

# 1 The evolutionary fates of a large segmental duplication in mouse

2 Andrew P Morgan<sup>1</sup>, J Matthew Holt<sup>2</sup>, Rachel C McMullan<sup>1</sup>, Timothy A Bell<sup>1</sup>, Amelia M-F Clayshulte<sup>1</sup>,  
3 John P Didion<sup>1</sup>, Liran Yadgary<sup>1</sup>, David Thybert<sup>3</sup>, Duncan T Odom<sup>4,5</sup>, Paul Flicek<sup>3,5</sup>, Leonard McMillan<sup>2</sup>,  
4 Fernando Pardo-Manuel de Villena<sup>1</sup>

5 <sup>1</sup> Department of Genetics, Carolina Center for Genome Sciences and Lineberger Comprehensive Cancer  
6 Center, University of North Carolina, Chapel Hill, NC

7 <sup>2</sup> Department of Computer Science, University of North Carolina, Chapel Hill, NC

8 <sup>3</sup> European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome  
9 Campus, Hinxton, Cambridge, CB10 1SD, United Kingdom

10 <sup>4</sup> University of Cambridge, Cancer Research UK Cambridge Institute, Robinson Way, Cambridge, CB2  
11 0RE, United Kingdom

12 <sup>5</sup> Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, United  
13 Kingdom

14 *Corresponding author:*

15 Fernando Pardo-Manuel de Villena

16 5049 Genetic Medicine Building

17 120 Mason Farm Road CB#7264

18 Chapel Hill, NC 27599-7264

## 19 ABSTRACT

20 Gene duplication and loss are major sources of genetic polymorphism in populations, and are important  
21 forces shaping the evolution of genome content and organization. We have reconstructed the origin and  
22 history of a 127 kbp segmental duplication, *R2d*, in the house mouse (*Mus musculus*). *R2d* contains a single  
23 protein-coding gene, *Cwc22*. *De novo* assembly of both the ancestral (*R2d1*) and the derived (*R2d2*) copies  
24 reveals that they have been subject to non-allelic gene conversion events spanning tens of kilobases. *R2d2*  
25 is also a hotspot for structural variation: its diploid copy number ranges from zero in the mouse reference  
26 genome to more than 80 in wild mice sampled from around the globe. Hemizyosity for high-copy-  
27 number alleles of *R2d2* is associated in *cis* with meiotic drive, suppression of meiotic crossovers, and  
28 copy-number instability, with a mutation rate in excess of 1 per 100 transmissions in laboratory  
29 populations. We identify an additional 57 loci covering 0.8% of the mouse genome with patterns of  
30 sequence variation similar to those at *R2d1* and *R2d2*. Our results provide a striking example of allelic  
31 diversity generated by duplication and demonstrate the value of *de novo* assembly in a phylogenetic  
32 context for understanding the mutational processes affecting duplicate genes.

33

## 34 INTRODUCTION

35 Duplication is an important force shaping the evolution of plant and animal genomes: it provides both a  
36 substrate for evolution and transient relief from selective pressure (Lynch and Conery 2000). Segmental  
37 duplications (SDs), defined as contiguous sequences which map to more than one physical position  
38 (Bailey and Eichler 2006), are a common feature of eukaryotic genomes and particularly those of  
39 vertebrates.

40 Like any sequence variant, a duplication first arises in a single individual in a population. The distinction  
41 between such copy-number variants (CNVs) and SDs is fluid and somewhat arbitrary: tracts of SDs are  
42 highly polymorphic in populations in species from *Drosophila* (Dopman and Hartl 2007) to mouse (She *et*  
43 *al.* 2008) to human (Bailey and Eichler 2006). Studies of parent-offspring transmissions have shown that  
44 SDs are prone to recurrent *de novo* mutations including some implicated in human disease (reviewed in  
45 Stankiewicz and Lupski 2002). Bursts of segmental duplication have preceded dramatic species radiations  
46 in primates. More broadly, blocks of conserved synteny in mammals frequently terminate at SDs (Bailey  
47 and Eichler 2006), suggesting that SDs could mediate the chromosomal rearrangements through which  
48 karyotypes diverge and reproductive barriers arise.

49 Notwithstanding their evolutionary importance, SDs are difficult to analyze. Repeated sequences with  
50 period longer than the insert size in a sequencing library and high pairwise similarity are likely to be  
51 collapsed into a single sequence during genome assembly. Efficient and sensitive alignment of high-  
52 throughput sequencing reads to duplicated sequence remains challenging (Treangen and Salzberg 2011).  
53 Genotyping of sites within SDs is difficult because variants between copies (paralogous variants) are  
54 easily confounded with variants within copies between individuals at a given copy (allelic variants).  
55 Latent paralogous variation may bias interpretations of sequence diversity and haplotype structure  
56 (Hurles 2002).

57 Paralogy also complicates phylogenetic inference. Ancestral duplication followed by differential losses  
58 along separate lineages may yield a local phylogeny that is discordant with the genome-wide phylogeny  
59 (Goodman *et al.* 1979). Within each duplicate copy, local phylogenies for adjacent intervals may also be  
60 discordant due to non-allelic gene conversion between copies (Dover 1982; Nagylaki and Petes 1982). As  
61 a result, over some fraction of the genome, sequences from individuals of the same species may be more  
62 closely related to sequences from individuals of an outgroup species than they are to each other.

63 In this manuscript we present a detailed analysis of a segmental duplication, *R2d*, in the house mouse  
64 (*Mus musculus*). *R2d* is a 127 kbp unit which contains the protein-coding gene *Cwc22* and flanking  
65 intergenic sequence. Although the C57BL/6J reference strain and other classical laboratory strains have a  
66 single haploid copy of the *R2d* sequence (in the *R2d1* locus), the wild-derived CAST/EiJ, ZALLENDE/EiJ,  
67 and WSB/EiJ strains have an additional 1, 16 and 33 haploid copies respectively in the *R2d2* locus. *R2d2* is  
68 the responder locus in a recently-described meiotic drive system on mouse chromosome 2 but is absent  
69 from the mouse reference genome (Waterston *et al.* 2002; Didion *et al.* 2015, 2016). We draw on a collection  
70 of species from the genus *Mus* sampled from around the globe to reconstruct the sequence of events  
71 giving rise to the locus' present structure (**Figure 1**). Using novel computational tools built around  
72 indexes of raw high-throughput sequencing reads, we perform local *de novo* assembly of phased  
73 haplotypes and explore patterns of sequence diversity within and between copies of *R2d*.

74 Both phylogenetic analyses and estimation of mutation rate in laboratory mouse populations reveal that  
75 *R2d2* and its surrounding region on chromosome 2 are unstable in copy number. Cycles of duplication,  
76 deletion and non-allelic gene conversion have led to complex phylogenetic patterns discordant with

77 species-level relationships within *Mus* which cannot be explained by known patterns of introgression  
78 between *Mus* species (Bonhomme *et al.* 2007; Yang *et al.* 2011).

79 Finally, we identify 57 other loci, covering 0.8% of the mouse genome, which share the key features of  
80 *R2d2*: elevated local sequence divergence; low recombination rate; and enrichment for segmental  
81 duplications. Previous studies of sequence variation in the mouse (Keane *et al.* 2011) have attributed this  
82 pattern to sorting of alleles segregating in the common ancestor of *M. musculus* and its sister species. We  
83 suggest instead that these loci have been subject to independent cycles of duplication and loss along *Mus*  
84 lineages. Marked enrichment for odorant, pheromone, and antigen-recognition receptors supports a role  
85 for balancing selection on the generation and maintenance of the extreme level of polymorphism  
86 observed at these loci.

## 87 RESULTS

### 88 *R2d* was duplicated in the common ancestor of *M. musculus* and *M. spretus*

89 In order to determine when the *R2d* CNV arose, we used quantitative PCR and/or depth of coverage in  
90 whole-genome sequencing to assay *R2d* copy number in a collection of samples spanning the phylogeny  
91 of the genus *Mus*. Samples were classified as having diploid copy number 2 (two chromosomes each with  
92 a single copy of *R2d*) or >2 (at least one chromosome with an *R2d* duplication).

93 We find evidence for >2 diploid copies in representatives of all mouse taxa tested from the Palearctic  
94 clade (Suzuki *et al.* 2004) (**Figure 1** and **Supplementary Table 1**): 236 of 525 *Mus musculus*, 1 of 1 *M.*  
95 *macedonicus*, 1 of 1 *M. spicilegus*, 1 of 1 *M. cypricus* and 8 of 8 *M. spretus* samples. However, we find no  
96 evidence of duplication in species from the southeast Asian clade, which is an outgroup to Palearctic  
97 mice: 0 of 2 *M. famulus*, 0 of 2 *M. fragilicauda*, 0 of 1 *M. cervicolor*, 0 of 1 *M. cookii* and 0 of 1 *M. caroli*  
98 samples. Outside the subgenus *Mus*, we found evidence for >2 diploid copies in none of the 9 samples  
99 tested from subgenus *Pyromys*. We conclude that the *R2d* duplication most likely occurred between the  
100 divergence of southeast Asian from Palearctic mice (~3.5 million years ago [Mya]) and the divergence of  
101 *M. musculus* from *M. spretus* (~2 Mya) (Suzuki *et al.* 2004; Chevret *et al.* 2005), along the highlighted  
102 branch of the phylogeny in **Figure 1A**. If the *R2d* duplication is ancestral to the divergence of *M. musculus*,  
103 then extant lineages of house mice which have 2 diploid copies of *R2d* — including the reference strain  
104 C57BL/6J (of predominantly *M. musculus domesticus* origin (Yang *et al.* 2007)) — represent subsequent  
105 losses of an *R2d* copy.

106 Duplication of the ancestral *R2d* sequence resulted in two paralogs residing in loci which we denote *R2d1*  
107 and *R2d2* (**Figure 1B**). Only one of these is present in the mouse reference genome, at chr2: 77.87 Mbp; the  
108 other copy maps approximately 6 Mbp distal (Didion *et al.* 2015), as we describe in more detail below.  
109 The more proximal copy, *R2d1*, lies in a region of conserved synteny with rat, rabbit, chimpanzee and  
110 human (Muffato *et al.* 2010) (**Supplementary Figure 1**); we conclude that it is the ancestral copy.

111 The sequence of the *R2d2* paralog was assembled *de novo* from whole-genome sequence reads (Keane *et al.*  
112 2011) from the strain WSB/EiJ (of pure *M. m. domesticus* origin (Yang *et al.* 2011)), which has diploid *R2d*  
113 copy number ~68 (Didion *et al.* 2015). We exploited the difference in depth of coverage for *R2d1* (1 haploid  
114 copy) and *R2d2* (33 haploid copies) to assign variants to *R2d1* or *R2d2*. Pairwise alignment of the *R2d2*  
115 contig against *R2d1* is shown in **Supplementary Figure 2**. The paralogs differ by at least 8 transposable-  
116 element (TE) insertions: 7 LINE elements specific to *R2d1* and 1 endogenous retroviral element (ERV)  
117 specific to *R2d2* (**Supplementary Table 2**). (Due to the inherent limitations of assembling repetitive

118 elements from short reads, it is likely that we have underestimated the number of young TEs in *R2d2*.  
119 The *R2d1*-specific LINEs are all < 2% diverged from the consensus for their respective families in the  
120 RepeatMasker database (<http://www.repeatmasker.org/cgi-bin/WEBRepeatMasker>), consistent with  
121 insertion within the last 2 My. The oldest *R2d2*-specific ERV we could detect is 0.7% diverged from its  
122 family consensus. TE insertions occurring since the ancestral *R2d* duplication are almost certainly  
123 independent, so these data are consistent with duplication <2 Mya. The *R2d* unit, minus paralog-specific  
124 TE insertions, is 127 kbp in size. *R2d* units in the *R2d2* locus are capped on both ends by (CTCC)<sub>n</sub>  
125 microsatellite sequences, and no read pairs spanning the breakpoint between *R2d2* and flanking sequence  
126 were identified.

127 In order to obtain a more precise estimate of the molecular age of the duplication event we assembled *de*  
128 *novo* a total of 16.9 kbp of intergenic and intronic sequence in 8 regions across the *R2d* unit from diverse  
129 samples and constructed phylogenetic trees. The trees cover 17 *R2d1* or *R2d2* haplotypes, 13 from inbred  
130 strains and 4 from wild mice. The sequence of *Mus caroli* (diploid copy number 2) is used as an outgroup.  
131 A concatenated tree is shown in **Figure 1C**. Using  $5.0 \pm 1.0$  million years before present (Mya) as the  
132 estimated divergence date for *M. caroli* and *M. musculus* (Suzuki *et al.* 2004; Chevret *et al.* 2005), Bayesian  
133 phylogenetic analysis with BEAST v1.8 (Drummond *et al.* 2012) yields 1.6 Mya (95% HPD 0.7 – 5.1 Mya)  
134 as the estimated age of the duplication event that gave rise to *R2d1* and *R2d2*. Although the assumption  
135 of a uniform molecular clock may not be strictly fulfilled for *R2d1* and *R2d2*, the totality of evidence –  
136 from presence/absence data across the mouse phylogeny, paralog-specific TE insertions, and sequence  
137 divergence between paralogs – strongly supports the conclusion that *R2d* was first duplicated within the  
138 last 2 My in the common ancestor of *M. musculus* and *M. spretus*.

139 For clarity, **Figure 1D** illustrates diploid copy number states that will be referenced in the remainder of  
140 the manuscript. Hereafter we refer to diploid copy numbers except when discussing inbred strains  
141 (which are effectively haploid).

#### 142 ***R2d* contains the essential gene *Cwc22***

143 The *R2d* unit encompasses one protein-coding gene, *Cwc22*, which encodes an essential mRNA splicing  
144 factor (Yeh *et al.* 2010). The gene is conserved across eukaryotes and is present in a single copy in most  
145 non-rodent species represented in the TreeFam database (<http://www.treefam.org/family/TF300510> (Li  
146 2006)). Five groups of *Cwc22* paralogs are present in mouse genomes: the copies in *R2d1* (*Cwc22<sup>R2d1</sup>*) and  
147 *R2d2* (*Cwc22<sup>R2d2</sup>*) plus retrotransposed copies in one locus at chr2: 83.9 Mbp and at two loci on the X  
148 chromosome (**Figure 2A**).

149 The three retrogenes are located in regions with no sequence similarity to each other, indicating that each  
150 represents an independent retrotransposition event. The copy on chr2 was subsequently expanded by  
151 further segmental duplication and now exists (in the reference genome) in 7 copies with >99.9% mutual  
152 similarity. The two retrotransposed copies on chrX are substantially diverged from the parent gene (<  
153 90% sequence similarity), lack intact open reading frames (ORFs), have minimal evidence of expression  
154 among GenBank cDNAs, and are annotated as likely pseudogenes (Pei *et al.* 2012). We therefore restricted  
155 our analyses to the remaining three groups of *Cwc22* sequences, all on chr2.

156 The canonical transcript of *Cwc22<sup>R2d1</sup>* (ENSMUST00000065889) is encoded by 21 exons on the negative  
157 strand. The coding sequence begins in the third exon and ends in the terminal exon (**Figure 2B**). Six of the  
158 seven protein-coding *Cwc22<sup>R2d1</sup>* transcripts in Ensembl v83 use this terminal exon, while one transcript  
159 (ENSMUST0000011824) uses an alternative terminal exon. Alignment of the retrogene sequence  
160 (ENSMUST00000178960) to the reference genome demonstrates that the retrogene captures the last 19

161 exons of the canonical transcript — that is, the 19 exons corresponding to the coding sequence of the  
162 parent gene.

### 163 **Copy number at *R2d2* is highly polymorphic in *M. musculus***

164 We previously demonstrated that haploid copy number of *R2d* ranges from 1 in the reference strain  
165 C57BL/6J and classical inbred strains A/J, 129S1/SvImJ, NOD/ShiLtJ, NZO/HILtJ; to 2 in the wild-derived  
166 strain CAST/EiJ; to 34 in the wild-derived strain WSB/EiJ. Using linkage mapping in two multiparental  
167 genetic reference populations, the Collaborative Cross (Collaborative Cross Consortium 2012) and  
168 Diversity Outbred (Svenson *et al.* 2012), we showed that, for the two strains with haploid copy number  
169 >1, one of the copies maps to *R2d1* while all extra copies map to the *R2d2* locus at chr2: 83 Mbp (Didion *et*  
170 *al.* 2015). *Cwc22* was recently reported to have diploid copy number as high as 83 in wild *M. m. domesticus*  
171 (Pezer *et al.* 2015). In whole-genome sequence data from more than sixty mice from both laboratory stocks  
172 and natural populations (**Supplementary Table 1**), we have observed zero instances in which the *R2d*  
173 copy in *R2d1* is lost. We conclude that diploid copy number >2 indicates at least one copy of *R2d* is  
174 present in *R2d2* (**Figure 1D**).

175 In order to understand the evolutionary dynamics of copy-number variation at *R2d2*, we investigated the  
176 relationship between copy number and the local phylogeny in the *R2d2* candidate region. In particular,  
177 we sought evidence for or against a single common origin for each of the copy-number states at *R2d2*  
178 which are derived with respect to the *M. spretus* – *M. musculus* common ancestor (**Figure 1D**). If a derived  
179 copy-number state has a single recent origin, it should be associated with a single haplotype at *R2d2*. If a  
180 derived copy-number state arises by recurrent mutation, the same copy number should be associated  
181 with multiple haplotype backgrounds and possibly in multiple populations.

182 The extent of *R2d* copy-number variation in *M. musculus*, as estimated on a continuous scale by qPCR, is  
183 shown in **Figure 3A**. (Note that the qPCR readout is proportional to copy number on the log scale.  
184 Extrapolation to integer copy number is imprecise for copy numbers greater than ~6.) We confirmed that  
185 *R2d2* maps to chr2: 83 Mbp by performing association mapping between SNP genotypes from the  
186 MegaMUGA array (Morgan *et al.* 2016) and the qPCR readout (**Figure 3B**).

187 To test the hypothesis that losses of *R2d2* (diploid copy number < 4, at least one chromosome with zero  
188 copies in *R2d2*, **Figure 1D**) have a single origin, we examined their distribution across the three well-  
189 differentiated subspecies of *M. musculus*. Losses of *R2d2* occur in all subspecies of *M. musculus*, in  
190 populations that span its geographic range (**Supplementary Table 3**). Based on this distribution and a  
191 previous observation that no common haplotype is shared in samples with low copy number in *M. m.*  
192 *domesticus* (Didion *et al.* 2016), we reject the hypothesis of single origin and conclude that *R2d2* has been  
193 lost multiple times on independent lineages in each subspecies.

