

1 **Single-library chromosome-scale diploid assemblies of vole genomes resolve a**
2 **species-specific duplication implicated in pair bonding**

3
4 Mohamed Abuelanin¹, Gulhan Kaya¹, Juniper A. Lake², Christine Lambert², Melody V. Wu³, Kristen M.
5 Berendzen^{8,9}, Ksenia Krasheninnikova⁴, Jonathan M.D. Wood⁴, Nancy G. Solomon⁵, Zoe R. Donaldson⁶,
6 Karen L. Bales⁷, Kerstin Howe⁴, Jonas Korlach², Devanand Manoli^{8,9}, Jessica Tollkuhn³, Megan Y.
7 Dennis^{1†}.

8
9 Running title: Single-library vole assemblies

10
11 ¹ Genome Center, MIND Institute, and Department of Biochemistry & Molecular Medicine, University of
12 California, Davis, CA 95616, USA.

13
14 ² Pacific Biosciences, Menlo Park, CA 94025, USA.

15
16 ³ Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA.

17
18 ⁴ Wellcome Sanger Institute, Hinxton CB10 1SA, UK.

19
20 ⁵ Department of Biology, Miami University of Ohio, Miami, OH, USA.

21
22 ⁶ Department of Molecular, Cellular, and Developmental Biology, Department of Psychology and
23 Neuroscience, University of Colorado Boulder, Boulder, CO, USA.

24
25 ⁷ Departments of Psychology; Neurobiology, Physiology, and Behavior; and California National Primate
26 Research Center, University of California, Davis, Davis, CA, USA.

27
28 ⁸ Department of Psychiatry and Behavioral Sciences, ⁹Weill Institute for Neuroscience, University of
29 California, San Francisco, San Francisco, CA, USA.

30
31 †Corresponding author:

32 Megan Y. Dennis, Ph.D.

33 University of California, Davis, School of Medicine

34 One Shields Avenue

35 Genome Center, 4303 GBSF

36 Davis, CA 95616

37 Email: mydennis@ucdavis.edu

38
39 Keywords: Long-read sequencing, Chromosome conformation capture, 3C, HiFi, CiFi, genome assembly,
40 genome scaffolding, prairie vole, meadow vole, mate-pair bonding, vasopressin

41 SUMMARY

42 High-quality reference genomes are essential to effectively characterize genomic drivers of speciation,
43 phenotypic diversity, and disease causality. Larger complex genomes often require integration of long-
44 read DNA sequencing with additional genomic data, such as chromosome conformation capture (Hi-C or
45 CiFi) to generate phased chromosome-scale assemblies, however this requires multiple sequencing
46 platforms (in the case of Hi-C) or the construction of multiple long-read sequencing libraries. Here, we
47 devise a strategy that combines PacBio HiFi and CiFi sequencing in a single library and run to efficiently
48 produce high-quality contiguous chromosome-scale diploid genome assemblies. We apply this approach
49 to liver tissue from single individuals of prairie vole (*Microtus ochrogaster*) and meadow vole (*Microtus*
50 *pennsylvanicus*), generating haplotype-resolved, chromosome-scale 2.3 Gbp genomes with QV~62, and
51 99.3% BUSCO completeness. Comparing the two new genomes identifies complex structural changes
52 impacting *Avpr1a*, previously implicated in pair bonding, including a species-specific duplication missing
53 from the existing prairie vole reference genome. These divergent genomic features offer new avenues of
54 investigation related to behavioral divergence between prairie and meadow voles. This single-library
55 approach facilitates a simplified and more affordable assembly workflow, producing near-complete
56 genomes of diverse species using one sequencing platform.

58 INTRODUCTION

59 High-quality reference genomes are foundational for characterizing the genomic basis of speciation,
60 phenotypic diversity, and disease susceptibility¹. However, repetitive and structurally complex regions,
61 including segmental duplications (SDs) and centromeres, have historically been underrepresented or
62 absent from reference assemblies. These regions are recognized as a primary source of new gene
63 functions for adaptive evolution², yet their accurate assembly remains technically challenging. Recent
64 advances in long-read sequencing, particularly PacBio HiFi technology, have transformed *de novo*
65 genome assembly by producing highly accurate reads exceeding 10 kbp³. When combined with trio-
66 binning approaches that leverage parental short-read data, genome-assembly approaches can generate
67 fully phased, haplotype-resolved assemblies⁴. Achieving chromosome-scale contiguity, however,
68 typically requires chromosome conformation capture (3C) data. Genome-wide 3C (Hi-C) depends on
69 short-read sequencing, which maps poorly to repetitive regions and necessitates multi-platform library
70 preparations^{5,6}. CiFi addresses these limitations by integrating 3C libraries with PacBio HiFi long-read
71 sequencing, validated on assembling the Mediterranean fruit fly genome (~600 Mbp) and demonstrating
72 effective scaffolding and improved signal-to-noise ratio compared with Hi-C⁷.

73
74 Here, we extend this combined HiFi-CiFi assembly approach to a further optimization of the workflow
75 and two mammalian genomes approximately four-fold larger (~2.3 Gbp): prairie vole (*Microtus*
76 *ochrogaster*) and meadow vole (*Microtus pennsylvanicus*). Prairie voles are among the few mammalian
77 species exhibiting social monogamy^{8,9}. They form enduring social attachments to their partners, known
78 as pair bonds¹⁰, and both parents care for offspring. In contrast, closely related meadow voles, which
79 diverged from a common ancestor with prairie voles <4 million years ago (mya)¹¹, do not show social
80 monogamy, and only mothers care for offspring. Thus, the neurogenetic differences between vole species
81 offer a powerful system to examine the mechanistic basis of complex social behaviors. Across taxa,
82 neuropeptide hormones represent a critical substrate in the evolution of social behaviors, modulating the
83 neural circuits and underlying physiology that control social interactions^{12,13}. Seminal comparative

84 behavioral and pharmacologic studies have implicated the nonapeptide hormones oxytocin and
85 vasopressin in prairie vole social attachment¹⁴⁻¹⁷. Oxytocin receptor (*Oxtr*) and vasopressin receptor 1a
86 (*Avpr1a*) exhibit divergent brain expression patterns in prairie voles relative to other vole species,
87 suggesting substantial regulatory divergence at these loci^{18,19}.

88

89 The existing prairie vole reference (MicOch1.0; Broad Institute, 2012) predates modern long-read
90 technologies—generated from a 94× coverage Illumina sequence of a female individual—and lacks
91 chromosome-scale contiguity and haplotype resolution. Notably, a >105 kbp prairie-vole duplication of
92 *Avpr1a*, identified by targeted BAC sequencing²⁰, was not present in MicOch1.0. Subsequent work on
93 the *Avpr1a* paralog has been hampered, including verification of its existence in prairie vole or other
94 related species, due to the complex nature of this locus. A chromosome-scale meadow vole assembly was
95 recently generated by the Vertebrate Genomes Project (VGP mMicPen1¹) but has not been compared to
96 the prairie vole in the context of social behavior genetics. Using a new workflow that combines HiFi and
97 CiFi DNA for a single library preparation and sequencing run per vole species, we produce haplotype-
98 resolved, chromosome-scale references enabling a systematic assessment of species divergence across the
99 *Avpr1a* locus.

100

101 RESULTS

102 HiFi-CiFi single-library sequencing and genome assembly

103 To generate HiFi and CiFi data from a single PacBio sequencing library, we isolated high-molecular-
104 weight genomic DNA (gDNA) and fixed chromatin from liver tissue of an individual male prairie vole
105 (*Microtus ochrogaster*) and male meadow vole (*Microtus pennsylvanicus*), respectively (Figure 1A). HiFi
106 gDNA was sheared to a target average fragment size of 16 kbp, while CiFi DNA was prepared using the
107 restriction enzyme HindIII and standard protocol⁷, yielding average fragment sizes of ~8 kbp (Figure S1).
108 The two preparations were pooled at a 90:10 molar ratio (HiFi: CiFi), constructed into a single library, and
109 sequenced on one Revio SMRT Cell per species, producing 102.7 Gbp (prairie vole) and 98.2 Gbp
110 (meadow vole) of data at a median read quality of QV38.

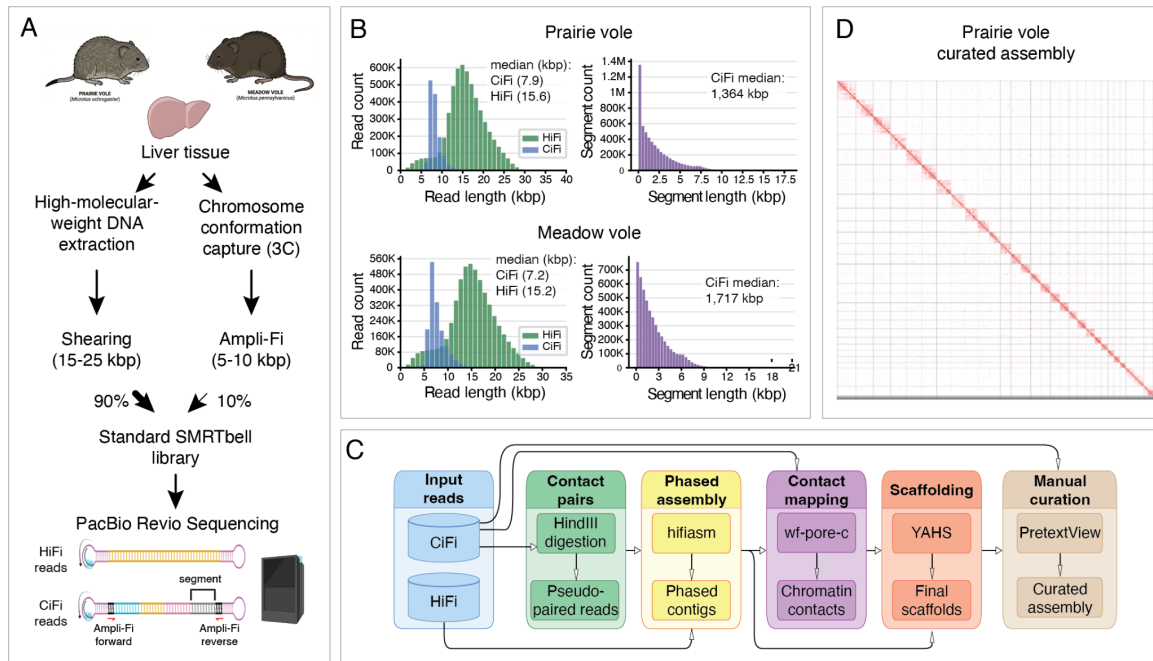
111

112 HiFi (88%) and CiFi (12%) reads were segregated using unique Ampli-Fi adapter sequences flanking
113 CiFi fragments (Methods, Figure 1A). HiFi reads showed median lengths of 15.6 kbp (prairie vole) and
114 15.2 kbp (meadow vole) at 35–40× predicted coverage, while CiFi reads were shorter, as expected, at 7.9
115 kbp (prairie vole) and 7.2 kbp (meadow vole) at 5× predicted coverage (Figure 1B, Table S1). CiFi reads
116 were then converted from multi-segment concatemers into chromatin-contact pairs: 1.37 million CiFi
117 reads yielded 23.4 million contact pairs for prairie vole (median segment size 1.4 kbp), and 1.54 million
118 reads yielded 8.1 million pairs for meadow vole (median segment size 1.7 kbp). The greater contact-pair
119 yield for prairie vole reflects slightly smaller segment sizes and longer CiFi read lengths. Together, these
120 results demonstrate that a combined HiFi and CiFi library produce high-quality data.

