

1 Population and species diversity at mouse centromere satellites

2

3 Uma P. Arora^{1,2}, Caleigh Charlebois¹, Raman Akinyanju Lawal¹, Beth L. Dumont^{1,2}

4 1 The Jackson Laboratory, 600 Main Street, Bar Harbor ME 04609

5 2 Tufts University, Graduate School of Biomedical Sciences, 136 Harrison Ave, Boston MA
6 02111

7

8 Correspondence: Uma.Arora@jax.org, Beth.Dumont@jax.org

9

10 **Abstract (150 words for elife)**

11 Centromeres are satellite-rich chromatin domains that are essential for chromosome segregation.
12 Centromere satellites evolve rapidly between species but little is known about population-level
13 diversity across these loci. We developed a *k*-mer based method to quantify centromere copy
14 number and sequence variation from whole genome sequencing data. We applied this method to
15 diverse inbred and wild house mouse (genus *Mus*) genomes and uncover pronounced variation in
16 centromere architecture between strains and populations. We show that patterns of centromere
17 diversity do not mirror the known ancestry of inbred strains, revealing a remarkably rapid rate of
18 centromere sequence evolution. We document increased satellite homogeneity and copy number
19 in inbred compared to wild mice, suggesting that inbreeding remodels mouse centromere
20 architecture. Our results highlight the power of *k*-mer based approaches for probing variation
21 across repetitive regions and provide the first in-depth, phylogenetic portrait of centromere
22 variation across *Mus musculus*.

23 Introduction

24
25 Centromeres are chromatin domains that are essential for chromosome segregation and the
26 maintenance of genome stability (Bakhoum, Thompson, Manning, & Compton, 2009; Fukagawa
27 & Earnshaw, 2014; McKinley & Cheeseman, 2016; Schalch & Steiner, 2017). Centromeres
28 serve as focal points for the assembly of the kinetochore complex, which provides the protein
29 interface linking chromosomes to microtubules during mitosis and meiosis (Bakhoum et al.,
30 2009; Fukagawa & Earnshaw, 2014; McKinley & Cheeseman, 2016; Schalch & Steiner, 2017).
31 Mutations that abolish or reduce centromere function can impair kinetochore assembly and lead
32 to spontaneous chromosome loss, cell cycle arrest, or chromosome mis-segregation (Holland &
33 Cleveland, 2009). Thus, the loss of centromere integrity can have adverse consequences for
34 genome stability, and represents an important mechanism leading to both cancer and infertility
35 (Aldrup-MacDonald, Kuo, Sullivan, Chew, & Sullivan, 2016; Barra & Fachinetti, 2018; Hudson
36 et al., 1998; Régnier et al., 2005; Zhang et al., 2016).

37
38 In most vertebrate species, centromeric DNA is comprised of tandem arrays of one or more
39 satellite repeat units (Malik & Henikoff, 2009; Rocchi, Archidiacono, Schempp, Capozzi, &
40 Stanyon, 2012; Ventura et al., 2007). As a consequence of this satellite-rich architecture,
41 centromeres are predisposed to high rates of structural mutation via replication slippage, unequal
42 exchange, and transposition (Barra & Fachinetti, 2018). These processes actively contribute to
43 size and sequence variability between species (Alexandrov, Kazakov, Tumeneva, Shepelev, &
44 Yurov, 2001; Alkan et al., 2011; Cacheux, Ponger, Gerbault-Seureau, Richard, & Escudé, 2016;
45 Musich, Brown, & Maio, 1980). For example, in mammals, centromere repeat sizes range from
46 6bp in the Chinese hamster (*Cricetulus griseus*) to 1419 bp in cattle (*Bos taurus taurus*), with
47 GC content ranging from 28-74% (Melters et al., 2013). The remarkable size and sequence
48 variability of centromeres, combined with their critical and highly conserved cellular roles in
49 chromosome segregation and genome stability, impose an enduring biological paradox.

50
51 Due to their inherent repeat-rich nature, centromeres persist as gaps on most reference genome
52 assemblies. To date, only a handful of mammalian centromeres (human chromosomes 8, X, and
53 Y) have been fully sequenced and assembled (Jain et al., 2018; Longsdon et al., 2020; Miga et
54 al., 2020). The near absence of high-quality reference sequences and the challenge of uniquely
55 anchoring short reads within repeat-rich regions pose significant barriers to the discovery and
56 analysis of genetic variation across these functionally critical regions. Consequently, the scope of
57 centromere structural and sequence diversity within and between populations remains largely
58 unknown.