194 We performed a similar analysis to test the hypothesis that *R2d2* alleles with high copy number (diploid  
195 copy number >4, **Figure 1D**; hereafter “*R2d2<sup>HC</sup>*”) have a single origin. First we observed that *R2d2<sup>HC</sup>*  
196 alleles are confined with few exceptions to *M. m. domesticus* (**Supplementary Table 3**). The best-  
197 associated SNP on the MegaMUGA array (JAX00494952) only weakly tags copy number ( $r^2 = 0.137$ ), but  
198 ascertainment bias on the MUGA platform (Morgan *et al.* 2016) makes local LD patterns difficult to  
199 interpret. To examine further, we constructed a neighbor-joining phylogenetic tree for the region  
200 containing *R2d2* (chr2: 83 – 84 Mb) using genotypes from the 600K-SNP Mouse Diversity Array (Yang *et*  
201 *al.* 2011). We restricted our attention to inbred strains or wild mice with homozygous, non-recombinant  
202 haplotypes in the target region. Twelve samples with *R2d2<sup>HC</sup>* alleles, both wild mice and laboratory

203 stocks, cluster in a single clade (**Figure 3C**). (A single *M. spretus* strain, SPRET/EiJ, also carries an *R2d2<sup>HC</sup>*  
204 allele, but see **Discussion**).

205 Next we expanded the analysis to include an additional 11 samples with *R2d2<sup>HC</sup>* alleles and evidence of  
206 heterozygosity around *R2d2*. The total set of 24 samples includes 7 wild-derived laboratory strains  
207 (DDO, RBA/DnJ, RBB/DnJ, RBF/DnJ, WSB/EiJ, ZALENDE/EiJ and SPRET/EiJ), 4 classical inbred strains  
208 (ALS/LtJ, ALR/LtJ, CHMU/LeJ and NU/J), a line derived from the ICR:HsD outbred stock (HR8; Swallow  
209 *et al.* 1998) and 12 wild-caught mice. All 24 samples with *R2d2<sup>HC</sup>* alleles share an identical haplotype  
210 across a single 21 kbp interval, chr2: 83,896,447 – 83,917,565 (GRCm38/mm10 coordinates) (**Figure 3D**).  
211 These analyses support a single origin for *R2d2<sup>HC</sup>* alleles within *M. m. domesticus*.

## 212 *Cwc22* is intact in and expressed from all *R2d* paralogs, and fast-evolving in rodents

213 To identify the coding sequence of *Cwc22<sup>R2d2</sup>* we first aligned the annotated transcript sequences of  
214 *Cwc22<sup>R2d1</sup>* from Ensembl to our *R2d2* contig. All 21 exons present in *R2d1* are present in *R2d2*. We created a  
215 multiple sequence alignment and phylogenetic tree of *Cwc22* cDNAs and predicted amino acid sequences  
216 from *Cwc22<sup>R2d1</sup>*, *Cwc22<sup>R2d2</sup>*, retro-*Cwc22*, and *Cwc22* orthologs in 19 other placental mammals, plus  
217 opossum, platypus and finally chicken as an outgroup (**Supplementary Figure 3**). An open reading frame  
218 (ORF) is maintained in all three *Cwc22* loci in mouse, including the retrogene. Information content of each  
219 column along the alignment (**Supplementary Figure 4**) reveals that sequence is most conserved in two  
220 predicted conserved domains, MIF4G and MA3, required for *Cwc22*'s function in mRNA processing (Yeh  
221 *et al.* 2010).

222 Next we examined public RNA-seq data from adult brain and testis in inbred strains with one or more  
223 copies of *R2d2* for evidence of transcription of each *Cwc22* family member. We identified several novel  
224 transcript isoforms specific to *R2d2* arising from two intron-retention events and one novel 3' exon  
225 (**Figure 4A**). The 18<sup>th</sup> intron is frequently retained in *Cwc22<sup>R2d2</sup>* transcripts, most likely due to an A>G  
226 mutation at the 5' splice donor site of exon 17 in *Cwc22<sup>R2d2</sup>*. The 12<sup>th</sup> intron is also frequently retained.  
227 While we could not identify any splice-region variants near this intron, it contains an ERV insertion that  
228 may interfere with splicing (**Figure 4A**). Both intron-retention events would create an early stop codon.  
229 Finally, we find evidence for a novel 3' exon that extends to the boundary of the *R2d* unit and is used  
230 exclusively by *Cwc22<sup>R2d2</sup>* (**Figure 4A**).

231 We estimated the expression of the various isoforms of *Cwc22<sup>R2d1</sup>*, *Cwc22<sup>R2d2</sup>* and retro-*Cwc22* in adult  
232 brain and testis. For brain we obtained reads from 8 replicates (representing both sexes) on 3 inbred  
233 strains, and for testis a single replicate on 23 inbred strains and estimated transcript abundance using the  
234 kallisto package (Bray *et al.* 2015). Briefly, kallisto uses an expectation-maximization (EM) algorithm to  
235 accurately estimate the abundance of a set of transcripts by distributing the "weight" of each read across  
236 all isoforms with whose sequence it is compatible. *Cwc22* is clearly expressed from all three paralogs in  
237 both brain and testis (**Figure 4B**). However, both the total expression and the pattern of isoform usage  
238 differ by tissue and copy number.

239 Maintenance of an ORF in all *Cwc22* paralogs for >2 My is evidence of negative selection against  
240 disrupting mutations in the coding sequence, but long branches within the rodent clade in  
241 **Supplementary Figure 3** suggest that *Cwc22* may also be under relaxed purifying selection or positive  
242 selection in rodents. The rate of evolution of *Cwc22* sequences in mouse is faster than in the rest of the tree  
243 ( $\chi^2 = 4.33$ , df = 1,  $p = 0.037$  by likelihood ratio test).

## 244 **Phylogenetic discordance in *R2d1* is due to non-allelic gene conversion**

245 The topology of trees across *R2d* is generally consistent: a long branch separating the single *M. caroli*  
246 sequence from the *M. musculus* sequences, and two clades corresponding to *R2d1*- and *R2d2*-like  
247 sequences. However, we observed that the affinities of some *R2d* paralogs change along the sequence  
248 (**Figure 5A**), a signature of non-allelic (*i.e.* inter-locus) gene conversion. In this context, we use “gene  
249 conversion” to describe a non-reciprocal “copy-and-paste” transfer of sequence from one donor locus into  
250 a different, homologous receptor locus, without reference to a specific molecular mechanism (Chen *et al.*  
251 2007).

252 To investigate further, we inspected patterns of sequence variation in whole-genome sequencing data  
253 from 15 wild-caught mice, 2 wild-derived inbred strains, and 22 classical inbred strains of mice with  
254 diploid *R2d* copy number 2. We first defined 1,411 pairwise single-nucleotide differences between *R2d2*  
255 and *R2d1* for which *R2d2* has the derived allele with respect to *M. caroli*. Then we tested for the presence  
256 of the derived allele, ancestral allele or both at each site in each sample. Finally we identified conversion  
257 tracts by manual inspection as clusters of derived variants shared with *R2d2* (**Supplementary Figure 5**).

258 This analysis revealed non-allelic gene conversion tracts on at least 9 chromosomes (**Figure 5B**). The  
259 conversion tracts range in size from approximately 1.2 kbp to 119 kbp. The boundaries of several tracts  
260 are shared within populations, suggesting that they are identical by descent. We excluded the possibility  
261 of complementary losses from *R2d1* and *R2d2* – which would leave similar patterns of sequence variation  
262 – by finding read pairs spanning the boundary between *R2d1* and flanking sequence, and between *R2d1*-  
263 like and *R2d2*-like tracts on the same chromosome (examples shown in **Figure 5C**).

264 The conversion tracts we detected are orders of magnitude longer than the 15 to 750 bp reported in recent  
265 studies of allelic gene conversion at recombination hotspots in mouse meiosis (Cole *et al.* 2010, 2014). We  
266 require the presence of *R2d2*-diagnostic alleles at two or more consecutive variants to declare a  
267 conversion event, and these variants occur at a rate of approximately 1 per 100 bp, so the smallest  
268 conversion tracts we could theoretically detect are on the order of 200 bp in size. Even if we require only a  
269 single variant to define a conversion tract, all samples without a long conversion tract share fewer than 55  
270 and most fewer than 10 (of 1,411) derived alleles with *R2d2*, of which all are also shared by multiple other  
271 samples from different populations (**Supplementary Figure 5**). This pattern indicates that those sites in  
272 fact represent either artifacts (from mis-assignment of ancestral and derived alleles) or homoplasy rather  
273 than short gene conversions.

274 Four conversion tracts partially overlap the *Cwc22* gene to create a sequence that is a mosaic of *R2d1*- and  
275 *R2d2*-like exons (**Figure 5B**). Recovery of *Cwc22* mRNA in an inbred strain with a mosaic sequence  
276 (PWK/PhJ, see section “*Cwc22* is intact and expressed”) indicates that its exons are intact and properly  
277 oriented in *cis*. The presence of both *R2d1*- and *R2d2*-like sequence in extant *M. musculus* lineages with a 2  
278 diploid copies of *R2d* further reinforces our conclusion that the duplication is indeed ancestral to the  
279 divergence of *M. musculus*.

280 In addition to exchanges between *R2d1* and *R2d2*, we identified an instance of exchange between *R2d2*  
281 and the adjacent retrotransposed copy of *Cwc22* in a single *M. m. domesticus* individual from Iran  
282 (IR:AHZ\_STND:015; **Supplementary Figure 6**). This individual carries a rearrangement that has inserted  
283 a 30 kbp fragment corresponding to the 3’ half of *Cwc22<sup>R2d2</sup>* into the retro-*Cwc22* locus, apparently  
284 mediated by homology between the exons of *Cwc22<sup>R2d2</sup>* and retro-*Cwc22*.

## 285 **High copy number at *R2d2* suppresses meiotic recombination**

286 Based on a previous observation that the rate of meiotic recombination is reduced near clusters of  
287 segmental duplications (Liu *et al.* 2014), we tested whether the region around *R2d2* has lower  
288 recombination when an *R2d2<sup>HC</sup>* allele is present. Understanding patterns of recombination at *R2d2* is  
289 important for interpreting levels of sequence and haplotype diversity in the surrounding region.

290 First we analyzed local recombination rate in the Diversity Outbred (DO) population. The DO is an  
291 outbred stock derived from eight inbred founder strains (including one, WSB/EiJ, with an *R2d2<sup>HC</sup>* allele)  
292 and maintained by random mating with 175 breeding pairs; at each generation, one male and one female  
293 offspring are chosen from each mating and randomly paired with a non-sibling to produce the next  
294 generation (Svenson *et al.* 2012). **Figure 6A** shows the cumulative distribution of 2,917 recombination  
295 events on central chromosome 2, stratified according to *R2d2* copy number of the participating  
296 haplotypes. The recombination map has a pronounced plateau in the region between *R2d1* and  
297 approximately 1 Mb distal to *R2d2* (dashed lines) for *R2d2<sup>HC</sup>* haplotypes, but not *R2d2<sup>LC</sup>* haplotypes. As a  
298 result, *R2d2<sup>HC</sup>* haplotype blocks overlapping *R2d2* are significantly shorter than *R2d2<sup>LC</sup>* haplotype blocks  
299 ( $p < 0.01$  by Wilcoxon rank-sum tests with Bonferroni correction) in 8 of the 10 generations sampled  
300 (**Figure 6B**). The difference arose early in the breeding of the DO and persists through the most recent  
301 generation for the randomized breeding scheme was maintained (FPMdV, unpublished).

302 Second we re-examined genotype data from 11 published crosses in which at least one parent was  
303 segregating for an *R2d2<sup>HC</sup>* allele. Whereas in the DO we used haplotype block length as a proxy for  
304 recombination rate, in these F<sub>2</sub> and backcross designs we can directly estimate the recombination fraction  
305 across *R2d2* and compare it to its expected value in the absence of an *R2d2<sup>HC</sup>* allele (**Supplementary**  
306 **Figure 7**). In 9 of 11 crosses examined, the observed recombination fraction is lower than the expected  
307 ( $p < 0.032$ , one-sided binomial test).

## 308 **The genomic region containing *R2d2* is structurally unstable but has low sequence diversity**

309 The extent of copy-number polymorphism involving *R2d2* suggests that it is intrinsically unstable.  
310 Consistent with these observations, we find that the rate of *de novo* copy-number changes at *R2d2* is  
311 extremely high in laboratory populations (**Figure 7**). In 183 mice sampled from the DO population we  
312 identified and confirmed through segregation analysis 8 new alleles, each with distinct copy number and  
313 each occurring in an unrelated haplotype (**Supplementary Table 4**). Without complete pedigrees and  
314 genetic material from breeders a direct estimate of the mutation rate in the DO is not straightforward to  
315 obtain. However, since the population size is known, we can make an analogy to microsatellite loci  
316 (Moran 1975) and estimate the mutation rate via the variance in allele sizes: 3.2 mutations per 100  
317 transmissions (3.2%) (95% bootstrap CI 1.1% – 6.0%).

318 Structural instability in this region of chromosome 2 extends outside the *R2d2* locus itself. Less than 200  
319 kbp distal to *R2d2* is another segmental duplication (**Figure 8B**, grey shaded region) – containing a  
320 retrotransposed copy of *Cwc22* – that is present in 7 tandem copies in the reference genome. That region,  
321 plus a further 80 kbp immediately distal to it, is copy-number polymorphic in wild *M. m. domesticus* and  
322 wild *M. m. castaneus* (**Figure 8B**). Instability of the region over longer timescale is demonstrated by the  
323 disruption, just distal to the aforementioned segmental duplication, of a syntenic block conserved across  
324 all other mammals (**Supplementary Figure 1**).

325 Despite the high mutation rate for structural variants involving *R2d2* and nearby sequences, sequence  
326 diversity at the nucleotide level is modestly reduced relative to diversity in *R2d1* and relative to the  
327 genome-wide average in *M. m. domesticus*. In a 200 kbp region containing the *R2d2* insertion site at its



328 proximal end,  $\hat{\pi}$  (an estimator of average heterozygosity) in *M. m. domesticus* reduced from approximately  
329 0.3% (comparable to previous reports in this subspecies, (Salcedo *et al.* 2007)) to nearly zero (**Figure 8B**).  
330 Divergence between *M. musculus* and *M. caroli* is similar to its genome-wide average of  $\sim 2.5\%$  over the  
331 same region.

332 Estimation of diversity *within* a duplicated sequence such as *R2d* is complicated by the difficulty of  
333 distinguishing allelic from paralogous variation. To circumvent this problem we split our sample of 26  
334 wild *M. m. domesticus* into two groups: those having *R2d1* sequences only, and those having both *R2d1*  
335 and *R2d2* sequences. Within each group we counted the number of segregating sites among all *R2d2*  
336 copies, using nearby fixed differences between *R2d1* and *R2d2* to phase sites to *R2d2* (see **Methods** for  
337 details), and used Watterson's estimator to calculate nucleotide diversity per site. Among *R2d1*  
338 sequences,  $\theta = 0.09\% \pm 0.03\%$  versus  $\theta = 0.04\% \pm 0.02\%$  among *R2d2* sequences (**Figure 8C**) and  
339  $\theta = 0.13\% \pm 0.04\%$  among *R2d2* sequences in *M. m. castaneus*.

### 340 **The “revolving door” affects at least 0.8% of the mouse genome**

341 The superposition of duplication, non-allelic gene conversion, and loss — dubbed the “genomic  
342 revolving door” (Demuth *et al.* 2006) — can obscure the true history of multicopy sequences. Over short  
343 evolutionary times, duplication and loss may be further confounded with incomplete lineage sorting, or  
344 the stochastic distribution of ancestral polymorphisms along descendant lineages (Pamilo and Nei 1988).

345 Most previous studies of gene gain and loss have focused on between-species comparisons using finished  
346 genome assemblies (Demuth *et al.* 2006; Bailey and Eichler 2006). We sought evidence for the “revolving  
347 door” effect over the last  $\sim 1$  My of evolution within *M. musculus*. To do so we identified regions where  
348 the divergence between sequenced samples and the reference genome assembly substantially exceeds  
349 what is predicted by those samples' ancestry. We developed a simple and unbiased method to estimate  
350 divergence from short sequencing reads without alignment. Briefly, we estimated the proportion of short  
351 subsequences (*k*-mers, with  $k = 31$ ) from windows along the reference assembly for which no evidence  
352 exists in the reads generated for a sample. This quantity can be rescaled to approximate the divergence  
353 between the reference sequence and the template sequence from which reads were generated (see  
354 **Methods**). Applied genome-wide to 31 kbp ( $= 1000 \times k$ ) windows, this method captures the distribution  
355 of sequence divergence between the reference assembly (chiefly *M. m. domesticus* in origin) in  
356 representative samples from the three subspecies of *M. musculus* and the outgroups *M. spretus* and *M.*  
357 *caroli* (**Figure 9A**).

358 Divergent regions were identified by fitting a hidden Markov model (HMM) to the windowed  
359 divergence profiles for 7 wild or wild-derived samples with available whole-genome sequence (**Figure 9B**  
360 and **Supplementary Table 6**). This analysis revealed a striking pattern: over most of the genome, the  
361 divergence estimates hover around the genome-wide expectation, but high-divergence windows are  
362 clustered in regions 100 kbp to 5 Mbp in size. Our method does not capture signal from structural  
363 variation (besides deletions, see **Methods**), and so probably underestimates the true level of sequence  
364 divergence. In any case the union of these divergent regions across all seven *M. musculus* samples  
365 analyzed covers 0.82% of the reference genome (mean 0.58% per sample). Not surprisingly, divergent  
366 regions are enriched for segmental duplications: 39.0% of sequence in divergent regions is comprised of  
367 segmental duplications versus a median of 3.4% (central 99% interval 0.2% – 16.7%) in random regions  
368 of equal size ( $p < 0.001$ ). Divergent regions also have significant overlap with previously-identified (Liu  
369 *et al.* 2014) regions of low recombination ( $p < 0.001$ ). Yet divergent regions are more gene-dense than the  
370 genomic background: they contain  $3.6 \times 10^{-2}$  genes/kbp relative to a median  $1.6 \times 10^{-2}$  genes/kbp (central  
371 99% interval  $1.1 \times 10^{-2}$  –  $2.7 \times 10^{-2}$  genes/kbp) genome-wide ( $p = < 0.001$ ). Divergent regions are strongly

372 enriched for genes related to odorant and pheromone sensing ( $p = 1.3 \times 10^{-8}$ ) and adaptive immunity  
373 ( $p = 3.1 \times 10^{-3}$ ).

374 As a representative example we focus on a divergent region at chr4: 110 - 115 Mbp (**Figure 9B**). This  
375 region contains the 11 members of the *Skint* family of T-cell-borne antigen receptors. The first member to  
376 be described, *Skint1*, functions in negative selection in the thymus of T-cells destined for the epidermis  
377 (Boyden *et al.* 2008). Coding sequence from 23 inbred strains (including CAST/EiJ, WSB/EiJ, PWD/PhJ and  
378 SPRET/EiJ) was reported in Boyden *et al.* (2008). A phylogenetic tree constructed from those sequences,  
379 plus *M. caroli* and rat (not shown) as outgroups, reveals the expected pattern of “deep coalescence”  
380 (**Figure 9C**): the *M. m. domesticus* sequences are paraphyletic, and some (group A) are more similar to  
381 SPRET/EiJ (*M. spretus*) than to their subspecific congeners (group B). Although the level of sequence  
382 divergence in the *Skint1* coding sequence was attributed to ancestral polymorphism maintained by  
383 balancing selection in the original report, selection would not be expected to maintain diversity in both  
384 coding and non-coding sequence equally as we observe in **Figure 9B**. The best explanation for the  
385 observed pattern of diversity at *Skint1* is therefore that groups A and B represent “pseudo-orthologs”  
386 (Koonin 2005) descended from an ancestral duplication followed by subsequent deletion of different  
387 paralogs along different lineages (**Figure 9D**). This conclusion is supported by the structure of the *Skint*  
388 region in the mouse reference genome assembly, which reflects the superposition of many duplications  
389 and rearrangements (**Figure 9E**).