121

122 In order to assess CiFi as a single-platform alternative to Hi-C for *de novo* genome assembly, we
123 generated parallel short-read Hi-C data from the same liver tissues using the standard DpnII restriction
124 enzyme (~49 Gbp, 20× for prairie vole; ~73 Gbp, 30× for meadow vole). HiFi reads with either CiFi-
125 derived contact pairs or Hi-C read pairs were used to generate phased diploid contig assemblies with
126 hifiasm⁴, then scaffolded with YAHS²¹ (Figure 1C). CiFi-phased contigs were longer in three of four

127 haplotypes and consolidated into fewer final scaffolds, approaching expected chromosome numbers more
 128 closely than Hi-C (Figures S2 and S3, Table S2). For example, prairie vole Hap1 yielded 63 scaffolds
 129 with CiFi versus 171 with Hi-C (Figure S3). Hi-C also accumulated 1.7–2× more scaffold joins across
 130 YAHS rounds yet produced more final scaffolds, and introduced approximately twice the gap sequence
 131 across all four assemblies (Figure S4), suggesting that many Hi-C joins did not persist in the final
 132 assembly. Where Hi-C showed a higher scaffold N50 in the meadow vole, this was accompanied by more
 133 gaps. Overall, CiFi achieved better or equivalent scaffold consolidation, supporting its utility as a single-
 134 platform approach for chromosome-resolved diploid *de novo* assembly.
 135



136
 137 **Figure 1. Single-library HiFi-CiFi workflow and sequencing stats for workflow prairie and meadow vole. (A)**
 138 Library workflow from liver tissue: high-molecular-weight DNA is used to generate HiFi whole-genome reads, and
 139 a 3C library generates CiFi reads; both are combined in a single SMRTbell library (90:10 HiFi: CiFi). **(B)**
 140 Distribution of read length (HiFi and CiFi) and segment length (CiFi) for prairie (top) and meadow (bottom) voles.
 141 **(C)** Computational pipeline for HiFi-CiFi-based genome assembly and scaffolding. **(D)** CiFi contact map of the
 142 prairie vole curated Hap1 assembly. See Figures S3, S4, and S5 for contact maps of pre- and post-curated assemblies
 143 for both haplotypes and vole species.
 144

145 We next downsampled CiFi reads from 100% (5× coverage) to 1% in each species to determine the
 146 minimum input required to generate contiguous genomes (Figure S5). Assembly sizes remained stable
 147 across all conditions and haplotypes, while scaffold N50 and counts were both sensitive to reduced CiFi
 148 depth. All four haplotypes maintained near-optimal scaffold N50 at 40% input. Together, these results
 149 indicate that approximately 5 Gbp of CiFi data (~2× genome coverage) is sufficient for chromosome-
 150 scale scaffolding in mammalian genomes of this size and complexity and achievable at even lower input
 151 depending on the species⁷.

152
 153 **Final curated vole assemblies**

154 Assemblies produced with HiFi and CiFi required minimal correction during manual curation (see
 155 Methods), including five scaffold breaks and four joins for prairie vole and, similarly, five scaffold breaks
 156 and ten joins for meadow vole (Figures 1D, S2, S3, and S6). The sex chromosomes were identified for

157 each species (represented in Hap1), with the remaining autosomes named according to descending size.
 158 Assessment of curated assemblies show both haplotypes of each diploid vole assembly to be contiguous
 159 (contig N50 > 39 Mbp; scaffold N50 > 89 Mbp), complete (genomic BUSCO \geq 99%), and high quality²²
 160 (QV ~62) (Tables 1 and S3, Figure S7).

161
 162

Table 1. Metric summary of prairie and meadow vole curated genomes

Metric	Prairie vole (this study)		Prairie vole (MicOch1.0)	Meadow vole (this study)		Meadow vole (mMicPen1)	
	Hap1*	Hap2	merged	Hap1*	Hap2	Hap1*	Hap2
Total size	2.53 Gbp	2.30 Gbp	2.29 Gbp	2.36 Gbp	2.15 Gbp	2.37 Gbp	2.16 Gbp
No. contigs	235	156	187,012	200	134	240	196
Contig N50	39.2 Mbp	51.0 Mbp	29.2 kbp	46.5 Mbp	43.6 Mbp	47.2 Mbp	63.6 Mbp
No. scaffolds	97	62	6,335	105	52	180	143
Scaffold N50	89.2 Mbp	89.4 Mbp	61.8 Mbp	125.8 Mbp	115.6 Mbp	125.3 Mbp	118.4 Mbp
auN	98.6 Mbp	97.5 Mbp	56.0 Mbp	119.3 Mbp	120.1 Mbp	118.7 Mbp	120.3 Mbp
Scaffold L50	11	10	14	8	8	8	8
Scaffold L90	21	22	70	18	21	18	21
Gaps (%)	0.073%	0.039%	8.001%	0.006%	0.173%	< 0.001%	< 0.001%
Genome BUSCO	99.7%	97.2%	98.9%	99.7%	97.0%	99.7%	97.1%

163 * Primary assemblies with both sex chromosomes

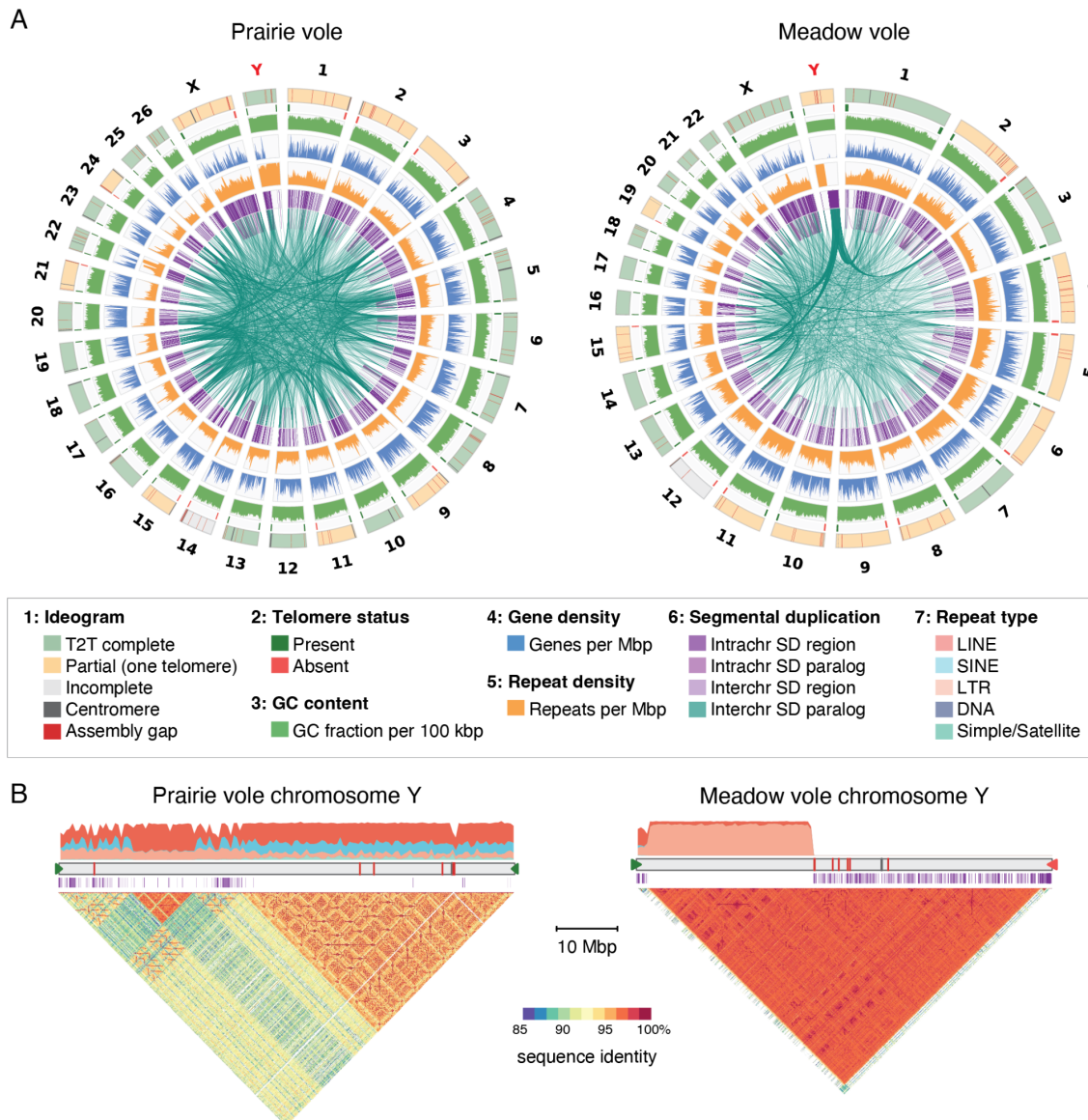
164

165 The prairie vole diploid assembly represents a complete karyotype ($2n = 54$)^{23,24} with 52 autosomes,
 166 chrX, and chrY (Total: 4.78 Gbp, Hap1: 2.50 Gbp). Telomeres were detected on all but three
 167 chromosomes (mean 3.5 gaps/chr). Of these, 31 represent T2T scaffolds and three gapless contigs (overall
 168 63%) (Figure 2, Table S4). In contrast, the current prairie vole reference (MicOch1.0)—a haploid
 169 assembly using short-read sequencing of a female individual—comprises 18 chromosomes (17 autosomes
 170 and chrX), 10 linkage groups, and four unlocalized scaffolds (1.66 Gbp). 632 Mbp of sequence is
 171 represented as 6,303 unplaced scaffolds versus ~30 Mbp in 69 unplaced scaffolds in prairie-vole Hap1 in
 172 this study. Comparing the assemblies also revealed major improvements in contiguity (e.g., contig N50:
 173 29.2 kbp MicOch1.0 vs. 39.2 Mbp Hap1) and reduced gap content (8.0% vs. 0.07%), with 238 previously
 174 unplaced MicOch1.0 scaffolds (520.8 Mbp) now incorporated in Hap1 chromosome-scale scaffolds
 175 (Figure S8, Tables S5 and S6).

176

177 Similarly, the meadow vole assembly matches the expected karyotype ($2n=46$)^{25,26} across 44 autosomes,
 178 chrX, and chrY, representing 4.47 Gbp total sequence (Hap1: 2.34 Gbp). All but two chromosomes have
 179 at least one telomere detected (mean 3.4 gaps/chr), with 17 T2T scaffolds and six gapless contigs (overall
 180 50%; Figure 2, Table S4). A majority of the meadow vole autosomes are telocentric (42/44), with the two
 181 sex chromosomes classified as subtelocentric^{25,26}, likely contributing to the higher proportion of single
 182 telomeres detected in the meadow versus prairie vole. This assembly is on par with the recent long-read
 183 VGP meadow vole reference (mMicPen1), which has all expected chromosomes at near equal contiguity
 184 (scaffold Hap1 N50: 125.8 Mbp vs. this study 125.3 Mbp) but a larger number of unanchored sequences
 185 ($n=156$ at 46.1 Mbp vs. this study $n=81$ at 19.7 Mbp; Figure S9).

186



187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

Figure 2. Prairie and meadow vole genome assembly characteristics. (A) Circos plots depicting chromosome characteristics and content of Hap1 curated assemblies for prairie (left) and meadow (right) vole assemblies. Tracks are described in the legend from outer ring (1: Ideogram) to inner ring (6: Segmental duplication (SD)). (B) Repeat and SD distribution of both genomes are depicted above and below the Hap1 chromosome Y ideogram for each vole species, respectively. Annotations on the ideogram use colors depicted in the legend from (A), including presence of centromere and assembly gap (legend 1), telomere status (legend 2), SD (purple; legend 6), and repeat types (legend 7). Below the ideogram are dot plots for each species' chromosome Y, with colors representing percent identities. Annotated ideograms of all Hap1 chromosomes are shown in Figures S10 and S11.