59
60 Defining levels of centromere diversity represents a crucial first step towards understanding the
61 potential phenotypic consequences of variation at these loci. Prior studies in humans have
62 identified centromere variants that associate with differences in the stability of kinetochore
63 protein binding, which can, in turn, influence the fidelity of chromosome segregation (Aldrup-
64 MacDonald et al., 2016). Studies in mice and monkeyflowers (*Mimulus*) have shown that
65 differences in centromere size can lead to biased, non-Mendelian chromosome transmission in
66 heterozygotes, a phenomenon known as centromere drive (Chmátal et al., 2014; Fishman &
67 Kelly, 2015; Iwata-Otsubo et al., 2017). However, owing to an incomplete catalog of centromere
68 diversity and the omission of variants in these regions from GWAS and linkage studies (Aldrup-

69 MacDonald et al., 2016; Langley, Miga, Karpen, & Langley, 2019), the contribution of
70 centromere variation to phenotypic variation, including disease, has yet to be fully realized.

71
72 House mice (genus *Mus*) provide an ideal system for ascertaining population level centromere
73 satellite diversity and evaluating its functional consequences for several reasons. First, prior
74 investigations have identified the focal *Mus musculus* centromere satellite repeat sequences, and
75 defined core features of house mouse centromere architecture (Kalitsis, Griffiths, & Choo, 2006;
76 Kipling, Wilson, Mitchell, Taylor, & Cooke, 1994; Narayanswami et al., 1992; Wong & Rattner,
77 1988). Specifically, the *Mus musculus* centromere is composed of two primary satellite domains.
78 The minor satellite domain is a tandem array of a 120-bp sequence that cumulatively extends over
79 ~1 Mb of sequence per chromosome. This satellite array delimits the region where the centromere-
80 specific histone variant CENP-A is bound and defines the core centromere (McKinley &
81 Cheeseman, 2016). The minor satellite region is flanked by a 234-bp major satellite repeat array
82 that extends over ~2 Mb of sequence per chromosome. The major satellite region forms the
83 pericentromeric heterochromatin, which is important for sister chromatid cohesion during cell
84 division (McKinley & Cheeseman, 2016; Peters et al., 2001). Second, mouse centromeric satellite
85 arrays are reported to be homogenous both within and between chromosomes (Kalitsis et al., 2006;
86 Wong & Rattner, 1988), a feature that simplifies the task of quantifying their variation in genomes.
87 This contrasts with the architecture of human centromeres, which are composed of distinct repeat
88 arrays that form higher order repeats that vary between chromosomes (Alexandrov et al., 2001;
89 Musich et al., 1980; Wong & Rattner, 1988). Third, whole-genome sequences from diverse wild
90 and inbred *Mus musculus*, as well as more divergent *Mus* taxa are publicly available (Adams,
91 Doran, Lilue, & Keane, 2015; Harr et al., 2016; Thybert et al., 2018). These resources enable
92 surveys of centromere diversity along several dimensions, including among inbred strains, within
93 natural populations, between subspecies, and between species. Finally, as the premiere mammalian
94 biomedical model system, house mice are equipped with experimental tools and detailed
95 phenotype catalogs that can be leveraged to test the functional consequences of centromere
96 variation.

97
98 Here, we harness these strengths of the *Mus musculus* model system to carry out the first
99 sequence-based analysis of centromere diversity and evolution in mice. We couple *k*-mer based
100 bioinformatic methods with experimental approaches to uncover remarkable variation in the size
101 and sequence composition of centromeres across a panel of diverse inbred strains and wild-
102 caught house mice. Overall, our study yields a portrait of centromere satellite diversity across a
103 group of closely related mammals and lays the groundwork for future functional studies on the
104 consequences of natural genetic variation at these essential chromatin domains.