## 390 DISCUSSION

391 In this manuscript we have reconstructed in detail the evolution of a multi-megabase segmental  
392 duplication (SD) in mouse, *R2d2*. Our findings illustrate the challenges involved in accurately  
393 interpreting patterns of polymorphism and divergence within duplicated sequence.

394 SDs are among the most dynamic loci in mammalian genomes. They are foci for copy-number variation  
395 in populations, but the sequences of individual duplicates beyond those present in the reference genome  
396 are often poorly resolved. Obtaining the sequence of this “missing genome,” as we have done for *R2d2*, is  
397 an important prerequisite to understanding the evolution of duplicated loci. Since each paralog follows a  
398 partially independent evolutionary trajectory, individuals in a population may vary both quantitatively  
399 (in the number of copies) and qualitatively (in which copies are retained). Cycles of duplication and loss  
400 may furthermore lead to the fixation of different paralogs along different lineages. This “genomic  
401 revolving door” leaves a signature of polymorphism far in excess of the genome-wide background, due  
402 to coalescence between alleles originating from distinct paralogs. We identify 57 additional regions  
403 covering 0.82% of the mouse genome with this property (**Figure 9** and **Supplementary Table 6**). These  
404 regions have gene density similar to unique sequence and are strongly enriched for genes involved in  
405 odorant sensing, pheromone recognition and immunity that play important roles in social behavior and  
406 speciation (Hurst *et al.* 2001). Excess polymorphism at these loci in mouse has previously been attributed  
407 to some combination of incomplete lineage sorting and diversifying selection (White *et al.* 2009; Keane *et*  
408 *al.* 2011). Our results suggest that inferences regarding the strength of selection on highly polymorphic  
409 loci in regions of subject to recurrent duplication and loss should be treated with caution.

410 Accurate deconvolution of recent duplications remains a difficult task that requires painstaking manual  
411 effort. Clone-based and/or single-molecule long-read sequencing remain the gold standard techniques.  
412 But short reads at sufficient depth nonetheless contain a great deal of information. We exploited the  
413 specific properties of *R2d2* in the WSB/EiJ mouse strain — many highly-similar copies of *R2d2* relative to  
414 the single divergent *R2d1* copy — to obtain a nearly complete assembly of *R2d2* from short reads

415 (Supplementary Figure 8). With the sequence of both the *R2d1* and *R2d2* paralogs in hand, we were able  
416 to recognize several remarkable features of *R2d2* that are discussed in detail below.

#### 417 Long-tract gene conversion.

418 Previous studies of non-allelic gene conversion in mouse and human have focused either on relatively  
419 small (<5 kbp) intervals within species, or have applied phylogenetic methods to multiple paralogs from a  
420 single reference genome (Dumont and Eichler 2013). This study is the first, to our knowledge, with the  
421 power to resolve large (>5 kbp) non-allelic gene conversion events on an autosome in a population  
422 sample. We identify conversion tracts up to 119 kbp in length, orders of magnitude longer than tracts  
423 arising from allelic conversion events during meiosis. Gene conversion at this scale can rapidly and  
424 dramatically alter paralogous sequences, including — as shown in Figure 5 — the sequences of essential  
425 protein-coding genes. This process has been implicated as a source of disease alleles in humans (Chen *et*  
426 *al.* 2007).

427 Importantly, we were able to identify non-allelic exchanges in *R2d1* as such only because we were aware  
428 of the existence of *R2d2* in other lineages. In this case the transfer of paralogous *R2d2* sequence into *R2d1*  
429 creates the appearance of deep coalescence among *R2d1* sequences. Ignoring the effect of gene conversion  
430 would cause us to overestimate the degree of polymorphism at *R2d1* by an order of magnitude, and  
431 would bias any related estimates of population-genetic parameters (for instance, of effective population  
432 size).

433 Our data are not sufficient to estimate the rate of non-allelic gene conversion between *R2d2* and  
434 homologous loci. At minimum we have observed two distinct events: one from *R2d2* into *R2d1*, and a  
435 second from *R2d2* into retro-*Cwc22*. From a single conversion event replacing most of *R2d1* with *R2d2*-like  
436 sequence, the remaining shorter conversion tracts could be generated by recombination with *R2d1*  
437 sequences. Because we find converted haplotypes in both *M. m. musculus* and *M. m. domesticus*, the single  
438 conversion event would have had to occur prior to the divergence of the three *M. musculus* subspecies  
439 and subsequently remain polymorphic in the diverged populations. We note that all conversion tracts we  
440 observed are polarized: *R2d2* is always the donor.

441 The other possibility is that non-allelic gene conversion between *R2d* sequences is recurrent. Recurrent  
442 gene conversion homogenizes duplicate sequences, coupling their evolutionary trajectories (“concerted  
443 evolution”, Dover 1982). The absolute sequence divergence (~2%) between *R2d1* and *R2d2* (Figure 1B)  
444 argues against a hypothesis of ongoing gene conversion between these loci. However, we cannot rule out  
445 a role for gene conversion in maintaining sequence identity between multiple copies of *R2d* located in  
446 *R2d2*. This would help explain the reduced diversity within *R2d2* versus *R2d1* (Figure 8C). For  
447 sequences with copy number greater than two — such as the *R2d2* cassette in *R2d2<sup>HC</sup>* alleles — gene  
448 conversion tends to slow the accumulation of new mutations in a copy-number dependent manner  
449 (Nagyaki and Petes 1982). New mutations arising in any single copy are prone to loss not only by drift  
450 but also by being “pasted-over” by gene conversion from the intact copies which outnumber the mutant  
451 (Melamed and Kupiec 1992).

452 *R2d2* is not unlike the male-specific region of the Y chromosome in mouse (Soh *et al.* 2014) and human  
453 (Rozen *et al.* 2003). The large palindromic repeats on chrY are homogenized by frequent non-allelic gene  
454 conversion (Hallast *et al.* 2013) such that they have retained >99% sequence identity to each other even  
455 after millions of years of evolution. Frequent non-allelic gene conversion has also been documented in  
456 arrays of U2 snRNA genes in human (Liao 1997), and in rRNA gene clusters (Eickbush and Eickbush  
457 2007) and centromeric sequences (Schindelbauer 2002; Shi *et al.* 2010) in several species.

#### 458 **Pervasive copy-number variation.**

459 Clusters of segmental duplications have long been known to be hotspots of copy-number variation in  
460 populations (Bailey and Eichler 2006; She *et al.* 2008) and *de novo* mutations in pedigrees (Egan *et al.* 2007).  
461 Recent large-scale sequencing efforts have revealed the existence of thousands of multiallelic CNVs  
462 segregating in human populations (Handsaker *et al.* 2015).

463 We have surveyed *R2d2* copy number in a large and diverse sample of laboratory and wild mice, and  
464 have shown that it varies from 0 to >80 in certain *M. m. domesticus* populations (**Figure 8A**). In a cohort of  
465 outbred mice expected to be hemizygous for an *R2d2<sup>HC</sup>* allele from WSB/EiJ (33 diploid copies) we  
466 estimate that large deletions, >2 Mbp in size, occur at a rate of 3.2% (95% bootstrap CI 1.1% – 6.0%) per  
467 generation. This estimate of the mutation rate for CNVs at *R2d2* should be regarded as a lower bound.  
468 The power of our copy-number assay to discriminate between copy numbers above ~25 is low, so that the  
469 assay is much more sensitive to losses than to gains. Even our lower-bound mutation rate exceeds that of  
470 the most common recurrent deletions in human (~1 per 7000 live births) (Turner *et al.* 2007) and is an  
471 order of magnitude higher than the most active CNV hotspots described to date in the mouse (Egan *et al.*  
472 2007).

473 However, the structural mutation rate appears to depend strongly on the diplotype configuration at *R2d2*.  
474 As **Figure 1D** shows, individuals heterozygous for an *R2d2<sup>HC</sup>* haplotype and an *R2d2*-null haplotype are  
475 in fact hemizygous for several megabases of DNA in *R2d2*. This has important consequences. High  
476 mutation rates are observed only in the context of populations in which hemizygosity for *R2d2<sup>HC</sup>* is  
477 common (**Figure 7**): highest in the DO, and to a lesser extent in wild *M. m. domesticus* populations  
478 harboring both *R2d2<sup>HC</sup>* and *R2d2*-null alleles. Homozygosity for *R2d2<sup>HC</sup>* is not associated with mutability:  
479 in 8 recombinant inbred lines from the Collaborative Cross which are homozygous for an *R2d2<sup>HC</sup>*  
480 haplotype, we observed zero new mutations in at least 400 meioses, through both the male and female  
481 germline (8 lines × 2 meioses/generation × 25 or more generations of inbreeding). Sex also appears to  
482 have a role in determining the mutation rate at *R2d2*: in a pedigree in which all females were hemizygous  
483 for *R2d2<sup>HC</sup>*, zero new mutations were observed in 1256 meioses (data not shown).

484 Taken together, these observations hint at a common structural or epigenetic mechanism affecting the  
485 resolution of double-strand breaks in large tracts of unpaired (*i.e.* hemizygous) DNA during male  
486 meiosis. At least one other study in mouse has hinted that hemizygous SDs on the sex chromosomes are  
487 unstable in inter-subspecific hybrids (Scavetta and Tautz 2010). Both the obligate-hemizygous sex  
488 chromosomes and large unpaired segments on autosomes are epigenetically marked for transcriptional  
489 silencing during male meiotic prophase (Laan 2004; Baarends *et al.* 2005), and are physically sequestered  
490 into a structure called the sex body. Repair of double-strand breaks within the sex body is delayed  
491 relative to the autosomes (Mahadevaiah *et al.* 2001) and involves a different suite of proteins (Turner *et al.*  
492 2004). We hypothesize that these male-specific pathway(s) are generally error-prone in the presence of  
493 non-allelic homologous sequences.

#### 494 **Origin and distribution of an allele subject to meiotic drive.**

495 Females heterozygous for a high- and low-copy allele at *R2d2* preferentially transmit the high-copy allele  
496 to progeny via meiotic drive (Didion *et al.* 2015). Meiotic drive can rapidly alter allele frequencies in  
497 laboratory and natural populations (Lindholm *et al.* 2016), and we recently showed that high-copy alleles  
498 of *R2d2* (*R2d2<sup>HC</sup>*) sweep through laboratory and natural populations despite reducing the fitness of  
499 heterozygous females (Didion *et al.* 2016). These “selfish sweeps” account for the marked reduction in  
500 within-population diversity in the vicinity of *R2d2* (**Figure 8B**).

501 The present study sheds additional light on the age, origins and fate of *R2d2<sup>HC</sup>* alleles. We find that *R2d2<sup>HC</sup>*  
502 alleles have a single origin in *M. m. domesticus*. They are present in several different “chromosomal races”  
503 — populations fixed for specific Robertsonian translocations between which gene flow is limited (Hauffe  
504 and Searle 1993) — indicating that they were likely present at intermediate frequency prior to the origin  
505 of the chromosomal races within the past 6,000 to 10,000 years (Nachman *et al.* 1994) and were dispersed  
506 through Europe as mice colonized the continent from the south and east (Boursot *et al.* 1993). The  
507 presence of *R2d2<sup>HC</sup>* in a non-*M. m. domesticus* sample (SPRET/Eij, *M. spretus* from Cadiz, Spain) is best  
508 explained by recent introgression following secondary contact with *M. m. domesticus* (Bonhomme *et al.*  
509 2007; Yang *et al.* 2011).

#### 510 **A new member of the *Cwc22* family.**

511 The duplication that gave rise to *R2d2* also created a new copy of *Cwc22*. Based on our assembly of the  
512 *R2d2* sequence, the open reading frame of *Cwc22<sup>R2d2</sup>* is intact and encodes a nearly full-length predicted  
513 protein that retains the two key functional domains characteristic of the *Cwc22* family. Inspection of  
514 RNA-seq data from samples with high copy number at *R2d2* reveals several novel transcript isoforms  
515 whose expression appears to be copy-number- and tissue-dependent. In testis, the most abundant  
516 isoform retains an intron containing an ERV insertion (red arrow in **Figure 4**), consistent with the well-  
517 known transcriptional promiscuity in this tissue. The most abundant isoforms in adult brain is unusual in  
518 that its stop codon is in an internal exon which is followed by a 7 kbp 3' UTR in the terminal exon.  
519 Transcripts with a stop codon in an internal exon are generally subject to nonsense-mediated decay  
520 (NMD) triggered by the presence of exon-junction complexes downstream the stop codon. Curiously,  
521 *Cwc22* is itself a member of the exon-junction complex (Steckelberg *et al.* 2012).

522 That an essential gene involved in such a central biochemical pathway should both escape NMD and be  
523 overexpressed more than tenfold is surprising. Preliminary data from the Diversity Outbred population  
524 shows that the *R2d2<sup>HC</sup>* allele is associated with elevated levels of both *Cwc22* transcripts and protein in  
525 adult liver (Gary Churchill, personal communication). Further studies will be required to determine the  
526 distribution of transcription and translation of *Cwc22* across isoforms, tissues and developmental stages.

#### 527 **Conclusions and future directions**

528 Our detailed analysis of the evolutionary trajectory of *R2d2* provides insight into the fate of duplicated  
529 sequences over short (within-species) timescales. The exceptionally high mutation rate and low  
530 recombination associated specifically with hemizygous *R2d2<sup>HC</sup>* alleles motivate hypotheses regarding the  
531 biochemical mechanisms which contribute to observed patterns of polymorphism at this and similar loci.  
532 Finally, the birth of a new member of the deeply conserved *Cwc22* gene family in *R2d2* provides an  
533 opportunity to test predictions regarding the evolution of young duplicate gene pairs.

## 534 **METHODS**

### 535 **Mice**

536 Wild *M. musculus* mice used in this study were trapped at a large number of sites across Europe, the  
537 United States, the Middle East, northern India and Taiwan (**Figure 8A**). Trapping was carried out in  
538 accordance with local regulations and with the approval of all relevant regulatory bodies for each locality  
539 and institution. Trapping locations are listed in **Supplementary Table 1**. Most samples have been  
540 previously published (Didion *et al.* 2016).

541 Tissue samples from the progenitors of the wild-derived inbred strains ZALENDE/EiJ (*M. m. domesticus*),  
542 TIRANO/EiJ (*M. m. domesticus*) and SPRET/EiJ (*M. spretus*) were provided by Muriel Davisson, as  
543 described in Didion *et al.* (2016).

544 Tissue samples from the high running (HR) selection and intercross lines were obtained as described in  
545 Didion *et al.* (2016).

546 Female Diversity Outbred mice used for estimating mutation rates at *R2d2* were obtained from the  
547 Jackson Laboratory and housed with a single FVB/NJ male. Progeny were sacrificed at birth by cervical  
548 dislocation in order to obtain tissue for genotyping.

549 All live laboratory mice were handled in accordance with the IACUC protocols of the University of North  
550 Carolina at Chapel Hill.

### 551 **DNA preparation**

552 *High molecular weight DNA.* High molecular weight DNA was obtained for samples genotyped with the  
553 Mouse Diversity Array or subject to whole-genome sequencing. Genomic DNA was extracted from tail,  
554 liver or spleen using a standard phenol-chloroform procedure (Sambrook and Russell 2006). High  
555 molecular weight DNA for most inbred strains was obtained from the Jackson Laboratory, and the  
556 remainder as a generous gift from Francois Bonhomme and the University of Montpellier Wild Mouse  
557 Genetic Repository.

558 *Low molecular weight DNA.* Low molecular weight DNA was obtained for samples to be genotyped on the  
559 MegaMUGA array (see “Microarray genotyping” below). Genomic DNA was isolated from tail, liver,  
560 muscle or spleen using Qiagen Genra Puregene or DNeasy Blood & Tissue kits according to the  
561 manufacturer’s instructions.

### 562 **Whole-genome sequencing and variant discovery**

563 *Inbred strains.* Sequencing data for inbred strains of mice except ZALENDE/EiJ and LEWES/EiJ was  
564 obtained from the Sanger Mouse Genomes Project website ([ftp://ftp-mouse.sanger.ac.uk/current\\_bams](ftp://ftp-mouse.sanger.ac.uk/current_bams)) as  
565 aligned BAM files. Details of the sequencing pipeline are given in Keane *et al.* (2011). Coverage ranged  
566 from approximately 25X to 50X per sample.

567 The strains LEWES/EiJ and ZALENDE/EiJ were sequenced at the University of North Carolina High-  
568 Throughput Sequencing Facility. Libraries were prepared from high molecular weight DNA using the  
569 Illumina TruSeq kit and insert size approximately 250 bp, and 2x100bp paired-end reads were generated  
570 on an Illumina HiSeq 2000 instrument. LEWES/EiJ was sequenced to approximately 12X coverage and  
571 ZALENDE/EiJ to approximately 20X. Alignment was performed as in Keane *et al.* (2011).

572 *Wild mice.* Whole-genome sequencing data from 26 wild *M. m. domesticus* individuals described in Pezer  
573 *et al.* (2015) was downloaded from ENA under accession #PRJEB9450. Coverage ranged from  
574 approximately 12X to 20X per sample. An additional two wild *M. m. domesticus* individuals, IT175 and  
575 ES446, were sequenced at the University of North Carolina to approximate coverage 8X each. Raw reads  
576 from an additional 10 wild *M. m. castaneus* described in Halligan *et al.* (2013), sequenced to approximately  
577 20X each, were downloaded from ENA under accession #PRJEB2176. Reads for a single *Mus caroli*  
578 individual sequenced to approximately 40X were obtained from ENA under accession #PRJEB2188.  
579 Reads for each sample were realigned to the mm10 reference using bwa-mem v0.7.12 with default  
580 parameters (Li 2013). Optical duplicates were removed with samblaster (Faust and Hall 2014).

581 *Variant discovery.* Polymorphic sites on chromosome 2 in the vicinity of *R2d2* (**Figure 8B**) were called  
582 using freebayes v0.9.21-19-gc003c1e (Garrison and Marth 2012) with parameters “—standard-filters”  
583 using the Sanger Mouse Genomes Project VCF files as a list of known sites (parameter “—@”). Raw calls  
584 were filtered to have quality score > 30, root mean square mapping quality > 20 (for both reference and  
585 alternate allele calls) and at most 2 alternate alleles.