Moving forward, we characterized primary Hap1 chromosomes, comprising single sets of autosomes and both sex chromosomes. Using RNA-seq data from liver, brain, and testes of adults and e11.5 (prairie vole) or e13.5 (meadow vole) embryos, we annotated ~20,500 protein-coding genes per species (Methods, Figure 2, Table S7). Repeat annotation revealed an increased proportion of interspersed repeats in prairie vole (e.g., L1 elements representing 16.7% vs. 13.0% of the genome), satellite repeats (8.62 Mbp vs. 511 kbp), and unclassified repeats (>40 Mbp more than meadow vole) (Figures 2A, S10, S11,

203 Table S8). Near-equal numbers of segmental duplications (SDs; regions >1 kbp at >90% sequence
204 identity^{27,28}) were identified in both genomes (~5.3K), with 72–79% intersecting an annotated gene
205 (prairie vole: 1,194 genes; meadow vole: 1,151 genes). Despite this similarity, prairie vole harbored
206 nearly double the duplicated sequence (54 Mbp, 2.15% of genome) compared to meadow vole (28 Mbp,
207 1.2%) (Figures 2A, S10, S11, Table S9). This expanded SD content is driven by larger duplicons (15.2
208 kbp vs. 8.21 kbp) and a greatly expanded interchromosomal SD fraction (34.7 Mbp, 63.7% of pairs vs.
209 9.7 Mbp, 46.2%).

210

211 We next examined the sex chromosomes given their known rapid and repeated structural remodeling in
212 *Microtus*²⁹. The most notable differences are evident on the Y chromosome (Figure 2B): both species
213 showed gene depletion and likely constitutive heterochromatin spanning either half (prairie vole) or the
214 entirety (meadow vole) of the chromosome. The remaining half of the prairie vole chrY comprises
215 tandemly arrayed SDs of a *Usp9y* mouse homolog, while the meadow vole chrY also carries considerable
216 SD content. Direct sequence comparison reveals no orthology (Figure 3A), indicating that homologous Y
217 chromosomes can reach complete sequence divergence over remarkably short evolutionary timescales.

218

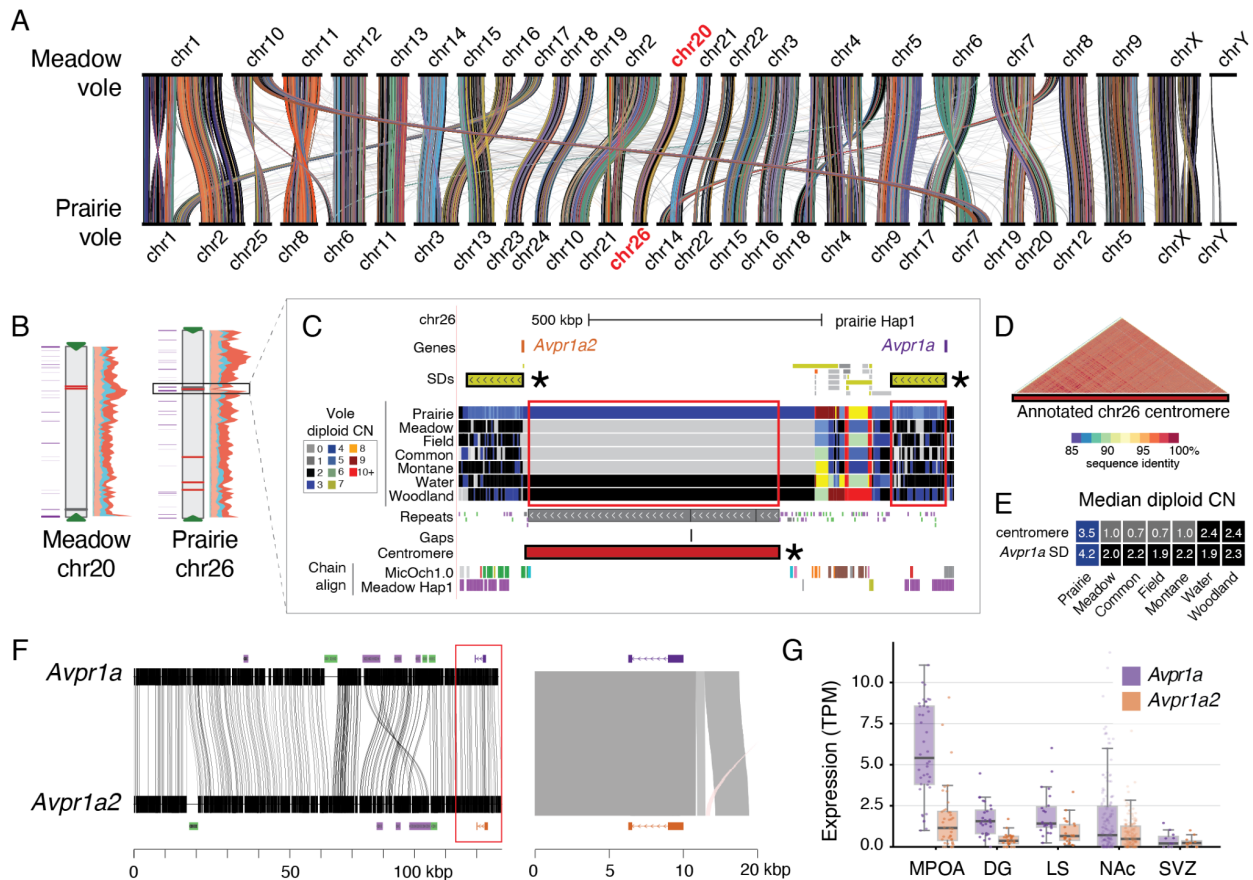
219 **Comparison of vole genomes identifies *Avpr1a* prairie-vole specific duplication**

220 To better match orthologous regions, we re-oriented Hap1 for both vole assemblies based on synteny to
221 the prairie vole reference (MicOch1.0) and linkage mapping (Tables S5, S10, and S11). Alignment
222 between the two species shows ~98% similarity with notable cytogenetic differences, including several
223 fusions and fissions, impacting 14 chromosomes in each genome (Figure 3A). Large-scale inversions are
224 evident on homologous chromosomes 1 and X, respectively. Mapping single-nucleotide, indel, and
225 structural variants identified ~462K coding variants using prairie vole as a reference³⁰, of which 17,466
226 were annotated as likely gene disrupting (LGD) impacting 2,418 genes (Table S12). Focusing on
227 candidate genes implicated in behavior—including those known to interact with nonapeptide hormones
228 oxytocin and vasopressin—show global reduced rates of substitutions relative to synonymous mutations
229 ($Ka/Ks \ll 1$; Figure S12) and no LGD variants obviously impacting function, suggesting purifying
230 selection (Tables S13 and S14).

231

232 Understanding that gene duplication is a common mechanism of trait innovation across the animal
233 kingdom³¹, we identified 306 and 389 gene duplications unique to prairie or meadow vole, respectively
234 (Table S15). Intersecting with our candidate genes list, we identified a prairie-vole-specific duplication of
235 *Avpr1a*, encoding vasopressin V1A receptor, a transmembrane G-protein coupled receptor that binds
236 arginine vasopressin. While this duplication was characterized via BAC sequencing nearly 15 years ago
237²⁰, both *Avpr1a* paralogs are missing in the current prairie vole reference (MicOch1.0) impeding genomic
238 analyses of these genes. Each prairie vole haplotype shows *Avpr1a* paralogs residing ~900 kbp apart on
239 chromosome 26, separated proximally by a large stretch of nearly identical satellite repeats not present at
240 the syntenic chr20 locus of the meadow vole (see chain alignment for Meadow vole Hap1; Figures 3B–
241 3D). The *Avpr1a* paralogs flank the chr26 annotated metacentric centromere, operationally defined here
242 as the largest stretch of satellite sequence per chromosome³², which is not observed in chromosome 20 of
243 the meadow vole (annotated as telocentric). A search for this sequence identifies matching satellite repeat
244 sequences on chromosomes 4 and X (~500 bp to 3 kbp in size at <93% sequence identity) for both Hap1
245 and Hap2 meadow vole assemblies but is found only at the chr26 *Avpr1a* locus in both prairie vole
246 haplotypes (547 kbp and 517 kbp in size, respectively).

247



248

249 **Figure 3. Prairie versus meadow vole cross-species comparison** (A) Assembly comparisons of synteny between
 250 prairie and meadow voles depicted as ribbon plots³³ colored by ancestral linkage groups (ALGs³⁴). (B) Repeat and
 251 segmental duplications (SDs) of *Avpr1a* flanking the prairie vole chr26 annotated centromeric region compared with
 252 meadow vole chr20. Color scheme described in Figure 2B. (C) A UCSC Genome Browser screenshot of the prairie
 253 vole Hap1 *Avpr1a* locus (chr26:13.0 Mbp–14.1 Mbp), with SDs of the paralogs and intervening centromeric repeat
 254 element highlighted with an asterisk. Windowed copy numbers (CN) of *Microtus* vole species are depicted as colors.
 255 Alignment chains for the current prairie vole reference (MicOch1.0) and the meadow vole Hap1 assembly (this
 256 study) are depicted at the bottom. (D) A dotplot of the annotated centromere with color representing sequence
 257 identities. (E) Both the annotated centromere and *Avpr1a* are at increased diploid CN uniquely in prairie vole. The
 258 two genotyped regions are highlighted as red boxes and colors match the legend in (C). (F) Hap1 genomic
 259 alignments³⁵ of prairie vole *Avpr1a* (purple) and *Avpr1a2* (red) including repeat annotations. (G) Transcriptomic
 260 analysis showing expression of both *Avpr1a* paralogs in medial preoptic area (MPOA), dentate gyrus (DG), lateral
 261 septum (LS), nucleus accumbens (NAc), and subventricular zone (SVZ). No expression was detected for either
 262 paralog in amygdala, hypothalamus, and ventral pallidum (not shown). TPM, transcripts per million reads.
 263

264 We next computed diploid copy number (CN) genome-wide from Illumina read depth³⁶ using data
 265 generated in this study from prairie vole, meadow vole, and woodland vole (*M. pinetorum*, also known as
 266 pine vole, also exhibiting monogamous pair bonding^{37,38}) alongside four publicly available *Microtus*
 267 genomes^{39,40} representing species with more promiscuous behaviors^{41–46}. While CN-estimates tend to be
 268 unreliable across repetitive loci due to repeat masking, prairie vole showed elevated CN across the chr26
 269 centromere annotation relative to all other voles, suggesting this region is prairie-vole specific (Figure
 270 3E). Also querying the CN of the 118-kbp SD spanning the complete *Avpr1a* gene and adjacent regions

271 confirmed this duplication to be prairie-vole specific with diploid CN of 4 compared to CN of 2 for all
272 other tested vole genomes.

