*.  
498 2004;32:D138-41.
- 499 67. Perteza G, Perteza M. GFF Utilities: GffRead and GffCompare. *F1000Research*. 2020;9.
- 500 68. Voigt O, Erpenbeck D, Wörheide G. A fragmented metazoan organellar genome: the two  
501 mitochondrial chromosomes of *Hydra magnipapillata*. *BMC genomics*. 2008;9:1:350.  
502 69. Smith DR. Not seeing the genomes for the DNA. *Briefings in functional genomics*. 2012;11:4:289-  
503 90.
- 504 70. Fink S, Excoffier L, Heckel G. Mitochondrial gene diversity in the common vole *Microtus arvalis*  
505 shaped by historical divergence and local adaptations. *Molecular Ecology*. 2004;13:11:3501-14.

- 506 71. Fischer MC, Foll M, Excoffier L, Heckel G. Enhanced AFLP genome scans detect local adaptation  
507 in high-altitude populations of a small rodent (*Microtus arvalis*). *Molecular Ecology*.  
508 2011;20:7:1450-62.
- 509 72. Frey JK. Genetics of allopatric populations of the montane vole (*Microtus montanus*) and  
510 Mogollon vole (*Microtus mogollonensis*) in the American Southwest. *Western North American*  
511 *Naturalist*. 2009;69:2:215-22.
- 512 73. Tołkacz K, Alsarraf M, Kowalec M, Dwużnik D, Grzybek M, Behnke JM, Bajer A. Bartonella  
513 infections in three species of *Microtus*: prevalence and genetic diversity, vertical transmission  
514 and the effect of concurrent *Babesia microti* infection on its success. *Parasites & vectors*. 2018  
515 Dec 1;11(1):491.
- 516 74. Silvestro D, Michalak I. raxmlGUI: a graphical front-end for RAxML. *Organisms Diversity &*  
517 *Evolution*. 2012;12:4:335-7.
- 518