#### 586 **Copy-number estimation**

587 *R2d* copy number was estimated using qPCR as described in Didion *et al.* (2016). Briefly, we used  
588 commercial TaqMan assays against intron-exon boundaries in *Cwc22* (Life Technologies assay numbers  
589 Mm00644079\_cn and Mm00053048\_cn) to determine copy number relative to reference genes *Tert* (cat. no.  
590 4458368, for target Mm00644079\_cn) or *Tfrc* (cat. no. 4458366, for target Mm00053048\_cn). Cycle  
591 thresholds for *Cwc22* relative to the reference gene were normalized across assay batches using linear  
592 mixed models with batch and target-reference pair treated as random effects. Control samples with  
593 known haploid *R2d* copy numbers of 1 (C57BL/6J), 2 (CAST/EiJ), 17 (WSB/EiJ×C57BL/6J)F<sub>1</sub> and 34  
594 (WSB/EiJ) were included in each batch.

595 Samples were classified as having 1, 2 or >2 haploid copies of *R2d* using linear discriminant analysis. The  
596 classifier was trained on the normalized cycle thresholds of the control samples from each plate, whose  
597 precise integer copy number is known, and applied to the remaining samples.

#### 598 **Microarray genotyping**

599 Genome-wide genotyping was performed using MegaMUGA, the second version of the Mouse Universal  
600 Genotyping Array platform (Neogen/GeneSeek, Lincoln, NE) (Morgan *et al.* 2016). Genotypes were called  
601 using the GenCall algorithm implemented in the Illumina BeadStudio software (Illumina Inc, Carlsbad,  
602 CA). For quality control we computed, for each marker *i* on the array:  $S_i = X_i + Y_i$ , where  $X_i$  and  $Y_i$  are  
603 the normalized hybridization intensities for the two alleles. The expected distribution of  $S_i$  was computed  
604 from a large set of reference samples. We excluded arrays for which the distribution of  $S_i$  was  
605 substantially shifted from this reference; in practice, failed arrays can be trivially identified in this manner  
606 (Morgan *et al.* 2016). Access to MegaMUGA genotypes was provided by partnership between the  
607 McMillan and Pardo-Manuel de Villena labs and the UNC Systems Genetics Core Facility.

608 Additional genotypes for inbred strains and wild mice from the Mouse Diversity Array were obtained  
609 from Yang *et al.* (2011).

#### 610 **De novo assembly of *R2d2***

611 Raw whole-genome sequencing reads for WSB/EiJ from the Sanger Mouse Genomes Project were  
612 converted to a multi-string Burrows-Wheeler transform and associated FM-index (msBWT) (Holt and  
613 McMillan 2014) using the msbwt v0.1.4 Python package (<https://pypi.python.org/pypi/msbwt>). The  
614 msBWT and FM-index implicitly represent a suffix array of sequencing to provide efficient queries over  
615 arbitrarily large string sets. Given a seed *k*-mer present in that string set, this property can be exploited to  
616 rapidly construct a de Bruijn graph which can in turn be used for local *de novo* assembly of a target  
617 sequence (**Supplementary Figure 8A**). The edges in that graph can be assigned a weight (corresponding  
618 to the number of reads containing the *k* + 1-mer implied by the edge) which can be used to evaluate  
619 candidate paths when the graph branches (**Supplementary Figure 8B**).

620 *R2d2* was seeded with the 30 bp sequence (TCTAGAGCATGAGCCTCATTTATCATGCCT) at the  
621 proximal boundary of *R2d1* in the GRCm38/mm10 reference genome. A single linear contig was

622 assembled by “walking” through the local de Bruijn graph. Because WSB/Eij has ~33 copies of *R2d2* and a  
623 single copy of *R2d1*, any branch point in the graph which represents a paralogous variant should have  
624 outgoing edges with weights differing by a factor of approximately 33. Furthermore, when two (or more)  
625 branch points occur within less than the length of a read, it should be possible to “phase” the underlying  
626 variants by following single reads through both branch points (**Supplementary Figure 8B**). We used  
627 these heuristics to assemble the sequence of *R2d2* (corresponding to the higher-weight path through the  
628 graph) specifically.

629 After assembling a chunk of approximately 500 bp the contig was checked for colinearity with the  
630 reference sequence (*R2d1*) using BLAT and CLUSTAL-W2 (using the EMBL-EBI web server:  
631 <http://www.ebi.ac.uk/Tools/msa/clustalw2/>).

632 Repetitive elements such as retroviruses are refractory to assembly with our method. Upon traversing  
633 into a repetitive element, the total edge weight (total number of reads) and number of branch points  
634 (representing possible linear assembled sequences) in the graph become large. It was sometimes possible  
635 to assemble a fragment of a repetitive element at its junction with unique sequence but not to assemble  
636 unambiguously across the repeat. Regions of unassemblable sequence were marked with blocks of Ns,  
637 and assembly re-seeded using a nearby *k*-mer from the reference sequence. The final contig is provided  
638 in FASTA format in **Supplementary File 1**.

639 The final contig was checked against its source msBWT by confirming that each 30-mer in the contig  
640 which did not contain an N was supported by at least 60 reads. A total of 16 additional haplotypes in 8  
641 regions of *R2d* totaling 16.9 kbp (**Supplementary Table 6**) were assembled in a similar fashion, using the  
642 WSB *R2d2* contig and the *R2d1* reference sequence as guides. Multiple sequence alignments from these  
643 regions are provided in **Supplementary File 1**.

#### 644 **Sequence analysis of *R2d2* contig**

645 *Pairwise alignment of *R2d* paralogs.* The reference *R2d1* sequence and our *R2d2* contig were aligned using  
646 LASTZ v1.03.54 (<http://www.bx.psu.edu/~rsharris/lastz/>) with parameters “—step=10 —seed=match12 —  
647 notransition —exact=20 —notrim —identity=95”.

648 *Transposable element (TE) content.* The *R2d2* contig was screened for TE insertions using the RepeatMasker  
649 web server (<http://www.repeatmasker.org/cgi-bin/WEBRepeatMasker>) with species set to “mouse” and  
650 default settings otherwise. As noted previously, we could not assemble full-length repeats, but the  
651 fragments we could assemble at junctions with unique sequence allowed identification of some candidate  
652 TEs to the family level. *R2d1*-specific TEs were defined as TEs annotated in the RepeatMasker track at the  
653 UCSC Genome Browser with no evidence (no homologous sequence, and no Ns) at the corresponding  
654 position in the *R2d2* contig. Candidate *R2d2*-specific TEs were defined as gaps  $\geq 100$  bp in size in the  
655 alignment to *R2d1* for which the corresponding *R2d2* sequence was flagged by RepeatMasker.

656 *Gene conversion tracts.* To unambiguously define gene conversion events without confounding from  
657 paralogous sequence, we examined 15 wild *M. m. domesticus* samples and 37 laboratory strains with  
658 evidence of 2 diploid copies of *R2d*. We confirmed that these copies of *R2d* were located at *R2d1* by  
659 finding read pairs spanning the junction between *R2d1* and neighboring sequence. Gene conversion tracts  
660 were delineated as clusters of derived alleles shared with *R2d2*. Using a pairwise alignment of *R2d2*  
661 and *R2d1* we identified single-nucleotide variants between the two sequences, and queried those sites in  
662 aligned reads for *Mus caroli*. If the *Mus caroli* and *R2d1* shared an allele, we recorded the site as a derived  
663 allele informative for the presence of *R2d2*. We used the resulting list of 1,411 informative sites to query  
664 aligned reads for the samples of interest and recorded, for each site and each sample, whether the derived



665 allele (*R2d2*), ancestral allele (*R2d1*) or both alleles were present. Conversion tracts were then identified by  
666 manual inspection. Boundaries of conversion tracts were defined at approximately the midpoint between  
667 the first *R2d1*- (or *R2d2*-) specific variant and the last *R2d2*- (or *R2d1*-) specific variant.

668 *Sequence diversity in R2d1 and R2d2.* Assembling individual copies of *R2d2* is infeasible in high-copy  
669 samples. Instead we treated each *R2d* unit as an independent sequence and used the number of  
670 segregating sites to estimate sequence diversity. Segregating sites were defined as positions in a collection  
671 of alignments (BAM files) with evidence of an alternate allele. To identify segregating sites we used  
672 freebayes v0.9.21-19-gc003c1e (Garrison and Marth 2012) with parameters “-ui -Kp 20 —use-best-n-alleles  
673 2 -m 8”. These parameters treat each sample as having ploidy up to 20, impose an uninformative prior on  
674 genotype frequencies, and limit the algorithm to the discovery of atomic variants (SNVs or short indels,  
675 not multinucleotide polymorphisms or other complex events) with at most 2 alleles at each segregating  
676 site. Sites in low-complexity sequence (defined as Shannon entropy < 1.6 in the 30 bp window centered on  
677 the site) or within 10 bp of another variant site were further masked, to minimize spurious calls due to  
678 ambiguous alignment of indels. To avoid confounding with the retrocopies of *Cwc22* outside *R2d*, coding  
679 exons of *Cwc22* were masked. Finally, sites corresponding to an unaligned or gap position in the pairwise  
680 alignment between *R2d1* and *R2d2* were masked.

681 To compute diversity in *R2d1* we counted segregating sites in 12 wild *M. m. domesticus* samples with 2  
682 diploid copies of *R2d* (total of 24 sequences), confirmed to be in *R2d1* by the presence of read pairs  
683 spanning the junction between *R2d1* and neighboring sequence. To compute diversity in *R2d2*, we  
684 counted segregating sites in 14 wild *M. m. domesticus* samples with >2 diploid copies of *R2d* (range 3 – 83  
685 per sample; total of 406 sequences) but excluded sites corresponding to variants among *R2d1* sequences.  
686 Remaining sites were phased to *R2d2* by checking for the presence of a 31-mer containing the site and the  
687 nearest *R2d1*-vs-*R2d2* difference in the raw reads for each sample using the corresponding msBWT.  
688 Sequence diversity was then computed using Watterson’s estimator (Watterson 1975), dividing by the  
689 number of alignable bases (128973) to yield a per-site estimate. Standard errors were estimated by 100  
690 rounds of resampling over the columns in the *R2d1*-vs-*R2d2* alignment.

## 691 **Analyses of *Cwc22* expression**

692 *RNA-seq read alignment.* Expression of *Cwc22* was examined in adult whole brain using data from  
693 Crowley *et al.* (2015), SRA accession #SRP056236. Paired-end reads (2x100bp) were obtained from 8  
694 replicates each of 3 inbred strains: CAST/EiJ, PWK/PhJ and WSB/EiJ. Raw reads were aligned to the  
695 mm10 reference using STAR v2.4.2a (Dobin *et al.* 2012) with default parameters for paired-end reads.  
696 Alignments were merged into a single file per strain for further analysis. Expression in adult testis was  
697 examined in 23 wild-derived inbred strains from Phifer-Rixey *et al.* (2014) SRA accession #PRJNA252743.  
698 Single-end reads (76bp) were aligned to the mm10 genome with STAR using default parameters for  
699 single-end, non-strand-specific reads.

700 *Transcript assembly.* Read alignments were manually inspected to assess support for *Cwc22* isoforms in  
701 Ensembl v83 annotation. To identify novel isoforms in *R2d2*, we applied the Trinity v0.2.6 pipeline  
702 (Grabherr *et al.* 2011) to the subset of reads from WSB/EiJ which could be aligned to *R2d1* plus their mates  
703 (a set which represents a mixture of *Cwc22<sup>R2d1</sup>* and *Cwc22<sup>R2d2</sup>* reads). De novo transcripts were aligned  
704 both to the mm10 reference and to the *R2d2* contig using BLAT, and were assigned to *R2d1* or *R2d2* based  
705 on sequence similarity. Because expression from *R2d2* is high in WSB/EiJ, *R2d2*-derived transcripts  
706 dominated the assembled set. Both manual inspection and the Trinity assembly indicated the presence of  
707 retained introns and an extra 3’ exon, as described in the **Results**. To obtain a full set of *Cwc22* transcripts  
708 including those of both *R2d1* and *R2d2* origin, we supplemented the *Cwc22* transcripts in Ensembl v83

709 with their paralogs from *R2d2* as determined by a strict BLAT search against the *R2d2* contig. We  
710 manually created additional transcripts reflecting intron-retention and 3' extension events described  
711 above, and obtained their sequence from the *R2d2* contig.

712 *Abundance estimation.* Relative abundance of *Cwc22* paralogs was estimated using kallisto v0.42.3 (Bray *et*  
713 *al.* 2015) with parameters “—bias” (to estimate and correct library-specific sequence-composition biases).  
714 The transcript index used for pseudoalignment and quantification included only the *Cwc22* targets.

## 715 **Phylogenetic analyses**

716 *Tree for R2d.* Multiple sequence alignments for 8 the regions in **Supplementary Figure 5** were generated  
717 using MUSCLE (Edgar 2004) with default parameters. The resulting alignments were manually trimmed  
718 and consecutive gaps removed. Phylogenetic trees were inferred with RAxML v8.1.9 (Stamatakis 2014)  
719 using the GTR+gamma model with 4 rate categories and *M. caroli* as an outgroup. Uncertainty of tree  
720 topologies was evaluated using 100 bootstrap replicates.

721 *Divergence time.* The time of the split between *R2d1* and *R2d2* was estimated using the Bayesian method  
722 implemented in BEAST v1.8.1r6542 (Drummond *et al.* 2012). We assumed a divergence time for *M. caroli*  
723 of 5 Mya and a strict molecular clock, and analyzed the concatenated alignment for our *de novo* assembled  
724 regions under the GTR+gamma model with 4 rate categories and allowance for a proportion of invariant  
725 sites. The chain was run for 10 million iterations with trees sampled every 1000 iterations.

726 *Local phylogeny around R2d2.* Genotypes for 173 SNPs in the region surrounding *R2d2* (chr2: 83 — 84 Mb)  
727 were obtained for 90 individuals representing both laboratory and wild mice genotyped with the Mouse  
728 Diversity Array (Yang *et al.* 2011) (**Supplementary Table 1**). Individuals with evidence of heterozygosity  
729 (>3 heterozygous calls) were excluded to avoid ambiguity in phylogenetic inference. A distance matrix  
730 for the remaining 62 samples was created by computing the proportion of alleles shared identical by state  
731 between each pair of samples. A neighbor-joining tree (**Figure 3C**) was inferred from the distance matrix  
732 and rooted at the most recent common ancestor of the *M. musculus*- and non-*M. musculus* samples.

733 *Cwc22 coding sequences.* To create the tree of *Cwc22* coding sequences, we first obtained the sequences of  
734 all its paralogs in mouse. The coding sequence of *Cwc22<sup>R2d1</sup>* (RefSeq transcript NM\_030560.5) was  
735 obtained from the UCSC Genome Browser and aligned to our *R2d2* contig with BLAT to extract the exons  
736 of *Cwc22<sup>R2d2</sup>*. The coding sequence of retro-*Cwc22* (genomic sequence corresponding to GenBank cDNA  
737 AK145290) was obtained from the UCSC Genome Browser. Coding and protein sequences of *Cwc22*  
738 homologs from non-*M. musculus* species were obtained from Ensembl (Cunningham *et al.* 2014). The  
739 sequences were aligned with MUSCLE and manually trimmed, and a phylogenetic tree estimated as  
740 described above.

741 We observed that the branches in the rodent clade of the *Cwc22* tree appeared to be longer than branches  
742 for other taxa. We used PAML (Yang *et al.* 2007) to test the hypothesis that *Cwc22* is under relaxed  
743 purifying selection in rodents using the branch-site model (null model “model = 2, NSsites = 2, fix\_omega  
744 = 1”; alternative model “model = 2, NSsites = 2, omega = 1, fix\_omega = 1”) as described in the PAML  
745 manual. This is a test of difference in evolutionary rate on a “foreground” branch ( $\omega_1$ ) — in our case, the  
746 rodent clade — relative to the tree-wide “background” rate ( $\omega_0$ ). The distribution of the test statistic is an  
747 even mixture of a  $\chi^2$  distribution with 1 df and a point mass at zero; to obtain the *p*-value, we calculated  
748 the quantile of the  $\chi^2$  distribution with 1 df and divided by 2.

## 749 Genome-wide sequence divergence

750 The msBWT of a collection of whole-genome sequencing reads can be used to estimate the divergence  
751 between the corresponding template sequence (*i.e.* genome) and a reference sequence as follows. Non-  
752 overlapping  $k$ -mers from the reference sequence are queried against the msBWT. (The value of  $k$  is  
753 chosen such that nearly all  $k$ -mers drawn from genomic sequence exclusive of repetitive elements.) Let  $x$   
754 be the count of reads containing an exact match to the  $k$ -mer or its reverse complement. If the template  
755 and reference sequence are identical, standard theory for shotgun sequencing (Lander and Waterman  
756 1988) holds that  $P(x > 0|\lambda) = 1 - e^{-\lambda}$ , where  $\lambda$  is the average sequencing coverage. We assume  
757  $P(x > 0|\lambda) \approx 1$ , which is satisfied in practice for high-coverage sequencing.

758 However, if a haploid template sequence contains at least one variant (versus the reference) within the  
759 queried  $k$ -mer, it will be the case that  $x = 0$ . We use this fact and assume that mutations arise along a  
760 sequence via a Poisson process to estimate the rate parameter  $\alpha$  from the proportion of  $k$ -mers that have  
761 read count zero. Let  $m$  be the number of mutations arising between a target and reference in a window of  
762 length  $L$ , and  $y$  the number of  $k$ -mers in that window with nonzero read count. Then  $P(m = 0|\alpha) = e^{-\alpha L}$   
763 and a simple estimator for  $\alpha$  is  $\hat{\alpha} = -\log\left(\frac{1-y}{L}\right)$ .

764 Interpretation of  $\alpha$  is straightforward in the haploid case: it is the per-base rate of sequence divergence  
765 between the template sequence and the reference sequence. In the diploid case it represents a lower  
766 bound on the sequence divergence of the two homologous chromosomes.

767 We applied this estimator with  $k = 31$  and  $L = 1000 \times k = 31$  kbp to msBWTs for 7 inbred strains (3 *M. m.*  
768 *domesticus*, 1 *M. m. musculus*, 1 *M. m. castaneus*, 1 *M. spretus*, 1 *M. caroli*) and 2 wild *M. m. domesticus*  
769 individuals (IT175, ES446) using the GRCm38/mm10 mouse reference sequence as the source of  $k$ -mer  
770 queries. As shown in **Figure 9A**, the mode of the distribution of divergence values matches what is  
771 expected based on the ancestry of the samples with respect to the reference. To identify divergent regions,  
772 we fit a discrete-time hidden Markov model (HMM) to the windows divergence values. The HMM had  
773 two hidden states: “normal” sequence, with emission distribution  $N(0.005, 0.005)$  and initial probability  
774 0.99; and “divergent” sequence, with emission distribution  $N(0.02, 0.005)$  and initial probability 0.01. The  
775 transmission probability between states was  $1 \times 10^{-5}$ . Posterior decodings were obtained via the Viterbi  
776 algorithm, as implemented in the R package HiddenMarkov ([https://cran.r-](https://cran.r-project.org/package=HiddenMarkov)  
777 [project.org/package=HiddenMarkov](https://cran.r-project.org/package=HiddenMarkov)).

778 Significance tests for overlap with genomic features were performed using the resampling algorithm  
779 implemented in the Genomic Association Tester (GAT) package for Python  
780 (<https://pypi.python.org/pypi/gat>). Segmental duplications were obtained from the genomicSuperDups  
781 table of the UCSC Genome Browser and genes from Ensembl v83 annotation.