273
274 Both curated prairie vole references show 97% nucleotide identity between *Avpr1a* paralogs, with
275 *Avpr1a2* harboring frameshift variants consistent with a truncated protein (218 aa vs. 420 aa full-length).
276 This result was supported by both Hap1 and Hap2 assemblies (Figure S13 and Table S16) and is
277 consistent with prior reports^{19,20}. If translated, the transcript would produce a truncated protein (218
278 amino acids (aa) versus 420 aa full length), sharing the first identical 199 aa with *Avpr1a* followed by
279 novel 19 aa, that includes the extracellular vasopressin-binding domain and the first four of seven
280 transmembrane domains but lacks the intracellular G-protein binding region⁴⁷. Comparing transcript
281 abundances using published RNA-seq data⁴⁸⁻⁵³ reveals consistent expression of both paralogs across five
282 of eight tested brain regions, with highest expression in the medial preoptic area, albeit reduced for
283 *Avpr1a2* (~1-4× relative to *Avpr1a*) (Figures 3G and S14); this could be influenced by a 600-bp indel
284 ~1.8 kbp upstream of *Avpr1a* paralogs (Figure 3F) and altered chromatin interactions outside of the
285 duplicated region⁵⁴, with the SD breakpoint only 4 kbp upstream. Even if *Avpr1a2* ultimately proves to
286 be a pseudogene, the corrected assembly enables expression and epigenomic analysis of the ancestral full-
287 length *Avpr1a* and its *cis*-regulatory landscape, previously absent from the prairie vole reference genome
288 entirely.

289

290 DISCUSSION

291 Assemblies enable robust and complete genome comparisons of variation contributing to divergence,
292 diversity, and disease. Most recent efforts to build gigabase-sized chromosome-scale genomes require
293 multiple technologies⁵⁵, typically comprising highly accurate HiFi long reads (~15-20 kbp at 30-60×
294 coverage), ultralong nanopore reads (>100 kbp at ~30× coverage), and short-read long-range information
295 to phase and scaffold at chromosome scale⁵⁶; even higher coverage is necessary to achieve complete
296 diploid T2T genomes⁵⁷. This “ideal” recipe is inaccessible to many researchers, requiring large amounts
297 of starting material, multiple libraries, and three sequencing platforms. Here, we generated ~100 Gbp of
298 HiFi and CiFi data from a single library and sequencing, producing contiguous (scaffold N50 90-125
299 Mbp), high-quality diploid assemblies (QV>60). The experimental approach is conceivably scalable to
300 tens of thousands of cells and nanograms of DNA/chromatin⁷. Based on CiFi downsampling (Figure S5),
301 a 2× CiFi and 30× HiFi coverage ratio is theoretically sufficient to assemble diploid genomes as large as
302 ~3.1 Gbp (human sized) on a single SMRT Cell. Compared with the “ideal” multi-platform recipe, this
303 approach reduces costs by at least threefold while still producing over half T2T-scale chromosomes.

304

305 Beyond introducing a simplified assembly approach, we provide a long-awaited resource expanding the
306 genomic toolkit for prairie voles, an emerging model organism for studying complex social behaviors
307 relevant to humans⁵⁸. This includes high-quality annotated assemblies for both prairie and meadow voles,
308 repeat and SD annotations, CN maps across additional vole species, chain alignments facilitating genome
309 liftover, and named gene orthologs enabling transcriptomic comparisons, all publicly accessible through a
310 UCSC Genome Browser hub (see Data Availability). These resources complement ongoing
311 neurobiological studies examining how social relationships are encoded in the brain.

312

313 Neuropeptide hormone pathways represent an important substrate through which neural circuits
314 mediating innate behaviors can rapidly evolve and diversify^{12,13}. Hence, divergence of the genes
315 encoding neuropeptide receptors or their regulatory regions may underlie species differences in behavior.
316 The highly contiguous assemblies generated here enabled discovery of a complex >350-kbp satellite
317 repeat element with centromere annotation and confirmed the presence of a >100 kbp SD impacting
318 *Avpr1a* (Figure 3). Neither feature is present in any of the four meadow vole haplotypes examined (this
319 study and mMicPen1), nor was elevated CN detected at these regions in the five other vole species tested,
320 including woodland vole, which also exhibits pair bonding. While this specificity to prairie vole argues
321 against these variants as cross-species drivers of pair bonding, they remain compelling candidates for
322 prairie-vole-specific function. Notably, if epigenetic signatures confirm the satellite repeat as a true
323 centromere, its position ~350 kbp downstream of *Avpr1a* would place the gene within the pericentromeric
324 region. Centromerization has been shown to result in increased H3K27me3-mediated repression, altered
325 chromatin accessibility, and elevated genomic instability⁵⁹—consequences that may alter vasopressin
326 signaling dynamics and could have contributed to the birth of the prairie-vole-specific *Avpr1a2* paralog.
327 This is a particularly compelling finding given that prairie voles exhibit elevated *Avpr1a* expression in the
328 ventral pallidum while meadow voles show higher expression in the lateral septum¹⁹. These assemblies
329 will serve as an important resource for continued investigation of *Avpr1a* paralog functions and divergent
330 brain expression patterns.

331
332 Our assessment of candidate behavioral genes (Tables S13 and S14) identified no obvious protein-coding
333 differences between prairie and meadow voles, suggesting that regulatory divergence, rather than coding
334 sequence change, plays an outsized role in behavioral differences between these species. These near-
335 complete genomes will enable comparative, genome-wide analyses of chromatin structure and gene
336 regulation, providing a framework for examining regulatory mechanisms underlying bond formation and
337 social memory. Ultimately, connecting vole genomics to human genetic studies offers a path toward
338 understanding how social behavior is encoded in the brain, and how its disruption contributes to
339 neuropsychiatric risk.

340 341 **Limitations of the study**

342 A notable limitation is that current assembly and scaffolding tools, including hifiasm and YAHS, do not
343 leverage the multi-contact information inherent to CiFi concatemers, instead requiring reduction into Hi-
344 C-like pairs. Development of algorithms designed to leverage these higher-order chromatin contacts will
345 likely yield further improvements in phasing, with CiFi showing an >8-fold greater ability to infer
346 haplotypes versus Hi-C when considering multicontacts⁷, as well as scaffolding, particularly in complex
347 genomic regions. Additionally, use of a single restriction enzyme (HindIII) introduces the possibility of
348 sequence-biased coverage (also inherent in short-read Hi-C); adoption of alternative fragmentation
349 approaches such as Omni-C could mitigate this in future implementations. Gene annotations might be
350 further refined by incorporating long-read isoform sequencing, enabling a fully integrated workflow for
351 assembly, scaffolding, and annotation from a single sequencing technology. Finally, although two
352 haplotypes were resolved per species through diploid assembly, single-individual sampling limits the
353 ability to distinguish fixed from polymorphic differences. Broader sampling within and across *Microtus*
354 species using a phylogenetic framework will be necessary to determine whether identified variants
355 segregate with behavioral phenotypes, an effort complicated by the uncertain evolutionary relationships
356 within this large clade of >60 species that has undergone rapid radiation over the last ~2 million years⁶⁰.

357 In summary, we present a combined HiFi and CiFi sequencing strategy and simplified bioinformatic
358 workflow enabling accurate, contiguous, chromosome-scale genome assembly from a single sequencing
359 platform and experiment. Applying this approach to prairie and meadow voles, we generate high-quality
360 diploid assemblies that reveal genome-wide differences between these behaviorally divergent species,
361 with particular focus on loci implicated in social behavior. Ultimately, the simplicity, affordability, and
362 minimal sample requirements of the combined HiFi–CiFi approach position it as a broadly accessible tool
363 for comparative genomics, with the future potential to make chromosome-scale assembly tractable for
364 rare, difficult-to-sample, and non-model organisms alike.

365

366 **STAR METHODS**

367 **EXPERIMENTAL MODEL AND STUDY PARTICIPANT DETAILS**

368

369 **Vole procedures**

370 Adult prairie vole (*M. ochrogaster*) and meadow vole (*M. pennsylvanicus*) individuals were maintained
371 under protocols approved by the Institutional Animal Care and Use Committee (IACUC) at the
372 University of California, San Francisco and Cold Spring Harbor Laboratories. Animals were euthanized
373 according to institutional guidelines, and tissues were rapidly dissected, flash-frozen in liquid nitrogen,
374 and stored at -80°C until DNA or RNA extraction.

375

376 **METHOD DETAILS**

377

378 **HiFi and CiFi library preparations and sequencing of prairie and meadow vole liver samples**

379 For whole-genome HiFi sequencing, high molecular weight (HMW) genomic DNA was extracted from
380 frozen liver tissue using the Monarch HMW DNA Extraction Kit for Tissue (New England Biolabs)
381 according to the manufacturer's protocol. DNA quality and fragment size distribution were assessed by
382 Qubit fluorometry and Femto Pulse analysis prior to downstream HiFi library preparation.

383 For CiFi library preparation, frozen liver tissue (~100 mg) was pulverized under liquid nitrogen and
384 crosslinked, followed by processing according to CiFi protocol Part 1 (HindIII restriction digestion and
385 proximity ligation)⁷. Crosslinks were reversed, and proximity-ligated 3C DNA was purified by phenol–
386 chloroform extraction prior to downstream size selection and amplification.

387

388 A starting amount of 2.2 μg of HMW DNA was sheared on a Hamilton NGS STAR MOA system and
389 size-selected on the Pippin HT with 10 kbp cut off. The pre-PCR CiFi DNA (1.5 μg) was size-selected
390 on the Pippin HT with 6.5 kbp cutoff and PCR amplified (50 ng input, 8 cycles) using the Ampli-Fi
391 protocol (PacBio, 103-648-000). The post-PCR CiFi DNA was size-selected on the Pippin HT with 5.5
392 kbp cutoff. They were subsequently mixed at a molar ratio of 90% sheared DNA for HiFi and 10% CiFi
393 DNA (translating to 650 ng and 35 ng, respectively, accounting for the size differences), followed by
394 standard library preparation using SMRTbell prep kit 3.0 (PacBio, 102-182-700). The SMRTbell library
395 was cleaned up with 1X SMRTbell cleanup beads. One SMRT Cell per species was run on the Revio
396 system at 250 pM on-plate loading concentration with SPRQ chemistry and 30-hours acquisition.

397

398 HiFi and CiFi data from the sequencing run were segregated and adapters were trimmed with lima
399 v2.14.0 (<https://github.com/PacificBiosciences/barcoding>) using a three-step process. First, the CiFi reads

400 were identified by their unique dual index adapters and separated from HiFi reads using very relaxed lima
401 settings (`--ccs --min-passes 0 --min-end-score 0 --min-score 5 --min-ref-span 0.2 --min-score-lead 0`) to
402 ensure that no CiFi reads remained in the HiFi data. Adapters were then trimmed from each dataset using
403 the recommended settings: `--hifi-preset SYMMETRIC` for the HiFi reads and `--hifi-preset`
404 `ASYMMETRIC --neighbors` for the CiFi reads. PCR duplicates were then removed and duplication rates
405 assessed for the CiFi reads using pbmarkdup v1.1.0 (<https://github.com/PacificBiosciences/pbmarkdup>).

406

407 **Illumina sequencing vole genomic samples**

408 Chromatin isolation and Hi-C library preparation was performed by Phase Genomics (WA, USA), using
409 100 mg flash-frozen liver tissue from the same individual male prairie and meadow voles as were used for
410 HiFi and CiFi library preparation. The Proximo Hi-C protocol (Phase Genomics, WA, USA)⁶¹ was used
411 to prepare the proximity ligation library and process it into an Illumina-compatible sequencing library.
412 Hi-C libraries were sequenced with 150-bp paired-end reads on an Illumina NextSeq500 at the CSHL
413 Next Generation Sequencing Core Facility.