105

106 **Materials and Methods**

107
108 Whole Genome Sequencing Data
109 Illumina whole genome sequences from 100 house mouse (*Mus*) genomes were obtained in
110 binary alignment map (bam) and fastq formats from public repositories (Supplementary Table 1).
111 These samples include 33 inbred house mouse strains of predominantly *Mus musculus*
112 *domesticus* ancestry (Adams et al., 2015), 27 wild *M. m. domesticus* mice from four populations,
113 22 wild *M. m. musculus* from three populations, ten wild *M. m. castaneus* from India, eight wild
114 *M. spretus* from Spain (Harr et al., 2016), a wild-derived inbred strain of *Mus caroli*
115 (CAROLI/EiJ), and a wild-derived inbred strain of *Mus pahari* (PAHARI/EiJ) (Thybert et al.,
116 2018). *Mus caroli* and *Mus pahari* sequence reads were mapped to the *Mus musculus* reference
117 (mm10) using bwa mem version 0.7.9 (Li & Durbin, 2010). Optical duplicates were removed
118 using the *rmdups* command in samtools version 1.8 (Li et al., 2009).

119
120 **Supplementary Table 1: House mouse (*Mus*) whole genome sequence samples.** Numbers
121 represented in the data source column correspond to the following data sources: (1) The Mouse
122 Genomes Project Release 1502/REL-1502 (Adams et al., 2015); (2) (Thybert et al., 2018) ; (3)
123 Wild Mouse Genomes Project (Harr et al., 2016).

124 *We excluded libraries from these Sanger inbred strains that generated reads with length <75bp.

125 126 *k*-mer Frequencies and Normalization

127 We computed the observed frequency of all *k*-mers in each mouse genome on a per-library basis.
128 Briefly, each sequenced read in a sample's fastq file was decomposed into its constituent
129 nucleotide words of length *k*, or *k*-mers using a custom Python script (*KmerComposition.py*). We
130 selected two lengths of *k*: $k = \{15, 31\}$. These *k* values were selected to balance computational
131 speed ($k=15$) and provide high sequence specificity ($k=31$). Each analyzed genome captured 440-
132 965 million unique 15-mers and 1.1 - 14.5 billion unique 31-mers.

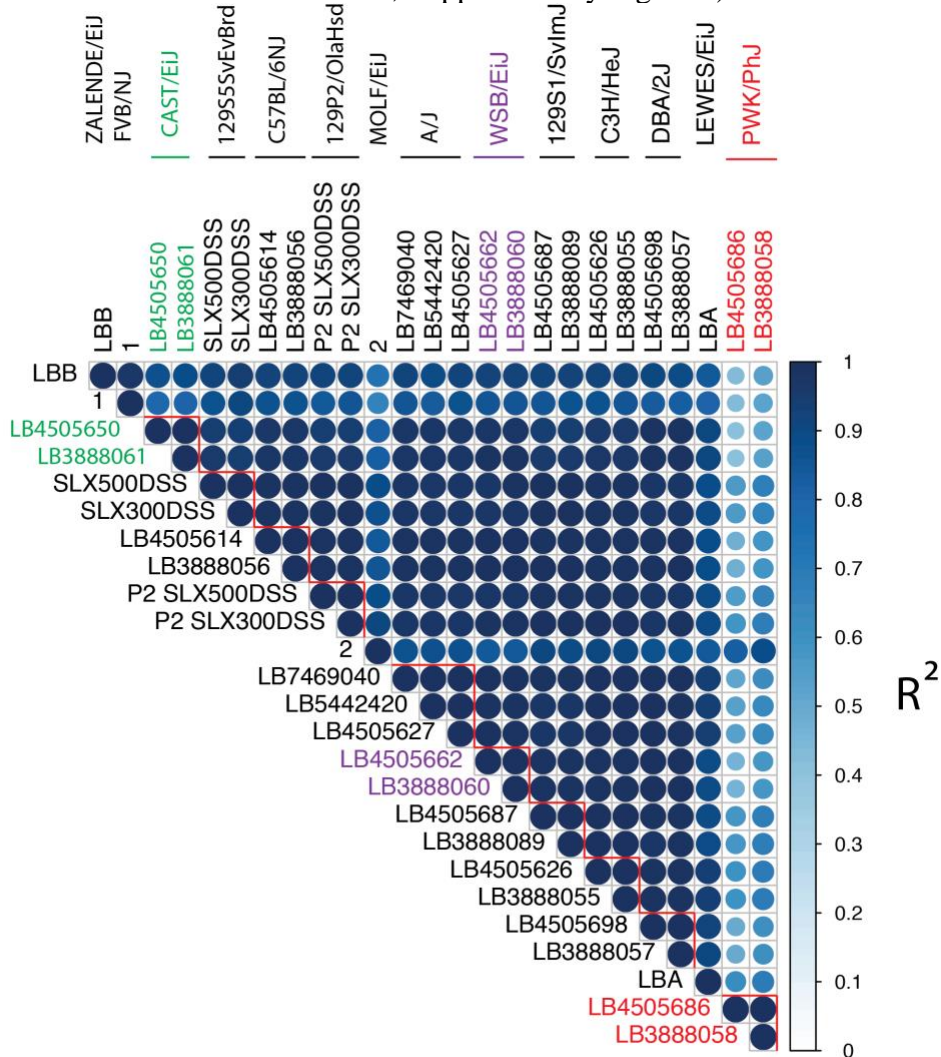
133
134 The efficiency of PCR amplification is not uniform with respect to GC-content, and this can lead
135 to biases in the nucleotide composition of sequencing libraries (Benjamini & Speed, 2012). If
136 uncorrected, such biases could cause false inference of differences in *k*-mer abundance between
137 independent libraries and samples. We implemented a GC-correction to rescale raw *k*-mer counts
138 by the extent of the observed GC-bias in each library. Briefly, we randomly selected a set of
139 ~100,000 *k*-mers that occur uniquely in the mouse reference genome (mm10). For each sample,
140 we modelled the observed counts of these unique *k*-mers as a function of their GC-content using
141 LOESS regression, with the span parameter set to 0.4. The LOESS regression produced a
142 predicted *k*-mer count for each GC-content bin; these values correspond to the magnitude and
143 direction of the empirical GC-bias in the sequencing library and represent the expected
144 "amplification" of a *k*-mer based on its GC-content. Finally, observed *k*-mer frequencies were
145 normalized by the LOESS predicted count for the corresponding GC-content bin:

$$146$$
$$147 \text{Normalized } k\text{-mer count} = \log_{10} \frac{\text{observed } k\text{-mer count}}{\text{LOESS predicted } k\text{-mer count}}$$

148 Normalized values were used for comparisons across libraries and samples.

149

150 We used reads derived from multiple independent sequencing libraries from a single inbred
 151 strain to confirm that our strategy was robust to potential artifacts introduced during library
 152 preparation. After GC-correction, we observe excellent concordance in centromere k -mer
 153 frequencies among replicate libraries for a given strain (Pearson correlation $0.990 < R^2 < 0.999$;
 154 Supplementary Figure 1). As expected, concordance between strains was generally weaker
 155 (Pearson correlation $0.41 < R^2 < 0.998$; Supplementary Figure 1).



156
 157 **Supplementary Figure 1: Concordance of GC-corrected k -mer counts among strains and**
 158 **replicate libraries within a strain.** Heatmap of pairwise Pearson correlations between GC-
 159 corrected consensus centromere 31-mer frequencies from replicate sequencing libraries across
 160 inbred *Mus musculus* strains. Both color intensity and circle size correspond to the magnitude of
 161 the R^2 correlation coefficient. Red lines delimit replicate libraries for single inbred strains.

162 Identification of highly variable k -mers across house mouse

163 To identify k -mers that differ in abundance across genomes, we selected a representative subset
 164 of $n=54$ diverse *Mus* samples (see Supplementary Table 1) and computed the variance in
 165 observed 15-mer frequencies across their genomes:
 166

167

$$\frac{\sum_{s=1}^n (F_s - \bar{F})^2}{n - 1},$$

168 where F is the absolute 15-mer frequency standardized by the read depth of strain s and \bar{F} is the
169 average normalized frequency of the 15-mer across the selected 54 strains. The 1000 15-mers
170 with the largest variance were plotted as a heatmap using the R package *pheatmap*. Clusters of
171 closely related 15-mers differing by a single nucleotide offset were manually assembled into
172 longer sequences.