## 782 Analyses of recombination rate at *R2d2*

783 To test the effect of *R2d2* copy number on local recombination rate examined recombination events  
784 accumulated during the first 16 generations of breeding of the Diversity Outbred (DO) population, in  
785 which the high-copy *R2d2* allele from WSB/Eij is segregating. Founder haplotype reconstructions were  
786 obtained for 4,640 DO individuals reported in (Didion *et al.* 2016), and recombination events were  
787 identified as junctions between founder haplotypes. We compared the frequency of junctions involving a  
788 WSB/Eij haplotype to junctions not involving a WSB/Eij haplotype over the region chr2: 75-90 Mb.  
789 Within each generation we also tested for differences in the lengths of haplotype blocks overlapping *R2d2*  
790 using one-sided Wilcoxon rank-sum tests (alternative hypothesis: WSB/Eij haplotypes longer than

791 others). Resulting  $p$ -values were subject to Bonferroni correction: for nominal significance level  $\alpha = 0.01$ ,  
792 the corrected threshold is  $p = \frac{0.01}{12} = 8.3 \times 10^{-4}$ .

793 We also estimated the difference between observed and expected recombination fraction in 11  
794 experimental crosses in which one of the parental lines was segregating for a high-copy allele at *R2d2*. We  
795 obtained expected recombination fractions from the standard mouse genetic map (Cox *et al.* 2009), which  
796 was constructed from crosses between strains lacking *R2d2<sup>HC</sup>* alleles. Genotype data was obtained from  
797 The Jackson Laboratory's Mouse Phenome Database QTL Archive  
798 (<http://phenome.jax.org/db/q?rtn=qtl/home>). Recombination fractions were calculated using R/qtl  
799 (<http://rqt.org/>). Confidence intervals for difference between observed and expected recombination  
800 fractions were calculated by 100 iterations of nonparametric bootstrapping over individuals in each  
801 dataset.

802 Results of these analyses are presented in **Supplementary Figure 7**.

## 803 DATA AVAILABILITY

804 All *de novo* assemblies used in this study are included in **Supplementary File 1**. The data structures on  
805 which the assemblies are based, and the interactive computational tools used for assembly, are publicly  
806 available at <http://www.csbio.unc.edu/CEGSseq/index.py?run=MsbwtTools>. Genotype data used for  
807 mapping the location of *R2d2* and defining associated haplotypes (**Figure 3B-C**) are available on Dryad  
808 (accession XXX).

## 809 ACKNOWLEDGMENTS

810 We thank all the scientists and personnel who collected and processed the wild mouse samples used in  
811 this study. In particular we thank Francois Bonhomme for providing samples from wild-derived inbred  
812 strains housed at the University of Montpellier Wild Mouse Genetic Repository, and Ted Garland for  
813 providing tissue samples from the HR selection lines and related crosses. This work was supported by  
814 National Institutes of Health grants P50GM076468 (FPMdV), U19AI100625 (FPMdV, APM),  
815 F30MH103925 (APM), T32GM067553 (JPD, APM), and by Vaadia-BARD Postdoctoral Fellowship Award  
816 FI-12 478-13 to LY. Additional support was provided by Cancer Research UK, the European Research  
817 Council, EMBO Young Investigator Programme (DTO), European Molecular Biology Laboratory (DTO,  
818 PF) and the Wellcome Trust (WT095908) (PF) and (WT098051) (PF, DTO); and finally by the European  
819 Union's Seventh Framework Programme (FP7/2007-2013) under grant agreement HEALTH-F4-2010-  
820 241504 (EURATRANS).

## 821 REFERENCES

- 822 Baarends W. M., Wassenaar E., Laan R. van der, Hoogerbrugge J., Sleddens-Linkels E., Hoeijmakers J. H.  
823 J., Boer P. de, Grootegoed J. A., 2005 Silencing of unpaired chromatin and histone H2A ubiquitination in  
824 mammalian meiosis. *Mol Cell Biol* **25**: 1041–1053.
- 825 Bailey J. A., Eichler E. E., 2006 Primate segmental duplications: Crucibles of evolution, diversity and  
826 disease. *Nat Rev Genet* **7**: 552–564.

- 827 Bonhomme F., Rivals E., Orth A., Grant G. R., Jeffreys A. J., Bois P. R., 2007 Species-wide distribution of  
828 highly polymorphic minisatellite markers suggests past and present genetic exchanges among house  
829 mouse subspecies. *Genome Biol* **8**: R80.
- 830 Boursot P., Auffray J. C., Britton-Davidian J., Bonhomme F., 1993 The evolution of house mice. *Annu Rev*  
831 *Ecol Syst* **24**: 119–152.
- 832 Boyden L. M., Lewis J. M., Barbee S. D., Bas A., Girardi M., Hayday A. C., Tigelaar R. E., Lifton R. P., 2008  
833 Skint1, the prototype of a newly identified immunoglobulin superfamily gene cluster, positively selects  
834 epidermal T cells. *Nat Genet* **40**: 656–662.
- 835 Bray N., Pimentel H., Melsted P., Patcher L., 2015 Near-optimal RNA-seq quantification. arXiv:  
836 1505.02710 [q-bio.QM]
- 837 Chen J.-M., Cooper D. N., Chuzhanova N., Férec C., Patrinos G. P., 2007 Gene conversion: Mechanisms,  
838 evolution and human disease. *Nat Rev Genet* **8**: 762–775.
- 839 Chevret P., Veyrunes F., Britton-Davidian J., 2005 Molecular phylogeny of the genus *Mus* (Rodentia:  
840 Murinae) based on mitochondrial and nuclear data. *Biol J Linn Soc* **84**: 417–427.
- 841 Cole F., Keeney S., Jasin M., 2010 Comprehensive, fine-scale dissection of homologous recombination  
842 outcomes at a hot spot in mouse meiosis. *Mol Cell* **39**: 700–710.
- 843 Cole F., Baudat F., Grey C., Keeney S., Massy B. de, Jasin M., 2014 Mouse tetrad analysis provides insights  
844 into recombination mechanisms and hotspot evolutionary dynamics. *Nat Genet* **46**: 1072–1080.
- 845 Collaborative Cross Consortium, 2012 The genome architecture of the Collaborative Cross mouse genetic  
846 reference population. *Genetics* **190**: 389–401.
- 847 Cox A., Ackert-Bicknell C. L., Dumont B. L., Ding Y., Bell J. T., Brockmann G. A., Wergedal J. E., Bult C.,  
848 Paigen B., Flint J., Tsaih S. W., Churchill G. A., Broman K. W., 2009 A new standard genetic map for  
849 the laboratory mouse. *Genetics* **182**: 1335–1344.
- 850 Crowley J. J., Zhabotynsky V., Sun W., Huang S., Pakatci I. K., Kim Y., Wang J. R., Morgan A. P., Calaway  
851 J. D., Aylor D. L., Yun Z., Bell T. A., Buus R. J., Calaway M. E., Didion J. P., Gooch T. J., Hansen S. D.,  
852 Robinson N. N., Shaw G. D., Spence J. S., Quackenbush C. R., Barrick C. J., Nonneman R. J., Kim K.,  
853 Xenakis J., Xie Y., Valdar W., Lenarcic A. B., Wang W., Welsh C. E., Fu C.-P., Zhang Z., Holt J., Guo Z.,  
854 Threadgill D. W., Tarantino L. M., Miller D. R., Zou F., McMillan L., Sullivan P. F., Pardo-Manuel de  
855 Villena F., 2015 Analyses of allele-specific gene expression in highly divergent mouse crosses identifies  
856 pervasive allelic imbalance. *Nat Genet* **47**: 353–360.
- 857 Cunningham F., Amode M. R., Barrell D., Beal K., Billis K., Brent S., Carvalho-Silva D., Clapham P.,  
858 Coates G., Fitzgerald S., Gil L., Giron C. G., Gordon L., Hourlier T., Hunt S. E., Janacek S. H., Johnson N.,  
859 Juettemann T., Kahari A. K., Keenan S., Martin F. J., Maurel T., McLaren W., Murphy D. N., Nag R.,  
860 Overduin B., Parker A., Patricio M., Perry E., Pignatelli M., Riat H. S., Sheppard D., Taylor K., Thormann  
861 A., Vullo A., Wilder S. P., Zadissa A., Aken B. L., Birney E., Harrow J., Kinsella R., Muffato M., Ruffier M.,  
862 Searle S. M. J., Spudich G., Trevanion S. J., Yates A., Zerbino D. R., Flicek P., 2014 Ensembl 2015. *Nucleic*  
863 *Acids Res* **43**: D662–D669.
- 864 Dallas J. F., 1992 Estimation of microsatellite mutation rates in recombinant inbred strains of mouse.  
865 *Mamm Genome* **3**: 452–456.

- 866 Demuth J. P., De Bie T., Stajich J. E., Christianini N., Hahn M. W., 2006 The evolution of mammalian gene  
867 families. *PLoS ONE* **1**: e85.
- 868 Didion J. P., Morgan A. P., Clayshulte A. M.-F., McMullan R. C., Yadgary L., Petkov P. M., Bell T. A., Gatti  
869 D. M., Crowley J. J., Hua K., Aylor D. L., Bai L., Calaway M., Chesler E. J., French J. E., Geiger T. R.,  
870 Gooch T. J., Garland T., Harrill A. H., Hunter K., McMillan L., Holt M., Miller D. R., O'Brien D. A., Paigen  
871 K., Pan W., Rowe L. B., Shaw G. D., Simecek P., Sullivan P. F., Svenson K. L., Weinstock G. M., Threadgill  
872 D. W., Pomp D., Churchill G. A., Pardo-Manuel de Villena F., 2015 A multi-megabase copy number gain  
873 causes maternal transmission ratio distortion on mouse chromosome 2. *PLoS Genet* **11**: e1004850.
- 874 Didion J. P., Morgan A. P., Yadgary L., Bell T. A., McMullan R. C., Solorzano L. O. de, Britton-Davidian J.,  
875 Bult C. J., Campbell K. J., Castiglia R., Ching Y.-H., Chunco A. J., Crowley J. J., Chesler E. J., Förster D. W.,  
876 French J. E., Gabriel S. I., Gatti D. M., Garland T., Giagia-Athanasopoulou E. B., Giménez M. D., Grize S.  
877 A., Gündüz I., Holmes A., Hauffe H. C., Herman J. S., Holt J. M., Hua K., Jolley W. J., Lindholm A. K.,  
878 López-Fuster M. J., Mitsainas G., Luz Mathias M. da, McMillan L., Ramalhinho M. G., Reherrmann B.,  
879 Rosshart S. P., Searle J. B., Shiao M.-S., Solano E., Svenson K. L., Thomas-Laemont P., Threadgill D. W.,  
880 Ventura J., Weinstock G. M., Pomp D., Churchill G. A., Pardo-Manuel de Villena F., 2016 R2d2 drives  
881 selfish sweeps in the house mouse. *Mol Biol Evol*: msw036.
- 882 Dobin A., Davis C. A., Schlesinger F., Drenkow J., Zaleski C., Jha S., Batut P., Chaisson M., Gingeras T. R.,  
883 2012 STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21.
- 884 Dopman E. B., Hartl D. L., 2007 A portrait of copy-number polymorphism in *Drosophila melanogaster*.  
885 *Proceedings of the National Academy of Sciences* **104**: 19920–19925.
- 886 Dover G., 1982 Molecular drive: A cohesive mode of species evolution. *Nature* **299**: 111–117.
- 887 Drummond A. J., Suchard M. A., Xie D., Rambaut A., 2012 Bayesian phylogenetics with BEAUti and the  
888 BEAST 1.7. *Mol Biol Evol* **29**: 1969–1973.
- 889 Dumont B. L., Eichler E. E., 2013 Signals of historical interlocus gene conversion in human segmental  
890 duplications. *PLoS ONE* **8**: e75949.
- 891 Edgar R. C., 2004 MUSCLE: Multiple sequence alignment with high accuracy and high throughput.  
892 *Nucleic Acids Res* **32**: 1792–1797.
- 893 Egan C. M., Sridhar S., Wigler M., Hall I. M., 2007 Recurrent DNA copy number variation in the  
894 laboratory mouse. *Nat Genet* **39**: 1384–1389.
- 895 Eickbush T. H., Eickbush D. G., 2007 Finely orchestrated movements: Evolution of the ribosomal RNA  
896 genes. *Genetics* **175**: 477–485.
- 897 Faust G. G., Hall I. M., 2014 SAMBLASTER: Fast duplicate marking and structural variant read  
898 extraction. *Bioinformatics* **30**: 2503–2505.
- 899 Garrison E., Marth G., 2012 Haplotype-based variant detection from short-read sequencing. *arXiv*.
- 900 Goodman M., Czelusniak J., Moore G. W., Romero-Herrera A. E., Matsuda G., 1979 Fitting the gene  
901 lineage into its species lineage, a parsimony strategy illustrated by cladograms constructed from globin  
902 sequences. *Syst Biol* **28**: 132–163.

- 903 Grabherr M. G., Haas B. J., Yassour M., Levin J. Z., Thompson D. A., Amit I., Adiconis X., Fan L.,  
904 Raychowdhury R., Zeng Q., Chen Z., Mauceli E., Hacohen N., Gnirke A., Rhind N., Palma F. di, Birren B.  
905 W., Nusbaum C., Lindblad-Toh K., Friedman N., Regev A., 2011 Full-length transcriptome assembly from  
906 RNA-seq data without a reference genome. *Nat Biotechnol* **29**: 644–652.
- 907 Hallast P., Balaesque P., Bowden G. R., Ballereau S., Jobling M. A., 2013 Recombination dynamics of a  
908 human Y-chromosomal palindrome: Rapid GC-biased gene conversion, multi-kilobase conversion tracts,  
909 and rare inversions. *PLoS Genet* **9**: e1003666.
- 910 Halligan D. L., Kousathanas A., Ness R. W., Harr B., Eöry L., Keane T. M., Adams D. J., Keightley P. D.,  
911 2013 Contributions of protein-coding and regulatory change to adaptive molecular evolution in murid  
912 rodents. *PLoS Genet* **9**: e1003995.
- 913 Handsaker R. E., Doren V. V., Berman J. R., Genovese G., Kashin S., Boettger L. M., McCarroll S. A., 2015  
914 Large multiallelic copy number variations in humans. *Nat Genet* **47**: 296–303.
- 915 Hauffe H. C., Searle J. B., 1993 Extreme karyotypic variation in a *Mus musculus domesticus* hybrid zone:  
916 The tobacco mouse story revisited. *Evolution* **47**: 1374.
- 917 Holt J., McMillan L., 2014 Merging of multi-string BWTs with applications. *Bioinformatics* **30**: 3524–3531.
- 918 Hurles M., 2002 Are 100,000 “SNPs” useless? *Science* **298**: 1509.
- 919 Hurst J. L., Payne C. E., Nevison C. M., Marie A. D., Humphries R. E., Robertson D. H. L., Cavaggioni A.,  
920 Beynon R. J., 2001 Individual recognition in mice mediated by major urinary proteins. *Nature* **414**: 631–  
921 634.
- 922 Keane T. M., Goodstadt L., Danecek P., White M. A., Wong K., Yalcin B., Heger A., Agam A., Slater G.,  
923 Goodson M., Furlotte N. A., Eskin E., Nellåker C., Whitley H., Cleak J., Janowitz D., Hernandez-Pliego P.,  
924 Edwards A., Belgard T. G., Oliver P. L., McIntyre R. E., Bhomra A., Nicod J., Gan X., Yuan W., Weyden L.  
925 van der, Steward C. A., Bala S., Stalker J., Mott R., Durbin R., Jackson I. J., Czechanski A., Guerra-  
926 Assunção J. A., Donahue L. R., Reinholdt L. G., Payseur B. A., Ponting C. P., Birney E., Flint J., Adams D.  
927 J., 2011 Mouse genomic variation and its effect on phenotypes and gene regulation. *Nature* **477**: 289–294.
- 928 Koonin, E. V., 2005 Orthologs, paralogs and evolutionary genomics. *Annu Rev Genet* **39**: 309–338.
- 929 Laan R. van der, 2004 Ubiquitin ligase Rad18Sc localizes to the XY body and to other chromosomal  
930 regions that are unpaired and transcriptionally silenced during male meiotic prophase. *J Cell Sci* **117**:  
931 5023–5033.
- 932 Lander E. S., Waterman M. S., 1988 Genomic mapping by fingerprinting random clones: A mathematical  
933 analysis. *Genomics* **2**: 231–239.
- 934 Li H., 2006 TreeFam: A curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res*  
935 **34**: D572–D580.
- 936 Li H., 2013 Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*.
- 937 Liao D., 1997 Concerted evolution of the tandemly repeated genes encoding human U2 snRNA (the  
938 RUN2 locus) involves rapid intrachromosomal homogenization and rare interchromosomal gene  
939 conversion. *EMBO J* **16**: 588–598.

- 940 Lindholm A. K., Dyer K. A., Firman R. C., Fishman L., Forstmeier W., Holman L., Johannesson H., Knief  
941 U., Kokko H., Larracuente A. M., Manser A., Montchamp-Moreau C., Petrosyan V. G., Pomiankowski A.,  
942 Presgraves D. C., Safronova L. D., Sutter A., Unckless R. L., Verspoor R. L., Wedell N., Wilkinson G. S.,  
943 Price T. A., 2016 The ecology and evolutionary dynamics of meiotic drive. *Trends Ecol Evol*.
- 944 Liu E. Y., Morgan A. P., Chesler E. J., Wang W., Churchill G. A., Villena F. P.-M. de, 2014 High-resolution  
945 sex-specific linkage maps of the mouse reveal polarized distribution of crossovers in male germline.  
946 *Genetics* **197**: 91–106.
- 947 Lynch M., 2000 The evolutionary fate and consequences of duplicate genes. *Science* **290**: 1151–1155.
- 948 Mahadevaiah S. K., Turner J. M., Baudat F., Rogakou E. P., Boer P. de, Blanco-Rodríguez J., Jasin M.,  
949 Keeney S., Bonner W. M., Burgoyne P. S., 2001 *Nat Genet* **27**: 271–276.
- 950 Melamed C., Kupiec M., 1992 Effect of donor copy number on the rate of gene conversion in the yeast  
951 *Saccharomyces cerevisiae*. *Mol Genet Genom* **235**: 97–103.
- 952 Moran P., 1975 Wandering distributions and the electrophoretic profile. *Theor Popul Biol* **8**: 318–330.
- 953 Morgan A. P., Fu C.-P., Kao C.-Y., Welsh C. E., Didion J. P., Yadgary L., Hyacinth L., Ferris M. T., Bell T.  
954 A., Miller D. R., Giusti-Rodríguez P., Nonneman R. J., Cook K. D., Whitmire J. K., Gralinski L. E., Keller  
955 M., Attie A. D., Churchill G. A., Petkov P., Sullivan P. F., Brennan J. R., McMillan L., Pardo-Manuel de  
956 Villena F., 2016 The Mouse Universal Genotyping Array: From substrains to subspecies. *G3* **6**.
- 957 Muffato M., Louis A., Poisnel C. E., Crollius H. R., 2010 Genomicus: A database and a browser to study  
958 gene synteny in modern and ancestral genomes. *Bioinformatics* **26**: 1119–1121.
- 959 Nachman M. W., Boyer S. N., Searle J. B., Aquadro C. F., 1994 Mitochondrial DNA variation and the  
960 evolution of robertsonian chromosomal races of house mice, *Mus domesticus*. *Genetics* **136**: 1105–1120.
- 961 Nagylaki T., Petes T. D., 1982 Intrachromosomal gene conversion and the maintenance of sequence  
962 homogeneity among repeated genes. *Genetics* **100**: 315–337.
- 963 Pamilo P., Nei M., 1988 Relationships between gene trees and species trees. *Mol Biol Evol* **5**: 568–583.
- 964 Pei B., Sisu C., Frankish A., Howald C., Habegger L., Mu X., Harte R., Balasubramanian S., Tanzer A.,  
965 Diekhans M., Reymond A., Hubbard T. J., Harrow J., Gerstein M. B., 2012 The GENCODE pseudogene  
966 resource. *Genome Biol* **13**: R51.
- 967 Pezer, Harr B., Teschke M., Babiker H., Tautz D., 2015 Divergence patterns of genic copy number  
968 variation in natural populations of the house mouse (*Mus musculus domesticus*) reveal three conserved  
969 genes with major population-specific expansions. *Genome Res* **25**: 1114–1124.
- 970 Phifer-Rixey M., Bomhoff M., Nachman M. W., 2014 Genome-wide patterns of differentiation among  
971 house mouse subspecies. *Genetics* **198**: 283–297.
- 972 Rozen S., Skaletsky H., Marszalek J. D., Minx P. J., Cordum H. S., Waterston R. H., Wilson R. K., Page D.  
973 C., 2003 Abundant gene conversion between arms of palindromes in human and ape Y chromosomes.  
974 *Nature* **423**: 873–876.