414

415 gDNA was extracted from dried museum pelt tissue of the woodland vole (*M. pinetorum*) using the
416 DNeasy Blood & Tissue Kit (Qiagen) following manufacturer instructions with minor modifications to
417 reduce PCR inhibitors. Extracted DNA was treated with RNase A and further purified by ethanol
418 precipitation. Sequencing libraries were prepared for short-read sequencing using prepared genomic DNA
419 for prairie vole, meadow vole, and woodland vole. Prepared libraries were sequenced (2×150 bp) by
420 Novogene on a NovaSeq 6000 platform to produce approximately 60 Gbp of raw data (~30× genome
421 coverage) for each species.

422

423 **Genome assembly and scaffolding**

424 HiFi reads were extracted from unaligned BAM files using samtools v1.21⁶² (`samtools fasta`). CiFi
425 reads were converted from BAM to FASTQ (`samtools collate -O -u | samtools fastq`) and then processed
426 into Hi-C-like paired-end reads using the cifi-toolkit (<https://github.com/mr-eyes/cifi-toolkit>) that
427 performs in silico HindIII digestion on each CiFi concatemer read, extracting the outermost restriction
428 fragments as R1 and R2 paired-end reads. Phased diploid assembly was performed with hifiasm v0.19.8⁴
429 using the `--dual-scaf` mode, which integrated either CiFi (input as paired-end reads) or Hi-C contact
430 information during graph resolution to produce haplotype-resolved primary contig graphs for each
431 haplotype (Hap1 and Hap2). HiFi FASTA and derived R1/R2 FASTQ files were provided as input (`--`
432 `h1`, `--h2`). GFA contig graphs were converted to FASTA using gfatools v0.5
433 (<https://github.com/lh3/gfatools>).

434

435 CiFi and Hi-C contact maps were generated by aligning the full CiFi BAM to each haplotype assembly
436 using the epi2me-labs/wf-pore-c Nextflow pipeline v1.3.0 (<https://github.com/epi2me-labs/wf-pore-c>)
437 with minimap2 in `-ax map-hifi` mode, HindIII as the restriction enzyme, and `--paired_end` enabled to
438 produce BED contact files. Contig scaffolding was performed with YAHS v1.2a.2²¹ using the porec BED
439 contact file as input, with contig error correction disabled (`--no-contig-ec`).

440

441 To evaluate the minimum CiFi sequencing depth required for chromosome-scale scaffolding, the CiFi
442 BAM was subsampled at 1%, 10%, 15%, 20%, 25%, 40%, 60%, 80%, and 100% of the original read
443 depth using `samtools view -s` with a fixed random seed. At each fraction, the full assembly and

444 scaffolding pipeline described above was executed independently, and scaffold N50 and auN were
445 compared across titration points.

446

447 **Genome assembly manual curation**

448 To facilitate manual assembly curation for each genome of the Prairie and the Meadow voles the two
449 scaffolded haplotypes were combined together and the corresponding CiFi datasets were mapped onto the
450 assemblies using wf-pore-c pipeline with parameters `--paired_end --cutter HindIII --
451 paired_end_minimum_distance 100 --paired_end_maximum_distance 200`. Hi-C maps were further
452 generated for the combined haplotypes of each genome with PretextView using the mock paired-end BAM
453 files as input and retaining non-uniquely mapping reads in the contact maps with --mapq 0.

454 Supplementary analysis were also embedded in the PretextView file using the Tree of Life curationpretext
455 pipeline - telomeres with standard vertebrate motif TTAGGG, N base scaffold gaps and mapped PacBio
456 long-read coverage. An AGP of the corrected contact maps were exported from PretextView and curated
457 haplotype assemblies generated using pretext-to-asm.

458

459 PretextView <https://github.com/sanger-tol/PretextView>

460 PretextViewMap <https://github.com/sanger-tol/PretextViewMap>

461 Curationpretext <https://github.com/sanger-tol/curationpretext>

462 PretextView-to-asm <https://github.com/sanger-tol/agp-tpf-utils>

463

464 **Telomere-to-telomere (T2T) completeness and gap assessment**

465 To evaluate the contiguity and completeness of our assemblies, we assessed telomere presence and
466 sequence gaps across all chromosome-scale scaffolds for both haplotypes of each species. Gaps (runs of
467 N bases) were identified using seqtk gap (<https://github.com/lh3/seqtk>). Telomeric repeat sequences were
468 detected using tidk v0.2 (Telomere Identification Toolkit; ⁶³), which scans for the canonical vertebrate
469 telomere motif (TTAGGG) in sliding windows across each scaffold. We searched for TTAGGG repeats
470 in 10 kbp windows and considered a telomere present at a chromosome terminus when the terminal
471 window contained at least 10 repeat units (forward and reverse strands combined). Chromosomes were
472 classified as T2T (telomeric repeats at both ends), partial (one end only), or incomplete (neither end). A
473 chromosome was considered fully T2T-resolved only when both telomere ends were detected and no
474 sequence gaps were present.

475

476 **Centromere annotation**

477 Centromeric regions were identified on each hap1 assembly using centroAnno v1.0.2 ³² in assembly
478 annotation mode (`-x anno-asm`). Each chromosome was processed individually; centroAnno
479 decomposed tandem satellite repeats into monomer units and reported their genomic coordinates.

480 Contiguous monomer annotations separated by fewer than 10 kbp were merged into repeat regions, and
481 the largest repeat region per chromosome was designated as the putative centromere. Centromeres were
482 detected on all chromosomes of both species.

483

484 **Segmental duplication and repeat annotations**

485 Species-specific *de novo* repeat libraries were constructed for each vole species using RepeatModeler
486 v2.0.7 ⁶⁴ with default parameters. A RepeatModeler database was built from each hap1 primary assembly,
487 and *de novo* repeat family consensus sequences were identified. The resulting species-specific libraries

488 were combined with the Dfam database and used as input for RepeatMasker v4.2.2⁶⁵ with the
489 RMBLAST search engine. RepeatMasker was executed via the Dfam TE Tools container
490 (<https://hub.docker.com/r/dfam/tetools>; RepeatModeler v2.0.7, RepeatMasker v4.2.1, RMBLAST
491 v2.14.1+). Segmental duplications were identified using BISER v1.4²⁸ with default parameters. BISER
492 was run on the soft-masked hap1 assemblies for each species, and output was generated in BEDPE
493 format.

494

495 **Transcriptome analysis of vole tissues**

496 Total RNA was prepared from flash frozen prairie and meadow vole tissues: liver, brain, and testes from
497 adults and e11.5 (prairie) or e13.5 (meadow) embryos. 50–75mg of tissue per sample was homogenized in
498 500 µl Trizol (Life Technologies) on ice using a Kimble Kontes Disposable Pellet Pestle (VWR). An
499 additional 500 µl Trizol was added followed by further homogenization using an 18-gauge needle on a
500 1ml syringe before proceeding to RNA extraction according to the Trizol protocol. The subsequent RNA
501 was DNase treated with a TURBO DNA-free kit (Life Technologies). 150 ng of total RNA was used as
502 input for library preparation with Encore Complete RNA-seq kits (NuGen), using ten cycles of
503 amplification. Multiplexed libraries were sequenced with 76-bp paired-end reads on the Illumina NextSeq
504 500 at the CSHL Next Generation Sequencing Core Facility.

505

506 *Gene annotation:* Illumina RNA-seq data from the four samples per species were used in the NCBI
507 Eukaryotic Genome Annotation Pipeline (EGAPx) v0.4.1-alpha⁶⁶ to perform gene annotation for
508 bothHap1 assemblies of both species.

509

510 *Gene expression analysis:* Transcript quantification was performed using Salmon v1.10.3⁶⁷ in mapping-
511 based mode. For each species, *Avpr1a* transcript-level expression was first quantified across four in-house
512 paired-end Illumina RNA-seq libraries (brain, liver, testes, and embryo) using a targeted Salmon index
513 containing the *Avpr1a* coding sequences with the hap1 genome assembly as a decoy. For the prairie vole,
514 the index included both *Avpr1a* paralogs to enable paralog-specific quantification. To contextualize
515 *Avpr1a* expression within the broader transcriptome, we performed transcriptome-wide Salmon
516 quantification for the prairie vole. The full EGAPx-annotated transcriptome was extracted from the prairie
517 hap1 assembly using gffread, and a decoy-aware Salmon index was built by concatenating the
518 transcriptome with the genome assembly. We quantified 453 publicly available prairie vole brain RNA-
519 seq libraries from previously published studies^{48–53} spanning eight brain regions—amygdala (AMY),
520 dentate gyrus (DG), hypothalamus (HT), lateral septum (LS), medial preoptic area (MPOA), nucleus
521 accumbens (NAc), subventricular zone (SVZ), and ventral pallidum (VP)—under BioProject accessions
522 PRJNA428754, PRJNA682808, PRJNA631040, PRJNA786347, PRJNA887096, PRJNA792575,
523 PRJNA1005323, and PRJEB89367. After excluding 147 technical replicates, 306 samples were retained
524 for analysis. Salmon was run with `--validateMappings`, `--gcBias`, and `--seqBias` correction, with
525 library type automatically inferred. Differential expression analysis across brain regions was performed
526 using pyDESeq2⁶⁸.

527

528 **Copy-number analysis**

529 Copy number (CN) was estimated using the FastCN pipeline³⁶ with mrsFAST v3.4.2. All species were
530 mapped to the prairie vole hap1 assembly as a common reference to enable direct cross-species
531 comparison at the same genomic coordinates. A four-layer masked reference was constructed from the

532 prairie vole hap1 assembly. Repetitive elements were masked using the RepeatMasker annotations
533 described above, tandem repeats were identified with Tandem Repeats Finder⁶⁹ run per chromosome, and
534 low-complexity regions were masked with WindowMasker/DUST⁷⁰. The masked genome was then
535 indexed with mrsFAST and all 50-mers were extracted and searched back against the genome; positions
536 where a 50-mer aligned 20 or more times were additionally masked (K50 masking). Gaps in the final
537 masked reference were extended by 36 bp on each side to account for the read-length shadow effect.
538 Control regions expected to be diploid (CN=2) were defined by excluding segmental duplications
539 (identified by BISER), centromeric regions, and target gene loci from the genome-wide window set.

540

541 For each species, the first 36 bp were extracted from each mate of the paired-end whole-genome
542 sequencing reads and mapped to the masked reference using mrsFAST with up to two mismatches
543 allowed. GC-corrected read depth was computed in 1 kb windows across the genome. Copy number was
544 calculated as $CN = 2 \times (\text{window depth} / \text{autosomal control mean})$. To prevent inflated copy number
545 values from masked regions with zero depth falling within control windows, we excluded zero-depth
546 windows from the control mean calculation. Control regions were further refined by retaining only
547 windows where all samples showed consistent CN near 2.0 (coefficient of variation < 0.2), and a final
548 normalization ensured the median CN at refined control windows equaled 2.0.

549

550 Whole-genome sequencing reads from seven *Microtus* species —prairie vole (*M. ochrogaster*, this study),
551 meadow vole (*M. pennsylvanicus*, this study), woodland vole (pine vole, *M. pinetorum*, this study),
552 montane vole (*M. montanus*; SRR12966109), common vole (*M. arvalis*; ERR3427942), field vole (*M.*
553 *agrestis*; SRR2167807), North American water vole (*M. richardsoni*; SRR12963053)— were mapped to
554 the prairie vole reference. Each species was then normalized independently by scaling CN values so the
555 median at refined diploid autosomal control regions equals 2.0, correcting for sequencing depth
556 differences between species. Per-gene CN was calculated as the mean CN across all non-zero 1 kb
557 windows overlapping each target locus. CN was evaluated at 15 target loci spanning the
558 vasopressin/oxytocin system (*Avpr1a*, *Avpr1b*, *Avp*, *Oxt*, and *Oxtr*), dopamine system (*Drd2*), estrogen
559 system (*Esr1*, *Esr2*, and *Cyp19a1*), stress system (*Crh*, *Crhr1*, and *Crhr2*), neural development genes
560 (*Chd8* and *Shank3*), and the androgen receptor (*Ar*).