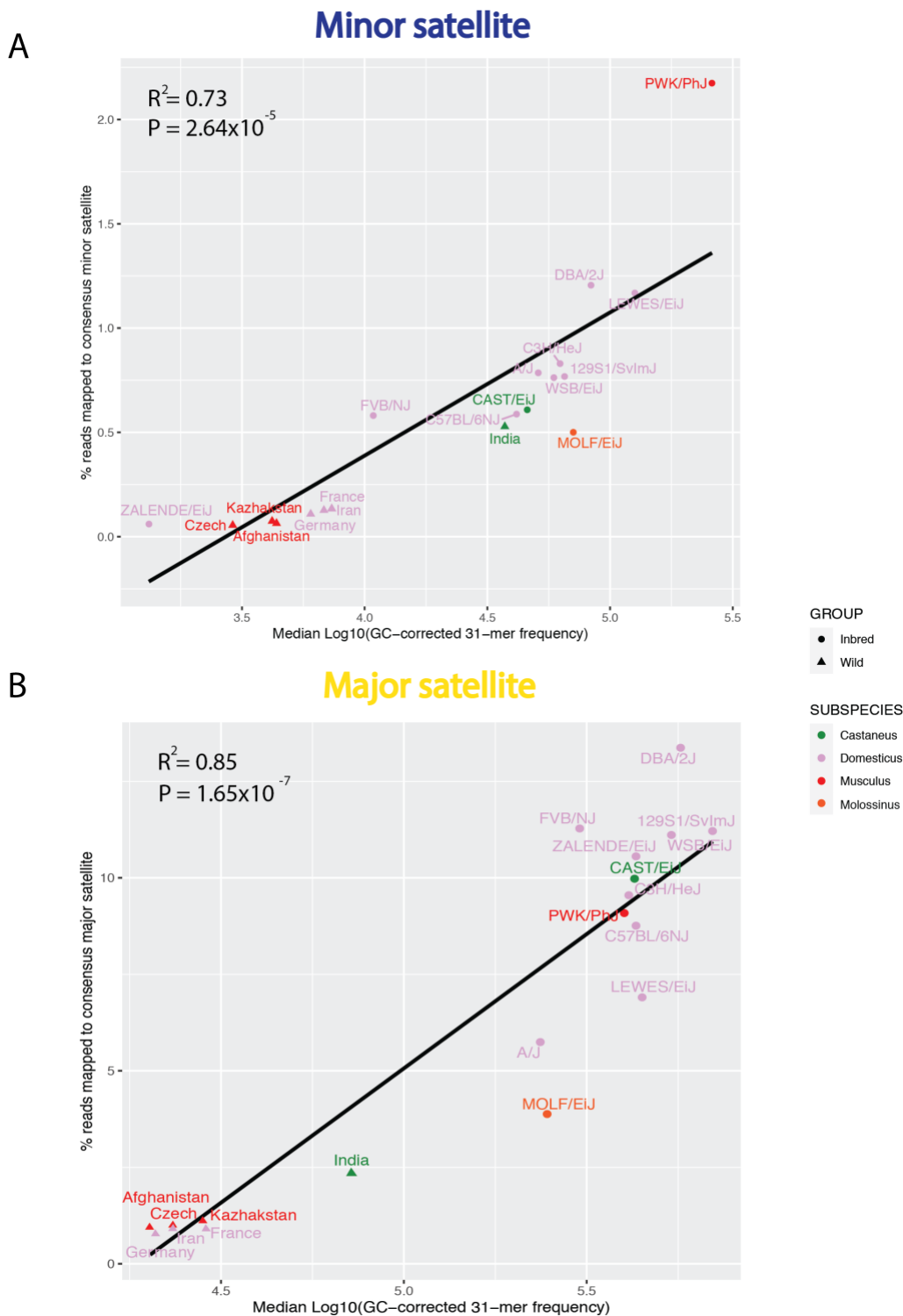
173

174 Quantifying centromere satellite abundance

175 We used a reference-informed approach to quantify the relative copy number of centromeric
176 satellites in each mouse genome. We first decomposed the minor and major satellite consensus
177 sequences into their constituent k -mers (Wong & Rattner, 1988). We then queried the GC-
178 corrected frequency of these centromere k -mers in each analyzed library and compared the
179 distribution of these centromere k -mer frequencies across libraries and samples.