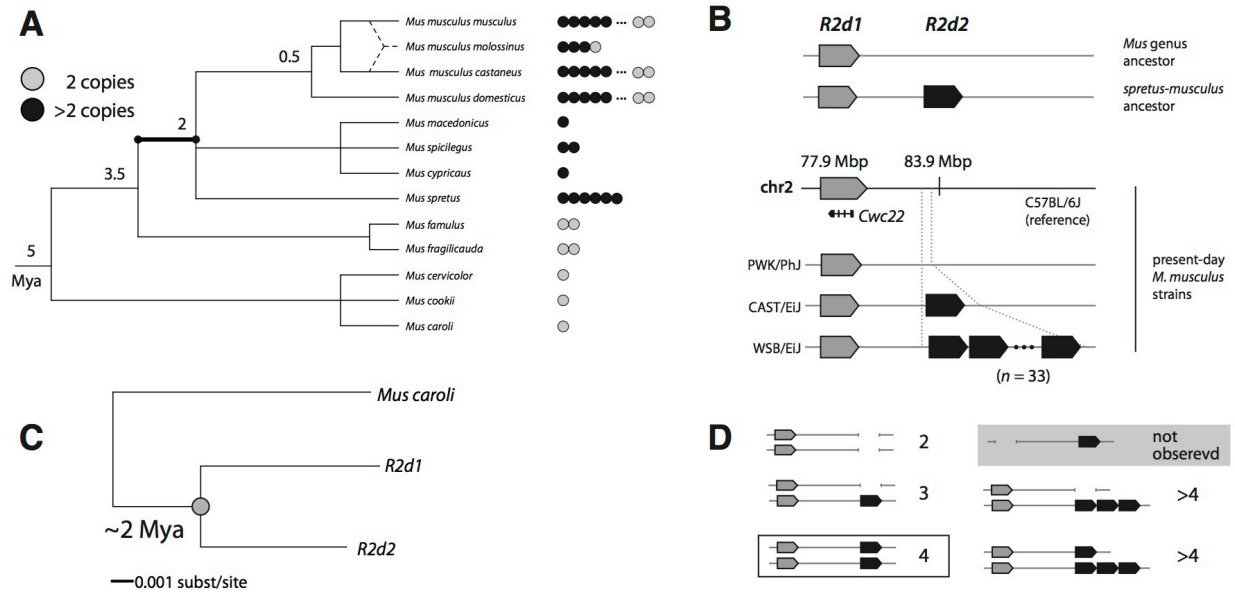


- 975 Salcedo T., Geraldles A., Nachman M. W., 2007 Nucleotide variation in wild and inbred mice. *Genetics*  
976 **177**: 2277–2291.
- 977 Sambrook J., Russell D. W. (Eds.), 2006 *Molecular cloning: A laboratory manual*. Cold Spring Harbor  
978 Laboratory Press.
- 979 Scavetta R. J., Tautz D., 2010 Copy number changes of CNV regions in intersubspecific crosses of the  
980 house mouse. *Mol Biol Evol* **27**: 1845–1856.
- 981 Schindelhauer D., 2002 Evidence for a fast, intrachromosomal conversion mechanism from mapping of  
982 nucleotide variants within a homogeneous alpha-satellite DNA array. *Genome Res* **12**: 1815–1826.
- 983 She X., Cheng Z., Zöllner S., Church D. M., Eichler E. E., 2008 Mouse segmental duplication and copy  
984 number variation. *Nat Genet* **40**: 909–914.
- 985 Shi J., Wolf S. E., Burke J. M., Presting G. G., Ross-Ibarra J., Dawe R. K., 2010 Widespread gene conversion  
986 in centromere cores. *PLoS Biol* **8**: e1000327.
- 987 Soh Y., Alföldi J., Pyntikova T., Brown L., Graves T., Minx P., Fulton R., Kremitzki C., Koutseva N.,  
988 Mueller J., Rozen S., Hughes J., Owens E., Womack J., Murphy W., Cao Q., de~Jong P., Warren W.,  
989 Wilson R., Skaletsky H., Page D., 2014 Sequencing the mouse Y chromosome reveals convergent gene  
990 acquisition and amplification on both sex chromosomes. *Cell* **159**: 800–813.
- 991 Stamatakis A., 2014 RAxML version 8: A tool for phylogenetic analysis and post-analysis of large  
992 phylogenies. *Bioinformatics* **30**: 1312–1313.
- 993 Stankiewicz P., Lupski J. R. 2002 Genome architecture, rearrangements and genomic disorders. *Trends*  
994 *Genet* **18**: 74–82.
- 995 Steckelberg A.-L., Boehm V., Gromadzka A., Gehring N., 2012 Cwc22 connects pre-mRNA splicing and  
996 exon junction complex assembly. *Cell Rep* **2**: 454–461.
- 997 Suzuki H., Shimada T., Terashima M., Tsuchiya K., Aplin K., 2004 Temporal, spatial, and ecological  
998 modes of evolution of Eurasian *Mus* based on mitochondrial and nuclear gene sequences. *Mol*  
999 *Phylogenet Evol* **33**: 626–646.
- 1000 Svenson K. L., Gatti D. M., Valdar W., Welsh C. E., Cheng R., Chesler E. J., Palmer A. A., McMillan L.,  
1001 Churchill G. A., 2012 High-resolution genetic mapping using the mouse Diversity Outbred population.  
1002 *Genetics* **190**: 437–447.
- 1003 Swallow J. G., Carter P. A., Jr. T. G., 1998 Artificial selection for increased wheel-running behavior in  
1004 house mice. *Behav Genet* **28**: 227–237.
- 1005 Treangen T. J., Salzberg S. L., 2011 Repetitive DNA and next-generation sequencing: Computational  
1006 challenges and solutions. *Nat Rev Genet*.
- 1007 Turner J. M., Aprelikova O., Xu X., Wang R., Kim S., Chandramouli G. V., Barrett J., Burgoyne P. S., Deng  
1008 C.-X., 2004 BRCA1, histone H2AX phosphorylation, and male meiotic sex chromosome inactivation. *Curr*  
1009 *Biol* **14**: 2135–2142.

- 1010 Turner D. J., Miretti M., Rajan D., Fiegler H., Carter N. P., Blayney M. L., Beck S., Hurles M. E., 2007  
1011 Germline rates of de novo meiotic deletions and duplications causing several genomic disorders. *Nat*  
1012 *Genet* **40**: 90–95.
- 1013 Waterston R. H., Chinwalla A. T., Cook L. L., Delehaunty K. D., Fewell G. A., Fulton L. A., Fulton R. S.,  
1014 Graves T. A., Hillier L. W., Mardis E. R., McPherson J. D., Miner T. L., Nash W. E., Nelson J. O., Nhan M.  
1015 N., Pepin K. H., Pohl C. S., Ponce T. C., Schultz B., Thompson J., Trevaskis E., Waterston R. H., Wendl M.  
1016 C., Wilson R. K., Yang S.-P., An P., Berry E., Birren B., Bloom T., Brown D. G., Butler J., Daly M., David R.,  
1017 Deri J., Dodge S., Foley K., Gage D., Gnerre S., Holzer T., Jaffe D. B., Kamal M., Karlsson E. K., Kells C.,  
1018 Kirby A., Kulbokas E. J., Lander E. S., Landers T., Leger J. P., Levine R., Lindblad-Toh K., Mauceli E.,  
1019 Mayer J. H., McCarthy M., Meldrim J., Meldrim J., Mesirov J. P., Nicol R., Nusbaum C., Seaman S., Sharpe  
1020 T., Sheridan A., Singer J. B., Santos R., Spencer B., Stange-Thomann N., Vinson J. P., Wade C. M.,  
1021 Wierzbowski J., Wyman D., Zody M. C., Birney E., Goldman N., Kasprzyk A., Mongin E., Rust A. G.,  
1022 Slater G., Stabenau A., Ureta-Vidal A., Whelan S., Ainscough R., Attwood J., Bailey J., Barlow K., Beck S.,  
1023 Burton J., Clamp M., Clee C., Coulson A., Cuff J., Curwen V., Cutts T., Davies J., Eyras E., Grafham D.,  
1024 Gregory S., Hubbard T., Hunt A., Jones M., Joy A., Leonard S., Lloyd C., Matthews L., McLaren S., McLay  
1025 K., Meredith B., Mullikin J. C., Ning Z., Oliver K., Overton-Larty E., Plumb R., Potter S., Quail M., Rogers  
1026 J., Scott C., Searle S., Shownkeen R., Sims S., Wall M., West A. P., Willey D., Williams S., Abril J. F., Guigó  
1027 R., Parra G., Agarwal P., Agarwala R., Church D. M., Hlavina W., Maglott D. R., Sapojnikov V.,  
1028 Alexandersson M., Pachter L., Antonarakis S. E., Dermitzakis E. T., Reymond A., Ucla C., Baertsch R.,  
1029 Diekhans M., Furey T. S., Hinrichs A., Hsu F., Karolchik D., Kent W. J., Roskin K. M., Schwartz M. S.,  
1030 Sugnet C., Weber R. J., Bork P., Letunic I., Suyama M., Torrents D., Zdobnov E. M., Botcherby M., Brown  
1031 S. D., Campbell R. D., Jackson I., Bray N., Couronne O., Dubchak I., Poliakov A., Rubin E. M., Brent M. R.,  
1032 Flicek P., Keibler E., Korf I., Batalov S., Bult C., Frankel W. N., Carninci P., Hayashizaki Y., Kawai J.,  
1033 Okazaki Y., Cawley S., Kulp D., Wheeler R., Chiaromonte F., Collins F. S., Felsenfeld A., Guyer M.,  
1034 Peterson J., Wetterstrand K., Copley R. R., Mott R., Dewey C., Dickens N. J., Emes R. D., Goodstadt L.,  
1035 Ponting C. P., Winter E., Dunn D. M., Niederhausern A. C. von, Weiss R. B., Eddy S. R., Johnson L. S.,  
1036 Jones T. A., Elnitski L., Kolbe D. L., Eswara P., Miller W., O'Connor M. J., Schwartz S., Gibbs R. A.,  
1037 Muzny D. M., Glusman G., Smit A., Green E. D., Hardison R. C., Yang S., Haussler D., Hua A., Roe B. A.,  
1038 Kucherlapati R. S., Montgomery K. T., Li J., Li M., Lucas S., Ma B., McCombie W. R., Morgan M., Pevzner  
1039 P., Tesler G., Schultz J., Smith D. R., Tromp J., Worley K. C., Lander E. S., Abril J. F., Agarwal P.,  
1040 Alexandersson M., Antonarakis S. E., Baertsch R., Berry E., Birney E., Bork P., Bray N., Brent M. R., Brown  
1041 D. G., Butler J., Bult C., Chiaromonte F., Chinwalla A. T., Church D. M., Clamp M., Collins F. S., Copley R.  
1042 R., Couronne O., Cawley S., Cuff J., Curwen V., Cutts T., Daly M., Dermitzakis E. T., Dewey C., Dickens  
1043 N. J., Diekhans M., Dubchak I., Eddy S. R., Elnitski L., Emes R. D., Eswara P., Eyras E., Felsenfeld A.,  
1044 Flicek P., Frankel W. N., Fulton L. A., Furey T. S., Gnerre S., Glusman G., Goldman N., Goodstadt L.,  
1045 Green E. D., Gregory S., Guigó R., Hardison R. C., Haussler D., Hillier L. W., Hinrichs A., Hlavina W.,  
1046 Hsu F., Hubbard T., Jaffe D. B., Kamal M., Karolchik D., Karlsson E. K., Kasprzyk A., Keibler E., Kent W.  
1047 J., Kirby A., Kolbe D. L., Korf I., Kulbokas E. J., Kulp D., Lander E. S., Letunic I., Li M., Lindblad-Toh K.,  
1048 Ma B., Maglott D. R., Mauceli E., Mesirov J. P., Miller W., Mott R., Mullikin J. C., Ning Z., Pachter L.,  
1049 Parra G., Pevzner P., Poliakov A., Ponting C. P., Potter S., Reymond A., Roskin K. M., Sapojnikov V.,  
1050 Schultz J., Schwartz M. S., Schwartz S., Searle S., Singer J. B., Slater G., Smit A., Stabenau A., Sugnet C.,  
1051 Suyama M., Tesler G., Torrents D., Tromp J., Ucla C., Vinson J. P., Wade C. M., Weber R. J., Wheeler R.,  
1052 Winter E., Yang S.-P., Zdobnov E. M., Whelan S., Worley K. C., Zody M. C., 2002 Initial sequencing and  
1053 comparative analysis of the mouse genome. *Nature* **420**: 520–562.
- 1054 Watterson G., 1975 On the number of segregating sites in genetical models without recombination. *Theor*  
1055 *Popul Biol* **7**: 256–276.

- 1056 White M. A., Ané C., Dewey C. N., Larget B. R., Payseur B. A., 2009 Fine-scale phylogenetic discordance  
1057 across the house mouse genome. *PLoS Genet* **5**: e1000729.
- 1058 Yang H., Bell T. A., Churchill G. A., Pardo-Manuel de Villena F., 2007 On the subspecific origin of the  
1059 laboratory mouse. *Nat Genet* **39**: 1100–1107.
- 1060 Yang H., Wang J. R., Didion J. P., Buus R. J., Bell T. A., Welsh C. E., Bonhomme F., Yu A. H.-T., Nachman  
1061 M. W., Pialek J., Tucker P., Boursot P., McMillan L., Churchill G. A., Pardo-Manuel de Villena F., 2011  
1062 Subspecific origin and haplotype diversity in the laboratory mouse. *Nat Genet* **43**: 648–655.
- 1063 Yeh T.-C., Liu H.-L., Chung C.-S., Wu N.-Y., Liu Y.-C., Cheng S.-C., 2010 Splicing factor Cwc22 is required  
1064 for the function of Prp2 and for the spliceosome to escape from a futile pathway. *Mol Cell Biol* **31**: 43–53.
- 1065

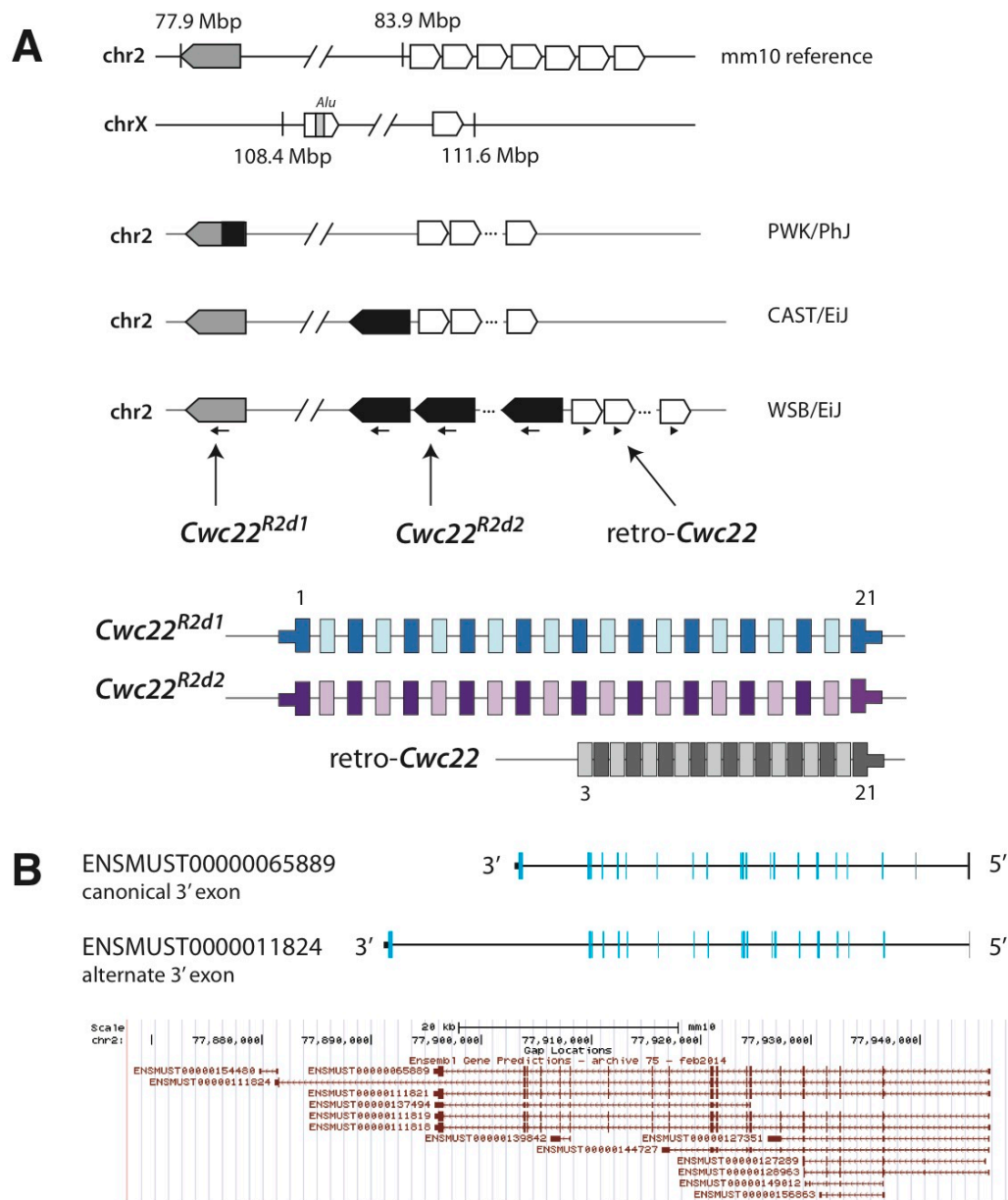
1066 **FIGURE LEGENDS**



1067

1068 **Figure 1. Origin and age of the *R2d2* duplication.** (A) *R2d* copy number across the phylogeny of the  
 1069 genus *Mus*. Each dot represents one individual; grey dots indicate diploid copy number 2 and black dots  
 1070 copy number >2. The duplication event giving rise to *R2d1* and *R2d2* most likely occurred on the  
 1071 highlighted branch. Approximate divergence times (REF: Suzuki 2004) are given in millions of years ago  
 1072 (Mya) at internal nodes. (B) Schematic structure of the *R2d1-R2d2* locus. The mouse reference genome  
 1073 (strain C57BL/6J, *M. m. domesticus*) contains a single copy of *R2d* at *R2d1*. Wild-derived inbred strains vary  
 1074 in haploid copy number from 1 (PWK/PhJ, *M. m. musculus*) to 2 (CAST/EiJ, *M. m. castaneus*) to 33  
 1075 (WSB/EiJ, *M. m. domesticus*). *R2d1* is located at approximately 77.9 Mbp and *R2d2* at 83.8 Mbp. (C)  
 1076 Concatenated tree constructed from *R2d1* (reference genome) and *de novo* assembled *R2d2* and *M. caroli*  
 1077 sequences assuming a strict molecular clock. The duplication node is indicated with a grey dot. (D)  
 1078 Relationship between observed *R2d* copy-number states and inferred structure of the *R2d1-R2d2* locus.  
 1079 The configuration of the *M. spretus* – *M. musculus* common ancestor (4 diploid copies) is boxed in black.  
 1080 We have yet to identify samples with diploid copy number 2 but no *R2d1* (grey shaded box).