561

562 Dot plots were generated using ODP v0.3.3³³, circular genome visualizations using pyCirclize
563 v1.9.1(<https://github.com/moshi4/pyCirclize>), and linear karyotype ideograms using matplotlib v3.10.8.
564 Assembly accuracy was evaluated using Yak⁴ by concatenating both haplotypes to calculate a combined
565 genome-wide adjusted QV against HiFi reads, while k-mer completeness was assessed by providing both
566 haplotypes in diploid mode to Merquy²².

567

568 **Comparative assembly and paralog analyses**

569 Ortholog groups were identified using Orthofinder (v2.5.5)⁷¹ with default parameters, using protein
570 sequences from both species. Orthogroups were classified as one-to-one, expanded in one of the species,
571 multi-copy in both, or species-specific. For the one-to-one ortholog pairs, we aligned CDS using parasail
572 (v2.6.2)⁷², then variants were called from the pairwise alignments. SNPs and indels were identified by
573 parsing the alignment CIGAR string. Variants were classified by comparing reference and alternate
574 codons using the standard genetic code: synonymous if both codons encode the same amino acid,
575 missense if different, and nonsense if the alternate codon is a stop codon. For genes with multiple

576 annotated isoforms, the longest one was selected as a representative. Ka/Ks (dN/dS) were calculated using
577 the Nei-Gojobori method implemented in BioPython v1.83⁷³. CDS length ratio <0.8 were excluded as
578 unreliable. Structural differences between ortholog CDS pairs were detected using BLASTN megablast
579 alignments (BLAST+ v2.16.0)⁷⁴. Frameshifts were defined as insertions/deletions events in the CDS
580 alignments with length not divisible by 3. Gene-level inversions were identified when BLASTN returned
581 alignments on opposite strands. Exon structure changes were detected by comparing the number,
582 boundaries, and sizes of exons between ortholog pairs in the EGAPx annotations. Truncations were
583 flagged when the CDS length ratio between species was < 0.8. *Avpr1a* coordinates for Hap2 assemblies
584 were obtained by aligning Hap1 regions \pm 50 kbp to Hap2 with minimap2 v2.30 with -x asm5, followed
585 by annotation transfer with miniprot v0.18. CDS and protein MSAs (6 sequences: 2 species \times 2
586 haplotypes, plus *Avpr1a2* \times 2 haplotypes) were generated with MACSE v2.07⁷⁵, which handles the +1 G
587 frameshift in *Avpr1a2* without breaking the reading frame. Gene DNA MSAs were generated with
588 MAFFT v7.525 (L-INS-i)⁷⁶. Pairwise SNPs were classified as synonymous, nonsynonymous, or intronic
589 using bcftools csq v1.22⁷⁷.

590

591 Pairwise comparisons and visualization of SD sequence containing *Avpr1a* (prairie vole Hap1
592 chr26:13970803-14087830) and *Avpr1a2* (prairie vole Hap1 chr26:13058418-13176477) was performed
593 using Micropeats³⁵. The annotated centromeric sequence (prairie vole Hap1 chr26:13182318-13728901)
594 was queried against both Hap1 and Hap2 of the prairie and meadow vole assemblies (this study and
595 mMicPen1) using minimap2 (v2.26)⁷⁸ selecting for matches 90% or higher. Dot plots generated using nf-
596 core/pairgenomealign⁷⁹.

597

598 **Scaffold orientation**

599 *Prairie vole*: To ensure compatibility with existing prairie vole genomic resources and maintain
600 consistent chromosome orientation conventions, scaffold orientation in the *de novo* haplotype-resolved
601 assemblies was standardized relative to the MicOch1.0 reference genome (GCF_000317375.1). The
602 assembly retains its original chromosome numbering (SUPER_1 through SUPER_26, SUPER_X,
603 SUPER_Y, renamed to chr1-chr26, chrX, chrY). The MicOch1.0 alignment was used exclusively to
604 determine whether each scaffold required reverse-complementation, without reassigning chromosome
605 identities.

606

607 Whole-genome alignments were conducted between each query haplotype assembly and MicOch1.0
608 using minimap2 v2.28⁷⁸ with the asm5 preset. Alignments were generated in both directions (query-to-
609 reference and reference-to-query) to facilitate orientation determination and to identify previously
610 unplaced MicOch1.0 scaffolds incorporated into chromosome-scale scaffolds in the new assembly.
611 Alignments shorter than 10 kb were excluded to minimize spurious matches in repetitive regions.

612

613 For each query scaffold, orientation was determined by quantifying the total aligned bases on the forward
614 (+) and reverse (-) strands relative to the best-matching MicOch1.0 chromosome or linkage group.
615 Scaffolds were designated for reverse-complementation if the majority of aligned bases mapped to the
616 reverse strand; otherwise, the original orientation was retained. Confidence was classified as high (>95%
617 of aligned bases on the dominant strand), medium (80-95%), or low/mixed (<80%, indicating mixed
618 strand orientations typical of alignments involving fragmented reference regions). As MicOch1.0 was
619 derived from a female individual, the Y chromosome scaffold (chrY) could not be oriented by alignment

620 and was therefore retained in its original orientation. Alignment-based orientation calls were
621 independently cross-validated using a radiation hybrid (RH) linkage map for the prairie vole²⁴. For each
622 linkage group, RH marker sequences were mapped to the *de novo* assembly using BLASTn, and the
623 Spearman rank correlation (ρ) between marker centimorgan position and physical position on the scaffold
624 was calculated. A positive ρ indicates concordance with the primary scaffold orientation, whereas a
625 negative ρ indicates the scaffold is in the reverse orientation relative to the genetic map. For scaffolds
626 with low alignment-based confidence (<80% dominant strand) and a strong linkage map orientation
627 signal ($|\rho| \geq 0.70$, ≥ 5 mapped markers), the linkage map call was prioritized as the primary evidence.
628 Cross-validation results are presented in Table S5.

629

630 Corrected assemblies were produced by reverse-complementing the designated scaffolds and renaming all
631 scaffolds according to a standardized nomenclature (chr1-chr26, chrX, chrY for chromosome-scale
632 scaffolds). A UCSC liftOver chain file was generated to facilitate coordinate conversion of annotation
633 files from the original to the corrected assembly. The complete set of orientation decisions, alignment
634 statistics, and confidence classifications for all 97 scaffolds is provided in Table S10.

635

636 The prairie vole Hap1 assembly consists of 97 scaffolds, including 28 chromosome-scale scaffolds (chr1-
637 chr26, chrX, chrY) and 69 minor unplaced scaffolds. All 26 autosomal scaffolds and chrX were assigned
638 to corresponding MicOch1.0 chromosomes or linkage groups, with alignment coverage ranging from
639 32.6% to 78.4%. As expected, chrY showed no alignment to the female-derived MicOch1.0 reference
640 (Table S10). Alignment-based orientation analysis revealed that 15 of 28 chromosome-scale scaffolds
641 were concordant with the MicOch1.0 convention (forward orientation), while 13 required reverse
642 complementation. Alignment confidence was classified as high for 7 scaffolds, medium for 9, low/mixed
643 for 11, and not assessable for 1 (chrY). The high frequency of low/mixed confidence values reflects the
644 substantial structural divergence between a chromosome-scale long-read assembly and a fragmented
645 short-read reference containing 6,336 sequences and approximately 631 Mbp of unplaced sequence.

646

647 Cross-validation with the RH linkage map was feasible for 26 of 28 chromosome-scale scaffolds (chrY
648 and chr22 lacked linkage markers), resulting in 29 scaffold-linkage group comparisons. Of these, 27 out
649 of 29 (93.1%) demonstrated concordant orientation calls between the two independent methods (Table
650 S5). Two comparisons were discordant: (i) chr10 was concordant with its primary linkage group (LG14; ρ
651 = +0.41, 11 markers) but discordant with a secondary linkage group (LGLG8; ρ = -0.98, 5 markers),
652 indicating that this scaffold incorporates sequences from multiple MicOch1.0 linkage groups in opposite
653 orientations; and (ii) chrX exhibited a weak linkage signal (ρ = +0.60, 11 markers) discordant with the
654 alignment call, consistent with limited recombination on the sex chromosome reducing the resolution of
655 marker-order correlations. For chr8, where alignment confidence was low/mixed (63.7% forward strand),
656 the RH linkage map provided strong evidence for reverse-complementation (ρ = -0.93, 13 markers),
657 which was used as the primary evidence source.

658

659 In the reference-to-query analysis, 238 previously unplaced MicOch1.0 scaffolds totaling 520.8 Mbp
660 were mapped within chromosome-scale scaffolds of the new assembly, with alignment coverage ranging
661 from 50.0% to 284.3% (median 78.8%) (Table S6). These resolved scaffolds were distributed across all
662 28 chromosomes, with the largest contributions to chrX (55.8 Mbp from 27 scaffolds), chr12 (78.8 Mb
663 from 10 scaffolds), and chr6 (46.0 Mb from 16 scaffolds). Among the 238 resolved scaffolds, 103

664 exceeded 1 Mb in length, and 149 exceeded 100 kb, indicating that the new assembly anchors a
665 substantial amount of previously unplaced sequence within a chromosomal context. An additional 69
666 minor scaffolds (hap1_scaffold_27 through hap1_scaffold_95; collectively 29.1 Mb) were evaluated. Of
667 these, 28 could be oriented by alignment to MicOch1.0 (16 reverse-complemented, 12 retained), while 41
668 showed no meaningful alignment and were retained in their original orientation. The complete scaffold-
669 to-chromosome assignments, orientation decisions, and validation results for all 97 scaffolds are provided
670 in Table S10.

671
672 *Meadow Vole*: Scaffold orientation in the meadow vole Hap1 assembly was standardized using the same
673 methodology described above, with the orientation-corrected prairie vole Hap1 assembly as the reference.
674 This ensures a consistent orientation convention across both species in this study. The minimap2 asm10
675 preset was used for cross-species alignment. The meadow vole Hap1 assembly comprises 105 scaffolds:
676 24 chromosome-scale (chr1-chr22, chrX, chrY; 2.34 Gb) and 81 minor unplaced (19.7 Mb). Alignment
677 coverage for the 23 evaluable chromosome-scale scaffolds ranged from 43.4% to 87.9%; chrY showed no
678 cross-species alignment and was retained in its original orientation. Of the 24 chromosome-scale
679 scaffolds, 12 were concordant and 12 required reverse-complementation, with confidence classified as
680 high (13), medium (4), low/mixed (6), or not assessable (1, chrY). The 6 low/mixed scaffolds reflect
681 cross-species rearrangements where a single meadow vole chromosome aligns to multiple prairie vole
682 chromosomes in mixed orientations (meadow 2n=46 vs. prairie 2n=54). Post-correction validation
683 showed 13 of 24 chromosome-scale scaffolds with $\geq 80\%$ forward-strand alignment; the remainder are
684 attributable to chromosomal fusions. Among the 81 minor scaffolds, 53 were oriented by alignment (29
685 reverse-complemented, 24 retained), and 28 showed no cross-species alignment (Table S11).

686

687 DATA AVAILABILITY

688 All data generated as part of this project can be found at the European Nucleotide Archive and NCBI
689 GenBank Project accession PRJEB108798. Previously published datasets have GenBank accessions listed
690 below.