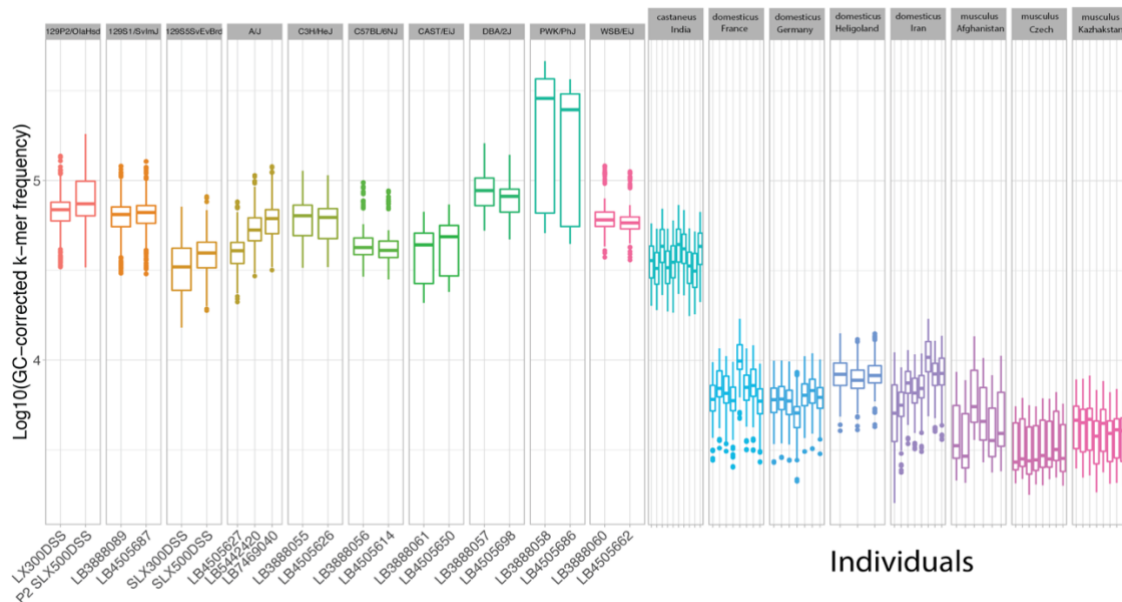
180

181 The relative copy number of a centromere satellite consensus sequence in a given mouse genome
182 was estimated from the median frequency of all constituent k -mers. For example, if the median
183 \log_{10} GC-corrected count for k -mers present in the major satellite in a given genome was 5, we
184 estimated 10^5 copies of the major satellite in that genome. This quantity is highly correlated with
185 the overall percentage of sequenced reads that map to the minor and major consensus sequences
186 (Pearson correlation; minor: $R^2 = 0.73$, $P=2.64 \times 10^{-5}$; major: $R^2 = 0.85$, $P=1.65 \times 10^{-7}$;
187 Supplementary Figure 2), suggesting that it provides a faithful readout of centromere satellite
188 copy number.



189
190 **Supplementary Figure 2: Centromere consensus 31-mer estimates of relative copy number**
191 **strongly correlate with the percentage of reads mapping to the centromere consensus.**
192 Correlation plots for median GC-corrected centromere consensus 31-mers frequency and the
193 percentage of reads mapping to the centromere consensus for the (A) minor and (B) major
194 satellite. Subspecies are represented by color. Inbred and wild-caught mice are distinguished by
195 shape.

196 We observe little variation in centromere satellite copy number among wild-caught mice
 197 sampled from a single population and replicate sequencing libraries from a single inbred strain
 198 (Supplementary Figure 3). We therefore combine all individuals from a given population to
 199 produce a single population-level copy number estimate. Similarly, for inbred strains with
 200 multiple sequencing libraries, GC-corrected *k*-mer counts were aggregated across libraries.



201 **Sequencing Libraries**
 202 **Supplementary Figure 3: Centromere consensus 31-mer frequencies exhibit low variance**
 203 **between independent sequencing libraries and among wild-caught individuals from a**
 204 **population.** Boxplots reveal the distribution of minor centromere satellite 31-mer frequencies for
 205 individual sequencing libraries or wild-caught individuals.