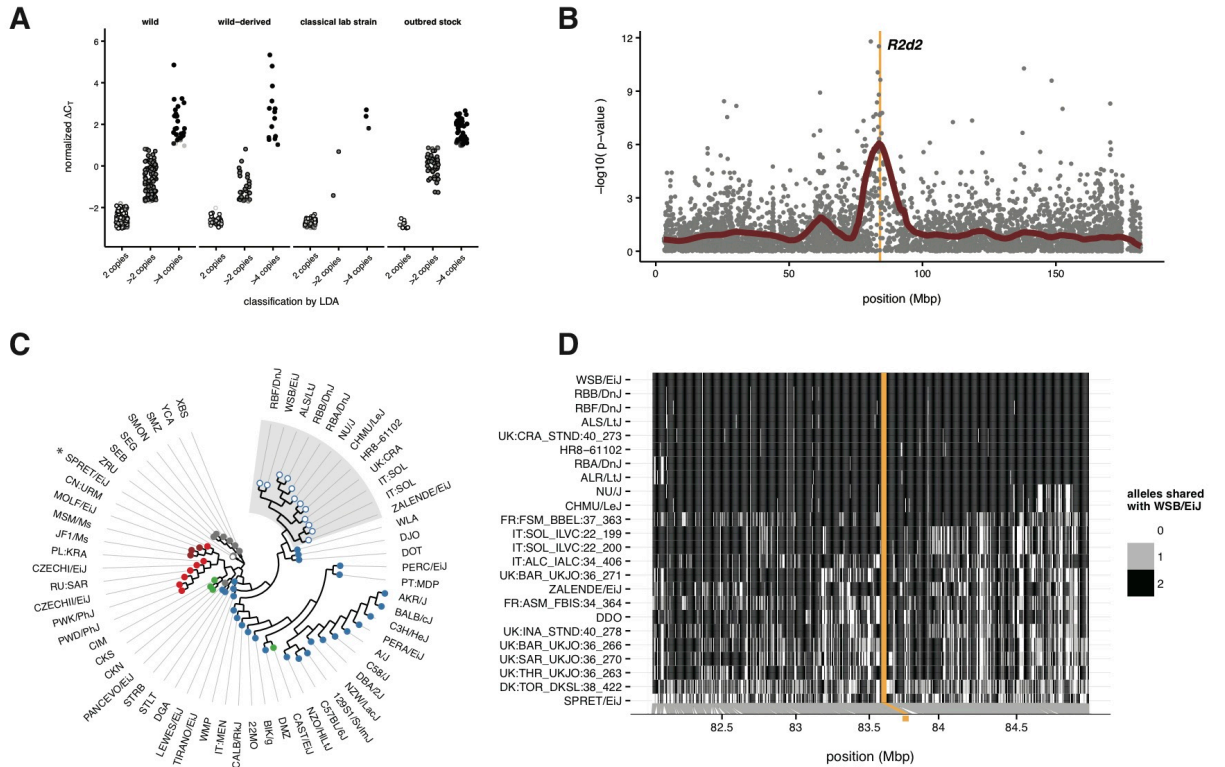
1081



1082

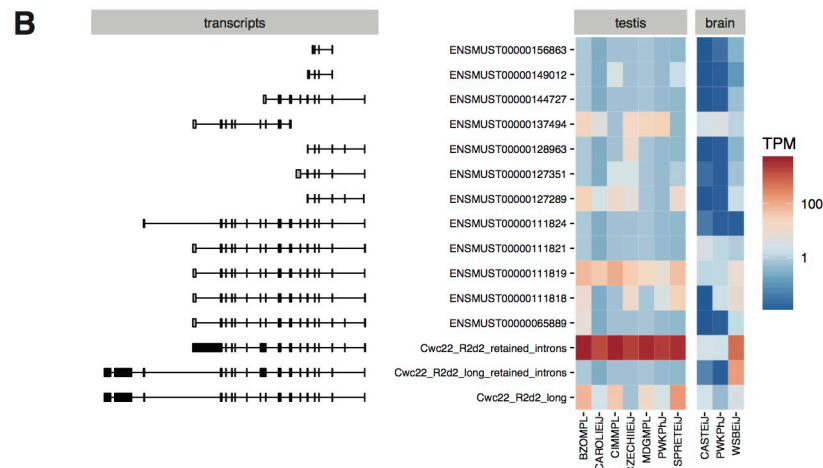
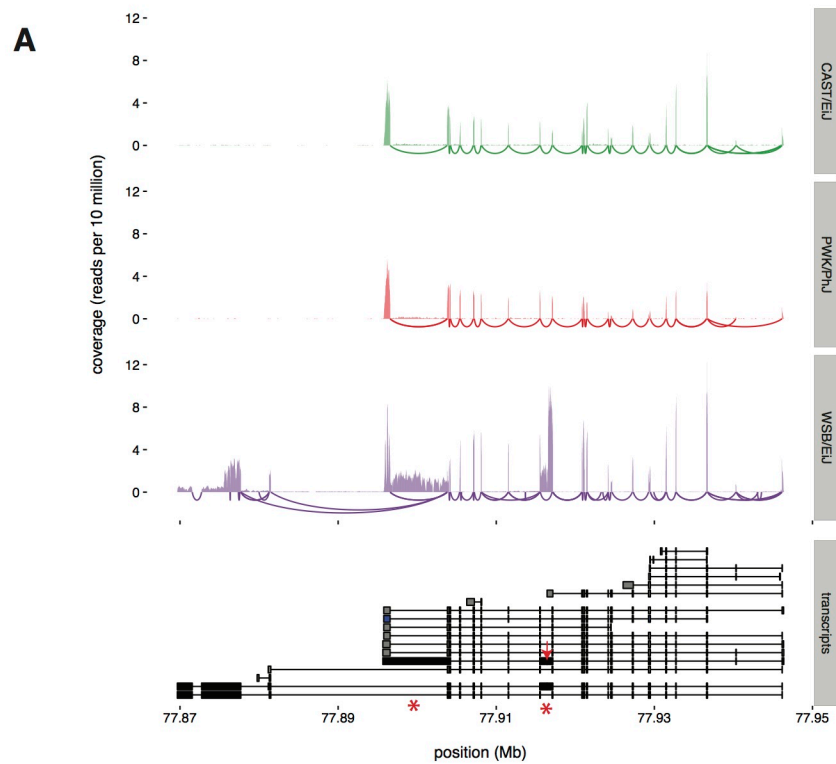
1083 **Figure 2. *Cwc22* paralogs in the mouse genome.** (A) Location and organization of *Cwc22* gene copies  
 1084 present in mouse genomes. The intact coding sequence of *Cwc22* exists in in both *R2d1* (grey shapes) and  
 1085 *R2d2* (black shapes). Retrotransposed copies (empty shapes) exist in two loci on chrX and one locus on  
 1086 chr2, immediately adjacent *R2d2*. Among the retrotransposed copies, coding sequence is intact only in the  
 1087 copy on chr2. Exon numbers are shown in grey above transcript models. (B) Alternate transcript forms of  
 1088 *Cwc22*, using different 3' exons. Coding exons shown in blue and untranslated regions in black. All  
 1089 Ensembl annotated transcripts are shown in the lower panel (from UCSC Genome Browser.)

1090



1091

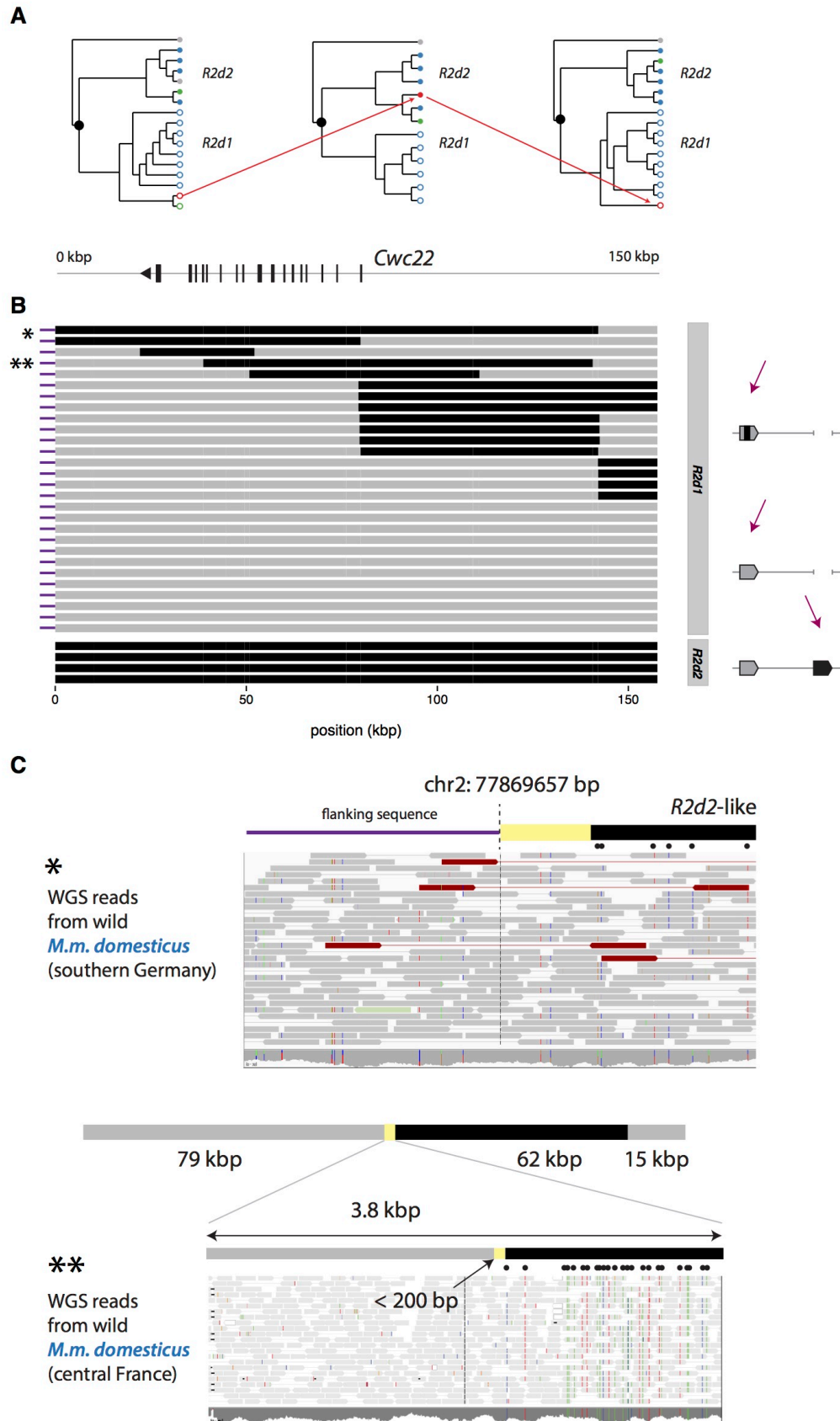
1092 **Figure 3. Copy-number variation of *R2d* in mouse populations worldwide.** (A) Copy-number variation  
 1093 as measured by quantitative PCR. The normalized  $\Delta Ct$  value is proportional to  $\log_2(\text{copy number})$ .  
 1094 Samples are classified as having 2 diploid copies, >2 copies or >4 copies of *R2d* using linear discriminant  
 1095 analysis (LDA). (B) Fine-mapping the location of *R2d2* in 83 samples genotyped on the Mouse Diversity  
 1096 Array (MDA). Grey points give nominal p-values for association between *R2d* copy number and  
 1097 genotype; red points show a smoothed fit through the underlying points. The candidate interval for *R2d2*  
 1098 from Didion *et al.* (2015), shown as an orange shaded box, coincides with the association peak. (C) Local  
 1099 phylogeny at chr2: 83–84 Mbp in 62 wild-caught mice and laboratory strains. Tips are colored by  
 1100 subspecies of origin: *M. m. domesticus*, blue; *M. m. musculus*, red; *M. m. castaneus*, green; *M. m. molossinus*,  
 1101 maroon; outgroup taxa, grey. Individuals with >4 diploid copies of *R2d* are shown as open circles. (D)  
 1102 Haplotypes of laboratory strains and wild mice sharing a high-copy allele at *R2d2*. All samples share a  
 1103 haplotype over the region shaded in orange.



1104

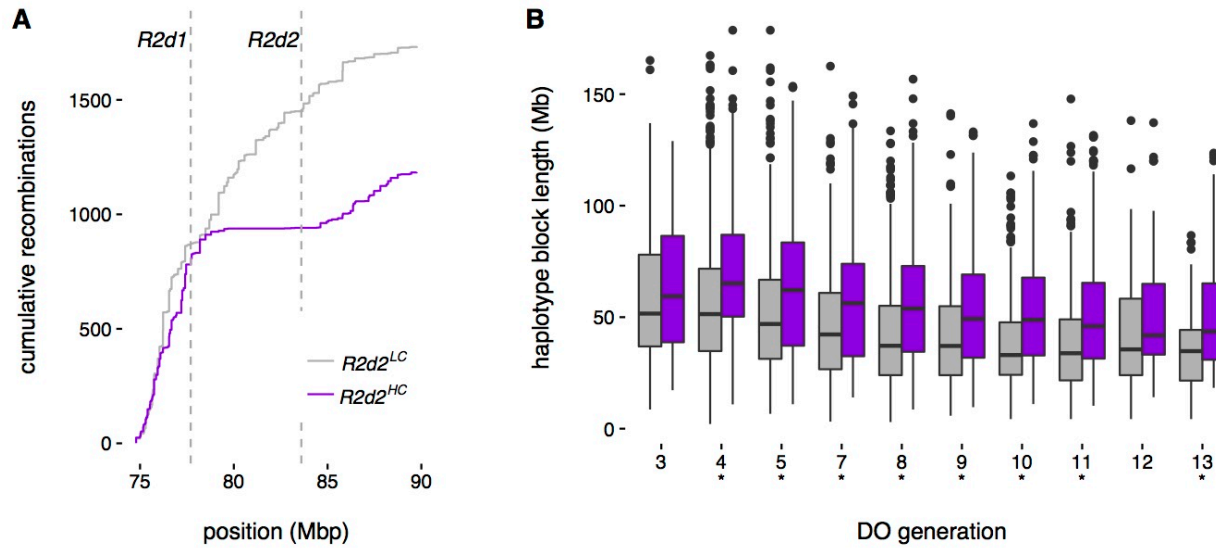
1105 **Figure 4. Expression of *Cwc22* isoforms.** (A) Read coverage and splicing patterns in *Cwc22* in adult  
1106 mouse brain from three wild-derived inbred strains. Swoops below x-axis indicate splicing events  
1107 supported by 5 or more split-read alignments. Known transcripts of *Cwc22*<sup>R2d1</sup> (grey, from Ensembl),  
1108 inferred transcripts from *Cwc22*<sup>R2d2</sup> (black) and the sequence of retro-*Cwc22* mapped back to the parent  
1109 gene (blue) are shown in the lower panel. Red stars indicate retained introns; red arrow indicates  
1110 insertion site of an ERV in *R2d2*. (B) Estimated relative expression of *Cwc22* isoforms (y-axis) in adult  
1111 mouse brain and testis in wild-derived inbred strains (x-axis). TPM, transcripts per million, on log10  
1112 scale.

1113



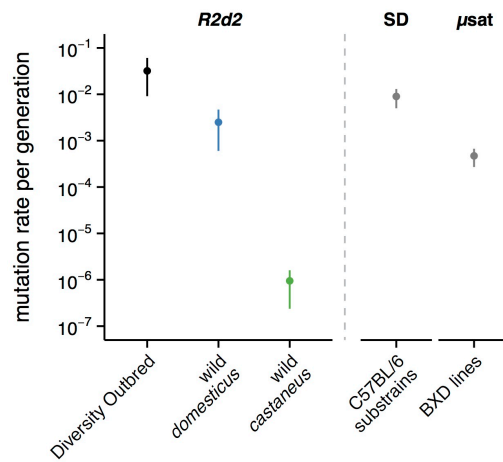


1115 **Figure 5. Signatures of non-allelic gene conversion between *R2d1* and *R2d2*.** (A) Phylogenetic trees for  
1116 three representative intervals across *R2d*. Sequences are labeled according to their subspecies of origin  
1117 using the same color scheme as in **Figure 1**; open circles are *R2d1*-like sequences and closed circles are  
1118 *R2d2*-like. Trees are drawn so that *M. caroli*, the outgroup species used to root the trees, is always  
1119 positioned at the top. The changing affinities of PWK/PhJ (red) and CAST/EiJ (green) along *R2d* are  
1120 evidence of non-allelic gene conversion. (B) *R2d* sequences from 20 wild-caught mice and 5 laboratory  
1121 inbred strains. Each track represents a single chromosome; grey regions are classified as *R2d1*-like based  
1122 on manual inspection of sequence variants, and black-regions *R2d2*-like. Upper panel shows sequences  
1123 from samples with a single copy of *R2d*, residing in *R2d1*. Lower panel shows representative *R2d2*  
1124 sequences for comparison. Asterisks indicate samples for which read alignments are shown in panel C.  
1125 (C) Upper panel, paire-end read alignments (visualized with IGV) across the proximal boundary (dashed  
1126 line) of *R2d1* in a sample with a conversion tract extending to the boundary. Positions of derived variants  
1127 shared with *R2d2* are indicated by black dots. Lower panel, read alignments across the boundary of a  
1128 non-allelic gene conversion tract. *R2d1* sequence from a single chromosome is a mosaic of *R2d1*-like  
1129 (grey) and *R2d2*-like (black) segments. A magnified view of read pairs in the 3.8 kbp surrounding the  
1130 proximal boundary of the tract shows read pairs spanning the junction. Black dots indicate the position of  
1131 derived alleles diagnostic for *R2d2*. The precise breakpoint lies somewhere in the yellow shaded region  
1132 between the last *R2d1*-specific variant and the first *R2d2*-specific variant.