691

692 Prairie and meadow vole reference genomes:

- 693 ● mMicPen1 reference: GCF_037038515.1
- 694 ● MicOch1.0 reference: GCF_000317375.1

695

696 RNA-seq of prairie vole brain regions: PRJNA428754, PRJNA682808, PRJNA631040, PRJNA786347,
697 PRJNA887096, PRJNA792575, PRJNA1005323, PRJEB89367

698

699 Illumina WGS sequence data of *Microtus* species: SRR12966109, ERR3427942, SRR2167807,
700 SRR12963053

701

702 Temporary hubs for UCSC Genome browser resources

- 703 ● Prairie vole: https://genome.ucsc.edu/s/mabuelanin/prairie_vole_cifi_hifi_hap1
- 704 ● Meadow vole: https://genome.ucsc.edu/s/mabuelanin/meadow_vole_cifi_hifi_hap1

705

706 **CODE AVAILABILITY**

707 Methods used for genome curation are available at <https://github.com/sanger-tol/>. A Zenodo doi will be
708 generated for code related to this manuscript at the time of publication. All source code and workflows
709 can be accessed through <https://github.com/mydennislabs/2026-voles-assembly>.

710

711 **ACKNOWLEDGEMENTS**

712 We thank Dr. Aaron Wenger and Jacob Brandenburg for coordination support in generating and
713 analyzing the HiFi-CiFi datasets. Also thanks to Drs. Michael Schatz and Daniela Soto for assistance with
714 earlier iterations of vole genome assemblies, and Michael Sherman for help with animal husbandry. This
715 work was supported, in part, by the U.S. National Science Foundation (CAREER 2145885 to M.Y.D.;
716 CAREER IOS-2045348 to Z.R.D.) and the National Institutes of Health (NIH) grants from the National
717 Institute of Mental Health (R01MH132818 to M.Y.D. and DP2MH119427 to Z.R.D.). K.K., J.W., and
718 K.H. are supported by Wellcome through the 220540/Z/20/A award that supports the Wellcome Sanger
719 Institute. Voles and tissue samples were supported by NIMH (R01MH123513) and NSF (1556974) grants
720 to D.S.M. CSHL voles studies were supported by an institutional grant to J.T. Woodland vole analysis
721 was supported by a UC-Davis Interdisciplinary Catalyst Award to K.L.B. and M.Y.D. Illumina
722 sequencing of RNA-seq and Hi-C libraries was performed at the Cold Spring Harbor Next Generation
723 Sequencing Shared Resource, which is supported by NIH Cancer Center Support Grant 5P30CA045508
724 Some images were created using BioRender.

725

726 **AUTHOR CONTRIBUTIONS**

727 G.K., C.L., M.V.W., N.G.S., Z.R.D., K.L.B., and J.T. generated biospecimens and/or data for the project.
728 M.A., J.L., K.K., J.W. K.H., and M.Y.D. performed bioinformatic and genomic analysis. K.L.B., D.M.
729 and J.T. provided biomaterials and provided vole expertise. J.K., D.M., J.T., and M.Y.D. devised the
730 project. M.A., J.K., D.M., and J.T., and M.Y.D. wrote the manuscript and all authors edited and approved
731 the manuscript.

732

733 **COMPETING INTERESTS**

734 J.L. & C.L. are employees and shareholders and J.K. is a consultant and shareholder of Pacific
735 Biosciences, a company developing single-molecule sequencing technologies.

736

737 **REFERENCES**

- 738 1. Rhie, A., McCarthy, S.A., Fedrigo, O., Damas, J., Formenti, G., Koren, S., Uliano-Silva, M., Chow,
739 W., Fungtammasan, A., Kim, J., et al. (2021). Towards complete and error-free genome assemblies
740 of all vertebrate species. *Nature* 592, 737–746.
- 741 2. Dennis, M.Y., and Eichler, E.E. (2016). Human adaptation and evolution by segmental duplication.
742 *Curr. Opin. Genet. Dev.* 41, 44–52.
- 743 3. Wenger, A.M., Peluso, P., Rowell, W.J., Chang, P.-C., Hall, R.J., Concepcion, G.T., Ebler, J.,
744 Fungtammasan, A., Kolesnikov, A., Olson, N.D., et al. (2019). Accurate circular consensus long-
745 read sequencing improves variant detection and assembly of a human genome. *Nat. Biotechnol.* 37,

- 746 1155–1162.
- 747 4. Cheng, H., Concepcion, G.T., Feng, X., Zhang, H., and Li, H. (2021). Haplotype-resolved de novo
748 assembly using phased assembly graphs with hifiasm. *Nat Methods* 18, 170–175.
- 749 5. Lieberman-Aiden, E., van Berkum, N.L., Williams, L., Imakaev, M., Ragozcy, T., Telling, A., Amit,
750 I., Lajoie, B.R., Sabo, P.J., Dorschner, M.O., et al. (2009). Comprehensive mapping of long-range
751 interactions reveals folding principles of the human genome. *Science* 326, 289–293.
- 752 6. Burton, J.N., Adey, A., Patwardhan, R.P., Qiu, R., Kitzman, J.O., and Shendure, J. (2013).
753 Chromosome-scale scaffolding of de novo genome assemblies based on chromatin interactions. *Nat.*
754 *Biotechnol.* 31, 1119–1125.
- 755 7. McGinty, S.P., Kaya, G., Sim, S.B., Makunin, A., Corpuz, R.L., Quail, M.A., Abuelanin, M.,
756 Lawniczak, M.K.N., Geib, S.M., Korlach, J., et al. (2025). CiFi: accurate long-read chromosome
757 conformation capture with low-input requirements. *Nat. Commun.* 17, 215.
- 758 8. Kleiman, D.G. (1977). Monogamy in mammals. *Q Rev Biol* 52, 39–69.
- 759 9. Lukas, D., and Clutton-Brock, T.H. (2013). The evolution of social monogamy in mammals. *Science*
760 341, 526–530.
- 761 10. Bales, K.L., Ardekani, C.S., Baxter, A., Karaskiewicz, C.L., Kuske, J.X., Lau, A.R., Savidge, L.E.,
762 Saylor, K.R., and Witzak, L.R. (2021). What is a pair bond? *Horm. Behav.* 136, 105062.
- 763 11. Stepan, S.J., and Schenk, J.J. (2017). Muroid rodent phylogenetics: 900-species tree reveals
764 increasing diversification rates. *PLoS One* 12, e0183070.
- 765 12. Robinson, K.J., Bosch, O.J., Levkowitz, G., Busch, K.E., Jarman, A.P., and Ludwig, M. (2019).
766 Social creatures: Model animal systems for studying the neuroendocrine mechanisms of social
767 behaviour. *J Neuroendocrinol* 31, e12807.
- 768 13. Elphick, M.R., Mirabeau, O., and Larhammar, D. (2018). Evolution of neuropeptide signalling
769 systems. *J Exp Biol* 221. <https://doi.org/10.1242/jeb.151092>.
- 770 14. Winslow, J.T., Hastings, N., Carter, C.S., Harbaugh, C.R., and Insel, T.R. (1993). A role for central
771 vasopressin in pair bonding in monogamous prairie voles. *Nature* 365, 545–548.
- 772 15. Lim, M.M., Wang, Z., Olazábal, D.E., Ren, X., Terwilliger, E.F., and Young, L.J. (2004). Enhanced
773 partner preference in a promiscuous species by manipulating the expression of a single gene. *Nature*
774 429, 754–757.
- 775 16. Carter, C.S., DeVries, A.C., and Getz, L.L. (1995). Physiological substrates of mammalian
776 monogamy: the prairie vole model. *Neurosci Biobehav Rev* 19, 303–314.
- 777 17. Sadino, J.M., and Donaldson, Z.R. (2018). Prairie voles as a model for understanding the genetic and
778 epigenetic regulation of attachment behaviors. *ACS Chem. Neurosci.* 9, 1939–1950.
- 779 18. Shapiro, L.E., and Insel, T.R. (1990). Infant's response to social separation reflects adult differences
780 in affiliative behavior: a comparative developmental study in prairie and montane voles. *Dev.*
781 *Psychobiol.* 23, 375–393.
- 782 19. Young, L.J., Nilsen, R., Waymire, K.G., MacGregor, G.R., and Insel, T.R. (1999). Increased

- 783 affiliative response to vasopressin in mice expressing the V1a receptor from a monogamous vole.
784 *Nature* 400, 766–768.
- 785 20. McGraw, L.A., Davis, J.K., Thomas, P.J., NISC Comparative Sequencing Program, Young, L.J., and
786 Thomas, J.W. (2012). BAC-based sequencing of behaviorally-relevant genes in the prairie vole.
787 *PLoS One* 7, e29345.
- 788 21. Zhou, C., McCarthy, S.A., and Durbin, R. (2023). YaHS: yet another Hi-C scaffolding tool.
789 *Bioinformatics* 39, btac808.
- 790 22. Rhie, A., Walenz, B.P., Koren, S., and Phillippy, A.M. (2020). Merqury: reference-free quality,
791 completeness, and phasing assessment for genome assemblies. *Genome Biol* 21, 245.
- 792 23. Hartke, G.T., Leipold, H.W., Huston, K., Cook, J.E., and Saperstein, G. (1974). Three mutations and
793 the karyotype of the prairie vole. White spotting, polydipsia, and muscular dystrophy in *Microtus*
794 *ochrogaster*. *J Hered* 65, 301–307.
- 795 24. McGraw, L.A., Davis, J.K., Young, L.J., and Thomas, J.W. (2011). A genetic linkage map and
796 comparative mapping of the prairie vole (*Microtus ochrogaster*) genome. *BMC Genet.* 12, 60.
- 797 25. Reich, L.M. (1981). *Microtus pennsylvanicus*. *Mammalian Species*, 1.
- 798 26. Schmid, W., and Leppert, M.F. (1968). [Karyotype, heterochromatin and DNA-content in 13 species
799 of voles (Microtinae, Mammalia-Rodentia)]. *Arch Julius Klaus Stift Vererbungsforsch*
800 *Sozialanthropol Rassenhyg* 43, suppl 88–91.
- 801 27. Bailey, J.A., Gu, Z., Clark, R.A., Reinert, K., Samonte, R.V., Schwartz, S., Adams, M.D., Myers,
802 E.W., Li, P.W., and Eichler, E.E. (2002). Recent segmental duplications in the human genome.
803 *Science* 297, 1003–1007.
- 804 28. Išerić, H., Alkan, C., Hach, F., and Numanagić, I. (2022). Fast characterization of segmental
805 duplication structure in multiple genome assemblies. *Algorithms Mol Biol* 17, 4.
- 806 29. Fredga, K. (1988). Aberrant chromosomal sex-determining mechanisms in mammals, with special
807 reference to species with XY females. *Philos Trans R Soc Lond B Biol Sci* 322, 83–95.
- 808 30. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P., and
809 Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biol.* 17, 122.
- 810 31. Ohno, S. (1970). *Evolution by gene duplication* (Allen & Unwin; Springer-Verlag).
- 811 32. Qi, J., Ma, J., Han, Z., Han, R., Yu, T., and Li, G. (2025). De novo annotation of centromere with
812 centroAnno. *bioRxiv*. <https://doi.org/10.1101/2025.02.19.639205>.
- 813 33. Schultz, D.T., Haddock, S.H.D., Bredeson, J.V., Green, R.E., Simakov, O., and Rokhsar, D.S.
814 (2023). Ancient gene linkages support ctenophores as sister to other animals. *Nature* 618, 110–117.
- 815 34. Simakov, O., Bredeson, J., Berkoff, K., Marletaz, F., Mitros, T., Schultz, D.T., O’Connell, B.L.,
816 Dear, P., Martinez, D.E., Steele, R.E., et al. (2022). Deeply conserved synteny and the evolution of
817 metazoan chromosomes. *Sci. Adv.* 8, eabi5884.
- 818 35. Parsons, J.D. (1995). Miropeats: graphical DNA sequence comparisons. *Comput Appl Biosci* 11,
819 615–619.