206 Quantifying within genome centromere satellite diversity
 207 To quantify centromere satellite diversity within a genome, we computed the average number of
 208 sequence differences between independent satellite repeats, a metric we term the *centromere*
 209 *diversity index* (CDI). We first mapped sequenced reads to the major and minor centromere
 210 consensus sequences using bwa version 0.7.9 (Li & Durbin, 2010). We then partitioned reads
 211 using samtools version 1.8 (Li et al., 2009) based on (1) whether they mapped to the major or
 212 minor satellite, (2) whether they mapped to the forward or reverse strand to prevent comparing
 213 sequences to their reverse complement, and (3) their mapped position along the consensus
 214 sequence. For each pair of reads mapping to an identical site in the same orientation on the major
 215 or minor satellite sequence, we computed the average number of observed sequence differences,
 216 d_{ij} . We then derived the CDI by averaging over all N tested read pairs:

$$\text{Centromere Diversity Index (CDI)} = \frac{\sum_{strand} \sum_L \sum_{ij} d_{ij}}{N}$$

217 where L is the length of the satellite repeat unit ($L=120$ and $L=234$ for the minor and major
 218 satellites respectively).

223 Quantifying consensus centromere satellite polymorphisms

224 To summarize the sequence polymorphism landscape across centromere satellite repeats, we
225 identified k -mers with a fixed edit distance (h) of the minor and major satellite consensus
226 sequence. For $k=15$, we allowed $h \leq 2$, and for $k=31$ we allowed $h \leq 5$. We used the frequencies
227 of these relaxed edit-distance k -mers, in conjunction with their positions across their respective
228 satellite consensus sequences, to derive a vector of relative nucleotide probabilities for each
229 position in the satellite consensus sequence. At a given position, we computed the total
230 frequency of k -mers with an “A”, “C”, “G”, or “T” at the focal position. These per-nucleotide k -
231 mer frequencies were then converted to relative probabilities summing to one and used to
232 populate a $4 \times N$ “polymorphism matrix” for each analyzed sample, where $N=120$ for the minor
233 satellite sequence and $N=234$ for the major satellite. Note that this approach ignores the
234 contribution of indel mutations to sequence polymorphism at centromere satellite repeats. We
235 then compared the percentage of non-consensus nucleotides for each strain across the minor and
236 major consensus satellite sequence.

237

238 Phylogenetic analysis of centromere diversity

239 A total of 56,500,187 high quality SNPs from 12 inbred *Mus musculus* genomes were used to
240 construct a Maximum Likelihood (ML) tree was constructed using RAxML version 8.2.12
241 (Stamatakis, 2014). The inbred strain SPRET/EiJ was included as an outgroup. An initial set of
242 20 ML trees was constructed using the GTRGAMMA substitution model. These trees were used
243 as input for subsequent branch length and topology refinements in order to estimate the tree with
244 the highest likelihood. We then used *GTRCAT* to derive bootstrap support values for the best ML
245 tree, with the number of random seeds set to 12345.

246

247 We applied Lynch’s phylogenetic comparative method to estimate the phylogenetic heritability
248 of centromere satellite copy number and CDI (Lynch, 1991). Under a neutral (*i.e.*, Brownian
249 motion) model of evolution, the extent of phenotypic divergence between species should be
250 proportional to their genetic divergence. We computed the phylogenetic variance-covariance
251 matrix from the *Mus musculus* ML phylogeny and then used this matrix to estimate the
252 proportion of variation in both major and minor satellite copy number and CDI that is explained
253 by the underlying tree. This quantity was then divided by the total variance in satellite copy
254 number and CDI to calculate the phylogenetic heritability (H_p^2) of each diversity parameter. The
255 significance of observed H_p^2 values was assessed by an *ad hoc* permutation test. We shuffled
256 observed satellite copy number and CDI values across the tree tips and then re-estimated H_p^2 on
257 each permuted dataset. Empirical P -values were determined from the quantile position of the
258 observed H_p^2 value along the distribution of 1000 permuted H_p^2 values.

259

260 All analyses were performed in R using the Analysis of Phylogenetics and Evolution (ape v5.3)
261 package (Paradis & Schliep, 2019).

262

263 Animal husbandry and ethical commitment

264 All animal experiments were approved by the Jackson Laboratory’s Animal Care and Use
265 Committee and carried out in compliance with National Institutes of Health guidelines.