1133

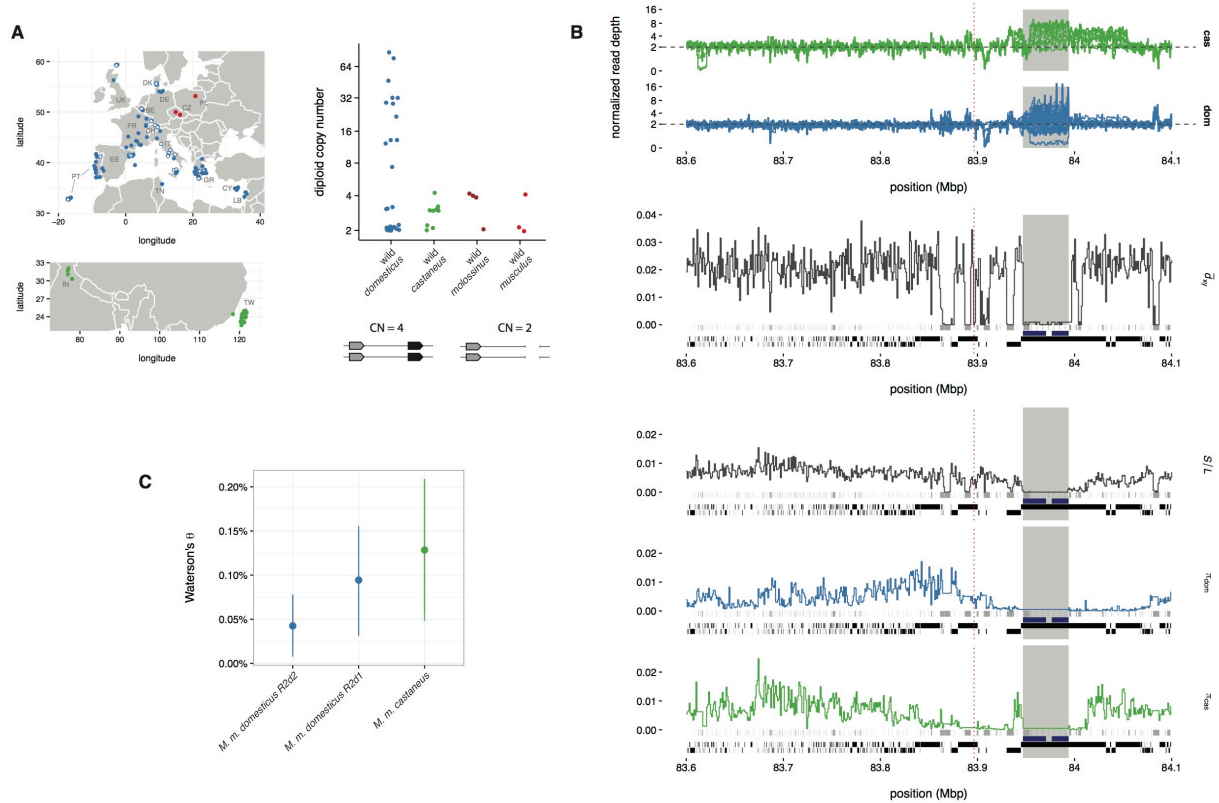
1134 **Figure 6. Suppression of crossing-over around  $R2d2$ .** (A) Cumulative number of unique recombination  
1135 events in the middle region of chr2 in genomes of 4,640 Diversity Outbred mice. Recombination events  
1136 involving the high-copy-number WSB/EiJ haplotype are shown in purple and all other events in grey.  
1137 Dashed vertical lines indicate the position of  $R2d1$  (left) and  $R2d2$  (right). (B) Distribution of haplotype  
1138 block sizes at  $R2d2$  in selected generations of the DO, for  $R2d2^{HC}$  (WSB/EiJ, purple) versus  $R2d2^{LC}$  (the  
1139 other seven founder haplotypes, grey). Asterisks indicate generations in which the length distributions  
1140 are significantly different by Wilcoxon rank-sum test.



1141

1142 **Figure 7. Rate of *de novo* copy-number changes at *R2d2*.** Estimates of per-generation mutation rate for  
1143 CNVs at *R2d2* ( $\pm 1$  bootstrap SE) in the Diversity Outbred population; among wild *M. m. domesticus*; and  
1144 among wild *M. m. castaneus*. For comparison, mutation rates are shown for the CNV with the highest rate  
1145 of recurrence in a C57BL/6J pedigree (Egan *et al.* 2007) and for a microsatellite whose mutation rate was  
1146 estimated in the BXD panel (Dallas 1992).

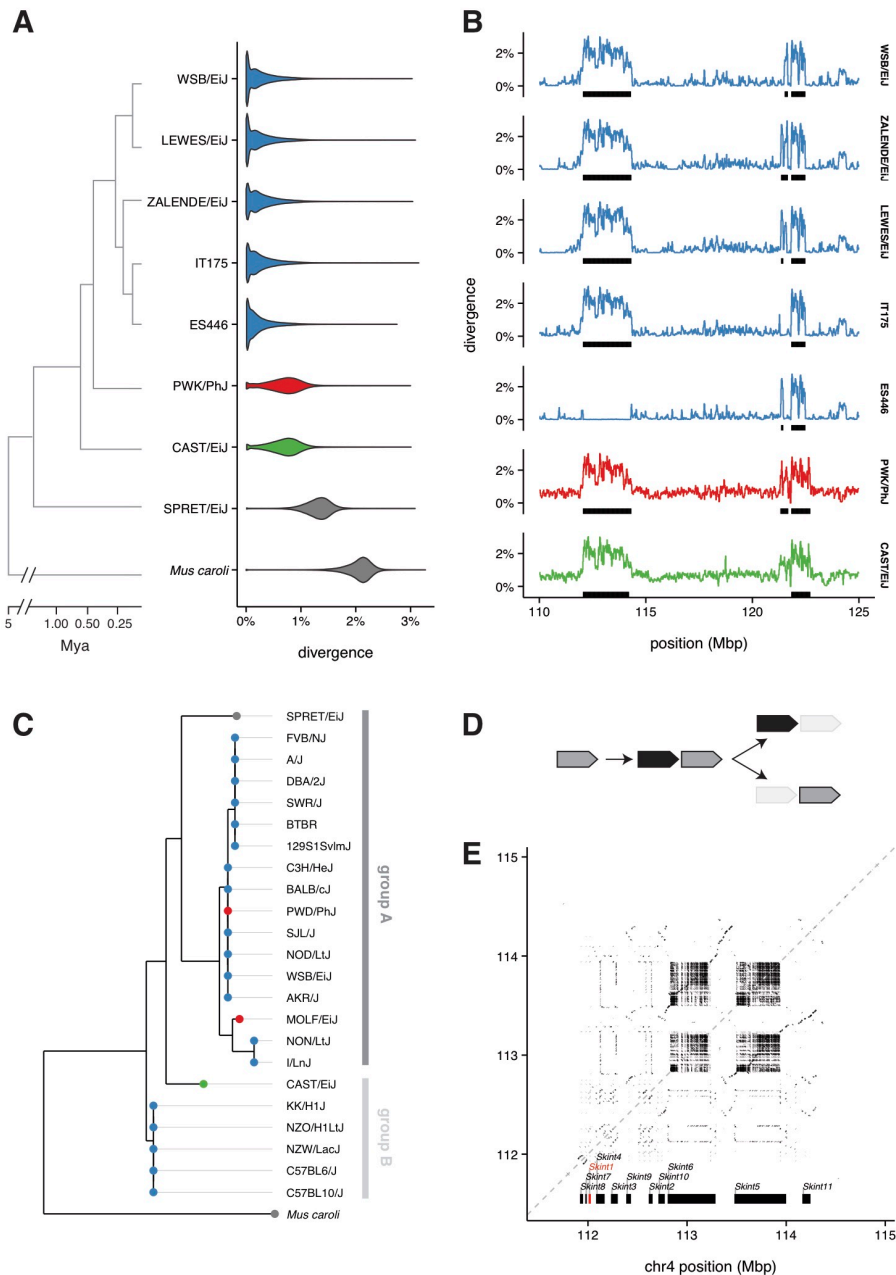
1147



1148

1149 **Figure 8. Sequence and structural diversity around R2d2.** (A) Geographic origin of wild mice used in  
 1150 this study, color-coded by subspecies (blue, *M. m. domesticus*; red, *M. m. musculus*; green, *M. m. castaneus*).  
 1151 Diploid copy number of the *R2d* unit is shown for wild samples for which integer copy-number estimates  
 1152 are available: 26 *M. m. domesticus* and 10 *M. m. castaneus* with whole-genome sequencing data, and  
 1153 representatives from *M. m. molossinus* and *M. m. musculus* for comparison. Schematic shows the  
 1154 *R2d1/R2d2* configurations corresponding to diploid copy numbers of 2 and 4. (B) Profiles of read depth  
 1155 (first two panels), average sequence divergence to outgroup species *M. caroli* ( $d_{xy}$ , third panel), number of  
 1156 segregating sites per base ( $S/L$ , fourth panel) and within-population average heterozygosity ( $\pi$ , fifth and  
 1157 sixth panels). The region shown is 500 kbp in size; the insertion site of *R2d2* is indicated by the red dashed  
 1158 line. Grey boxes along baseline show positions of repetitive elements (from UCSC RepeatMasker track);  
 1159 black boxes show non-recombining haplotype blocks. Blue bars indicate the position of 7 tandem  
 1160 duplications in the mm10 reference sequence with >99% mutual identity, each containing a copy of retro-  
 1161 *Cwc22*. Grey shaded region indicates duplicate sequence absent from *M. caroli*. (C) Estimated per-site  
 1162 nucleotide diversity within *M. m. domesticus* R2d1, *M. m. domesticus* R2d2 and *M. m. castaneus* R2d2.

1163



1164

1165 **Figure 9. A region of excess sequence divergence on chromosome 4.** (A) Genome-wide sequence  
 1166 divergence estimates for representative samples from the sub-genus *Mus*. (B) Estimated sequence  
 1167 divergence in 31 kbp windows across distal chr4 for the samples in panel A. Divergent regions identified  
 1168 by the hidden Markov model (HMM) are indicated with black bars along the horizontal axis. (C)  
 1169 Phylogenetic tree constructed from *Skint1* coding sequences reported in Boyden *et al.* (2008) (D)  
 1170 Schematic representation of the process of gene duplication, followed by differential loss of paralogs  
 1171 along independent lineages. (E) Dotplot of self-alignment of sequence from the region of distal chr4  
 1172 containing the *Skint* gene family. Positions of *Skint* genes are indicated along the horizontal axis; *Skint1*  
 1173 highlighted in red.

## 1174 SUPPLEMENTARY MATERIAL

1175 **Supplementary Figure 1.** Conservation of synteny between mouse and four other mammals around  
1176 *Cwc22<sup>R2d1</sup>* (upper panel) indicates that the *R2d1* sequence remains in its ancestral position. Chevrons  
1177 represent genes, alternating white and grey, and are oriented according to the strand on which the gene is  
1178 encoded. *Cwc22<sup>R2d2</sup>* is novel in the mouse but its position relative to genes with conserved order is shown  
1179 in the lower panel. Note that synteny is disrupted in mouse and rat distal to *R2d2*.

1180 **Supplementary Figure 2.** Pairwise alignment of *R2d2* contig (top) to the *R2d1* reference sequence  
1181 (bottom). Dark boxes show position of repetitive elements present in both sequences; syntenic positions  
1182 are connected by grey anchors, and blank space represents aligned bases in both sequences. Orange boxes  
1183 indicate position of repetitive elements present in the *R2d1* sequence but not detected in *R2d2*; blue boxes  
1184 indicate position of elements in *R2d2* but not *R2d1*. *Cwc22* transcripts are shown below the alignment.

1185 **Supplementary Figure 3.** Phylogenetic tree constructed from amino acid sequences for mammalian  
1186 *Cwc22* homologs (including all three mouse paralogs) with chicken as an outgroup. Node labels indicate  
1187 support in 100 bootstrap replicates.

1188 **Supplementary Figure 4.** Alignment of amino acid sequences from mouse *Cwc22<sup>R2d1</sup>*, *Cwc22<sup>R2d2</sup>* and retro-  
1189 *Cwc22*, plus *Cwc22* orthologs from 19 other placental mammals plus opossum, platypus and chicken as  
1190 outgroups. Residues are colored according to biochemical properties and gaps are shown in grey.  
1191 Information content of each column in the alignment, measured as the Jenson-Shannon divergence, is  
1192 plotted in the lower panel.

1193 **Supplementary Figure 5.** Diagnostic variants used to identify gene conversion tracts between *R2d2* and  
1194 *R2d1*. Each column represents a single variant (total of 1,411) between *R2d1* and *R2d2* for which *R2d2* has  
1195 the derived allele, and each row a single individual with whole-genome sequence data available (named  
1196 in **Supplementary Table 1**). Points are colored according to the “allele type” of each variant detected in 8  
1197 or more reads in each sample: 0:0 (red), neither *R2d1* nor *R2d2* allele present; 1:0, *R2d1* but not *R2d2*  
1198 (orange); 1:1, both *R2d1* and *R2d2* (green); 0:1, *R2d2* but not *R2d1* (blue). Individuals are grouped  
1199 according to diploid *R2d* copy number.

1200 **Supplementary Figure 6.** Partial loss of *R2d2* with structural rearrangement. (A) Inferred structure of the  
1201 *R2d1-R2d2* region in IR:AHZ\_STND:015, a wild *M. m. domesticus* individual from Iran. *R2d1* is present on  
1202 both chromosomes but only a fragment of *R2d2* remains on one chromosome, and it has been transposed  
1203 into the retro-*Cwc22* array. (B) Normalized depth of coverage (2 = normal diploid level) across *R2d*.  
1204 Regions in grey represent reads from *R2d1* alone, while region in black captures reads from *R2d1* and  
1205 *R2d2*, as shown by arrows from panel A. (C) Position of read pairs (red; not drawn to scale) with soft-  
1206 clipped alignments to *R2d1*. The proximal read aligns in the 3' UTR of *Cwc22*, and the distal read across  
1207 an exon-intron boundary within the gene body. Note the “outward”-facing direction of the alignments.  
1208 (D) Positions of the mates of the reads in panel C. Note that the x-axis is reversed so that the exons of  
1209 retro-*Cwc22* (encoded on the plus strand) parallel those of *Cwc22* (encoded on the minus strand). The 3'  
1210 read maps across the boundary of the 3' UTR of *Cwc22* and the ERV mediating the retrotransposition  
1211 event. The 5' read maps across two exon-exon boundaries in retro-*Cwc22*, so there is no ambiguity  
1212 regarding its alignment to the retro-transposed copy. (E) Inferred structure of *Cwc22* paralogs in this  
1213 sample. Note that one of the copies of retro-*Cwc22* is now a mosaic of retrotransposed and *Cwc22<sup>R2d2</sup>*-  
1214 derived sequence.

1215 **Supplementary Figure 7.** Difference between expected and observed recombination fraction between  
1216 markers flanking *R2d2* in experimental crosses in which at least one parent is segregating for a high-copy

1217 allele of *R2d2*. Thick and thin vertical bars show 90% and 95% confidence bounds, respectively, obtained  
1218 by non-parametric bootstrap.

1219 **Supplementary Figure 8. Targeted *de novo* assembly using the multi-string Burrows-Wheeler**  
1220 **Transform (msBWT).** (A) The msBWT and its associated FM-index implicitly represent a suffix array of  
1221 sequencing reads, such that read suffixes sharing a  $k$ -mer prefix are adjacent in the data structure. This  
1222 allows rapid construction of a local de Bruijn graph starting from a  $k$ -mer seed (dark blue) and extending  
1223 by successive  $k$ -mers (light blue) containing the  $(k - 1)$ -length suffix of the previous  $k$ -mer. A  $(k - 1)$ -  
1224 length prefix with more than one possible suffix (red and orange) creates a branch point. Adjacent nodes  
1225 in the graph with in-degree and out-degree one can be collapsed into a single node, yielding a simplified  
1226 graph, which can then be traversed to obtain linear contig(s). (B) Paralogs of *R2d* can be disentangled  
1227 using the local de Bruijn graph by exploiting differences in copy number. Edges in the graph are  
1228 weighted by read count, and linear contigs for the *R2d1* and *R2d2* paralogs obtained by traversing the  
1229 graph in a manner that minimizes the variance in edge weights along possible paths. Phase-informative  
1230 reads (those overlapping multiple paralogous variants) provide a second source of evidence.

1231

1232 **Supplementary Table 1.** List of mouse samples used in this study, with their taxonomic designation,  
1233 geographic origin, karyotype (STND, standard:  $2n = 40$ ; all others chromosomal races with Robertsonian  
1234 translocations) and *R2d2* copy-number classification.

1235 **Supplementary Table 2.** Transposable-element insertions private to *R2d1* or *R2d2*. Coordinates are  
1236 offsets with respect to the start position of *R2d* (for *R2d1*: chr2: 77,869,657 in the reference genome; for  
1237 *R2d2*: the beginning of the *de novo* assembled contig in **Supplementary File 1.**)

1238 **Supplementary Table 3.** Frequency table of copy-number status by geographic origin for wild-caught  
1239 and wild-derived *Mus musculus* individuals used in this study, stratified by subspecies.  
1240 “Europe/Mediterranean” includes continental Europe, the United Kingdom and countries in the  
1241 Mediterranean basin (Tunisia, Cyprus, Israel). “Asia” includes Asia, the Middle East and countries in the  
1242 Indian Ocean basin (Madagascar).

1243 **Supplementary Table 4.** Individuals from the Diversity Outbred population carrying *de novo* copy-  
1244 number mutations at *R2d2*. Each was expected to be heterozygous for the WSB/EiJ allele (33 haploid  
1245 copies).

1246 **Supplementary Table 5.** Regions of excess divergence between wild or wild-derived mice and the mouse  
1247 reference genome (GRCm38/mm10 build).

1248 **Supplementary Table 6.** Regions of *R2d* targeted for *de novo* assembly in inbred strains.

1249 **Supplementary File 1.** Compressed archive containing *R2d2* contig (from WSB/EiJ) and multiple  
1250 sequence alignments from selected regions in **Supplementary Table 6.**