- 820 36. Pendleton, A.L., Shen, F., Taravella, A.M., Emery, S., Veeramah, K.R., Boyko, A.R., and Kidd, J.M.
821 (2018). Comparison of village dog and wolf genomes highlights the role of the neural crest in dog
822 domestication. *BMC Biol* *16*, 64.
- 823 37. FitzGerald, R.W., and Madison, D.M. (1983). Social organization of a free-ranging population of
824 pine voles, *Microtus pinetorum*. *Behav. Ecol. Sociobiol.* *13*, 183–187.
- 825 38. Oliveras, D., and Novak, M. (1986). A comparison of paternal behaviour in the meadow vole
826 *Microtus pennsylvanicus*, the pine vole *M. pinetorum* and the prairie vole *M. chrogaster*. *Anim.*
827 *Behav.* *34*, 519–526.
- 828 39. Duckett, D.J., Sullivan, J., Pirro, S., and Carstens, B.C. (2021). Genomic Resources for the North
829 American Water Vole () and the Montane Vole (). *GigaByte* *2021*, gigabyte19.
- 830 40. Gouy, A., Wang, X., Kapopoulou, A., Neuenschwander, S., Schmid, E., Excoffier, L., and Heckel,
831 G. (2024). Genomes of *Microtus* Rodents Highlight the Importance of Olfactory and Immune
832 Systems in Their Fast Radiation. *Genome Biol Evol* *16*. <https://doi.org/10.1093/gbe/evae233>.
- 833 41. Wang, Z., Young, L.J., De Vries, G.J., and Insel, T.R. (1998). Voles and vasopressin: a review of
834 molecular, cellular, and behavioral studies of pair bonding and paternal behaviors. *Prog Brain Res*
835 *119*, 483–499.
- 836 42. Jannett, F.J. (1982). Nesting patterns of adult voles, *Microtus montanus*, in field populations. *J.*
837 *Mammal.* *63*, 495–498.
- 838 43. Madison, D.M., and McShea, W.J. (1987). Seasonal changes in reproductive tolerance, spacing, and
839 social organization in meadow voles: A microtine model. *Am. Zool.* *27*, 899–908.
- 840 44. Agrell, J. (1995). A shift in female social organization independent of relatedness: an experimental
841 study on the field vole (*Microtus agrestis*). *Behav. Ecol.* *6*, 182–191.
- 842 45. Schweizer, M., Excoffier, L., and Heckel, G. (2007). Fine-scale genetic structure and dispersal in the
843 common vole (*Microtus arvalis*). *Mol Ecol* *16*, 2463–2473.
- 844 46. Jeppsson, B. (1990). Effects of density and resources on the social system of water voles. In *Social*
845 *Systems and Population Cycles in Voles* (Birkhäuser Basel), pp. 213–226.
- 846 47. Fink, S., Excoffier, L., and Heckel, G. (2007). High variability and non-neutral evolution of the
847 mammalian *avpr1a* gene. *BMC Evol Biol* *7*, 176.
- 848 48. Duclot, F., Liu, Y., Saland, S.K., Wang, Z., and Kabbaj, M. (2022). Transcriptomic analysis of
849 paternal behaviors in prairie voles. *BMC Genomics* *23*, 679.
- 850 49. Tripp, J.A., Berrio, A., McGraw, L.A., Matz, M.V., Davis, J.K., Inoue, K., Thomas, J.W., Young,
851 L.J., and Phelps, S.M. (2021). Comparative neurotranscriptomics reveal widespread species
852 differences associated with bonding. *BMC Genomics* *22*, 399.
- 853 50. Duclot, F., Sailer, L., Koutakis, P., Wang, Z., and Kabbaj, M. (2022). Transcriptomic Regulations
854 Underlying Pair-bond Formation and Maintenance in the Socially Monogamous Male and Female
855 Prairie Vole. *Biol Psychiatry* *91*, 141–151.
- 856 51. Waddell, N.J., Liu, Y., Chitaman, J.M., Kaplan, G.J., Wang, Z., and Feng, J. (2023). Transcription
857 and DNA methylation signatures of paternal behavior in hippocampal dentate gyrus of prairie voles.

- 858 Sci Rep *13*, 11020.
- 859 52. Sadino, J.M., Bradeen, X.G., Kelly, C.J., Brusman, L.E., Walker, D.M., and Donaldson, Z.R. (2023).
860 Prolonged partner separation erodes nucleus accumbens transcriptional signatures of pair bonding in
861 male prairie voles. *Elife* *12*. <https://doi.org/10.7554/eLife.80517>.
- 862 53. Danoff, J.S., Carter, C.S., Gordevičius, J., Milčiūtė, M., Brooke, R.T., Connelly, J.J., and Perkeybile,
863 A.M. (2024). Maternal oxytocin treatment at birth increases epigenetic age in male offspring. *Dev*
864 *Psychobiol* *66*. <https://doi.org/10.1002/dev.22452>.
- 865 54. Karageorgiou, C., Gokcumen, O., and Dennis, M.Y. (2024). Deciphering the role of structural
866 variation in human evolution: a functional perspective. *Curr Opin Genet Dev* *88*, 102240.
- 867 55. Li, H., and Durbin, R. (2024). Genome assembly in the telomere-to-telomere era. *Nat Rev Genet* *25*,
868 658–670.
- 869 56. Wang, T., Antonacci-Fulton, L., Howe, K., Lawson, H.A., Lucas, J.K., Phillippy, A.M., Popejoy,
870 A.B., Asri, M., Carson, C., Chaisson, M.J.P., et al. (2022). The Human Pangenome Project: a global
871 resource to map genomic diversity. *Nature* *604*, 437–446.
- 872 57. Hansen, N.F., Dwarshuis, N., Ji, H.J., Rhie, A., Loucks, H., Logsdon, G.A., Vollger, M.R., Storer,
873 J.M., Kim, J., Adam, E., et al. (2025). A complete diploid human genome benchmark for
874 personalized genomics. *bioRxiv*. <https://doi.org/10.1101/2025.09.21.677443>.
- 875 58. Berendzen, K.M., and Manoli, D.S. (2022). Rethinking the Architecture of Attachment: New
876 Insights into the Role for Oxytocin Signaling. *Affect Sci* *3*, 734–748.
- 877 59. Naughton, C., Huidobro, C., Catacchio, C.R., Buckle, A., Grimes, G.R., Nozawa, R.-S., Purgato, S.,
878 Rocchi, M., and Gilbert, N. (2022). Human centromere repositioning activates transcription and
879 opens chromatin fibre structure. *Nat Commun* *13*, 5609.
- 880 60. Barbosa, S., Paupério, J., Pavlova, S.V., Alves, P.C., and Searle, J.B. (2018). The *Microtus* voles:
881 Resolving the phylogeny of one of the most speciose mammalian genera using genomics. *Mol*
882 *Phylogenet Evol* *125*, 85–92.
- 883 61. Liu, Z., Roesti, M., Marques, D., Hiltbrunner, M., Saladin, V., and Peichel, C.L. (2022).
884 Chromosomal fusions facilitate adaptation to divergent environments in threespine stickleback. *Mol.*
885 *Biol. Evol.* *39*. <https://doi.org/10.1093/molbev/msab358>.
- 886 62. Bonfield, J.K., Marshall, J., Danecek, P., Li, H., Ohan, V., Whitwham, A., Keane, T., and Davies,
887 R.M. (2021). HTSlib: C library for reading/writing high-throughput sequencing data. *Gigascience*
888 *10*, giab007.
- 889 63. Brown, M.R., Manuel Gonzalez de La Rosa, P., and Blaxter, M. (2025). Tidk: A toolkit to rapidly
890 identify telomeric repeats from genomic datasets. *Bioinformatics* *41*, btaf049.
- 891 64. Flynn, J.M., Hubley, R., Goubert, C., Rosen, J., Clark, A.G., Feschotte, C., and Smit, A.F. (2020).
892 RepeatModeler2 for automated genomic discovery of transposable element families. *Proc. Natl.*
893 *Acad. Sci. U. S. A.* *117*, 9451–9457.
- 894 65. Tarailo-Graovac, M., and Chen, N. (2009). Using RepeatMasker to identify repetitive elements in
895 genomic sequences. *Curr. Protoc. Bioinformatics Chapter 4*, 4.10.1–4.10.14.

- 896 66. Thibaud-Nissen, F., DiCuccio, M., Hlavina, W., Kimchi, A., Kitts, P.A., Murphy, T.D., Pruitt, K.D.,
897 and Souvorov, A. (2016). P8008 The NCBI Eukaryotic Genome Annotation Pipeline. *J. Anim. Sci.*
898 *94*, 184–184.
- 899 67. Patro, R., Duggal, G., Love, M.I., Irizarry, R.A., and Kingsford, C. (2017). Salmon provides fast and
900 bias-aware quantification of transcript expression. *Nat. Methods* *14*, 417–419.
- 901 68. Muzellec, B., Teleńczuk, M., Cabeli, V., and Andreux, M. (2023). PyDESeq2: a python package for
902 bulk RNA-seq differential expression analysis. *Bioinformatics* *39*, btad547.
- 903 69. Benson, G. (1999). Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids*
904 *Res.* *27*, 573–580.
- 905 70. Morgulis, A., Gertz, E.M., Schäffer, A.A., and Agarwala, R. (2006). WindowMasker: window-based
906 masker for sequenced genomes. *Bioinformatics* *22*, 134–141.
- 907 71. Emms, D.M., and Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for comparative
908 genomics. *Genome Biol.* *20*, 238.
- 909 72. Daily, J. (2016). Parasail: SIMD C library for global, semi-global, and local pairwise sequence
910 alignments. *BMC Bioinformatics* *17*, 81.
- 911 73. Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I., Hamelryck,
912 T., Kauff, F., Wilczynski, B., et al. (2009). Biopython: freely available Python tools for
913 computational molecular biology and bioinformatics. *Bioinformatics* *25*, 1422–1423.
- 914 74. Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L.
915 (2009). BLAST+: architecture and applications. *BMC Bioinformatics* *10*, 421.
- 916 75. Ranwez, V., Douzery, E.J.P., Cambon, C., Chantret, N., and Delsuc, F. (2018). MACSE v2: Toolkit
917 for the alignment of Coding Sequences accounting for frameshifts and stop codons. *Mol. Biol. Evol.*
918 *35*, 2582–2584.
- 919 76. Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version 7:
920 improvements in performance and usability. *Mol. Biol. Evol.* *30*, 772–780.
- 921 77. Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., Whitwham, A., Keane,
922 T., McCarthy, S.A., Davies, R.M., et al. (2021). Twelve years of SAMtools and BCFtools.
923 *Gigascience* *10*, giab008.
- 924 78. Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* *34*, 3094–
925 3100.
- 926 79. Plessy, C., Mansfield, M.J., Bliznina, A., Masunaga, A., West, C., Tan, Y., Liu, A.W., Grašič, J., Del
927 Río Pisula, M.S., Sánchez-Serna, G., et al. (2024). Extreme genome scrambling in marine planktonic
928 *Oikopleura dioica* cryptic species. *Genome Res.* *34*, 426–440.