266

267 The following inbred mouse strains were obtained from The Jackson Laboratory: CAST/EiJ,
268 LEWES/EiJ, PWK/PhJ, WSB/EiJ, and PAHARI/EiJ. Mice were housed in a low barrier room

269 and provided food and water *ad libitum*. Mice were euthanized by CO₂ asphyxiation or cervical
270 dislocation in accordance with recommendations from the American Veterinary Medical
271 Association.

272

273 Mouse embryonic fibroblasts cultures

274 Primary mouse embryonic fibroblasts (MEFs) were isolated from E12.5-E13.5 embryos of male
275 and female mice from four inbred strains: CAST/EiJ, LEWES/EiJ, PWK/PhJ, and WSB/EiJ.
276 MEFs were cultured in MEF media composed of Dulbecco's Modified Eagle medium (DMEM)
277 supplemented with 10% FBS (Lonza), 100ug/mL Primocin (Invivogen) and 1xGlutaMAX
278 (Thermo Fisher Scientific /GIBCO). MEFs were cultured in 150mm tissue culture-treated plates
279 (Thermo Fisher Scientific) at 37°C in a humidified atmosphere with 5% CO₂.

280

281 Metaphase chromosome spreads and FISH

282 MEFs were used for the preparation of metaphase spreads. Briefly, MEFs were cultured to ~
283 80% confluency at 37°C in a humidified atmosphere with 5% CO₂ in MEF media. Cells were
284 subsequently serum starved on MEF media without FBS and exposed to 0.02 ug/ml Colcemid
285 (Thermo Fisher Scientific/GIBCO) for 12 hours to synchronize and arrest cells in metaphase.
286 MEFs were subsequently shaken off and resuspended in hypotonic solution (56 mM KCl) for 60
287 min. The harvested cells were then gradually fixed in 3:1 Methanol:Glacial Acetic Acid under
288 constant agitation. Cells were pelleted by centrifugation, the fixative decanted off, and re-fixed
289 for a total of 3-4 times. Following the final fixation round, cells were suspended in a 1-2 mL
290 volume of fixative and dropped onto slides from a height of ~1m. Slides were allowed to air dry
291 for approximately 10 minutes and then stored at -20C until hybridization.

292

293 Commercially synthesized oligos corresponding to the *Mus musculus* major satellite, *Mus*
294 *musculus* minor satellite, and the predicted *Mus pahari* centromere sequences were fluorescently
295 labelled via nick translation (Supplementary Table 2). Briefly, 250 - 1000 ng of synthesized
296 DNA was combined with nick translation buffer (200 mM Tris pH 7.5, 500 mM MgCl₂, 5mM
297 Dithiothreitol, 500 mg/mL Bovine Serum Albumin), 0.2 mM dNTPs, 0.2mM fluorescent
298 nucleotides, 1U DNase (Promega) and 1U DNA Pol I (Thermo Fisher Scientific). Three
299 fluorescent nucleotides were used: Fluorescein-12-dUTP (Thermo Fisher Scientific),
300 ChromaTide Texas Red-12-dUTP (Thermo Fisher Scientific/Invitrogen), and Alexa Fluor 647-
301 aha-dUTP (Thermo Fisher Scientific/Invitrogen). The reaction mixture was incubated at 14.5 C
302 for 90 minutes, and then terminated by addition of 10mM EDTA. Probes ranged from 50-200 bp
303 in size, as assessed by gel electrophoresis.

304

305 **Supplementary Table 2: Sequences and primers used for FISH experiments.**

306

307 Probes were used in FISH reactions on MEF metaphase cell spreads. Probes were denatured in
308 hybridization buffer (50% formamide, 10% Dextran Sulfate, 2x saline-sodium citrate (SSC),
309 mouse Cot-1 DNA) at 72°C for 10 min and then allowed to re-anneal at 37°C until slides were
310 ready for hybridization. Slides were dehydrated in a sequential ethanol series (70%, 90%, 100%;
311 each 5 min) and dried at 42°C. Slides were then denatured in 70% formamide/2x SSC at 72°C
312 for 3 min, and immediately quenched in ice cold 70% ethanol for 5 minutes. Slides were
313 subjected to a second ethanol dehydration series (90%, 100%; each 5 min) and air dried. The
314 probe hybridization solution was then applied to the denatured slide. The hybridized region was

315 then cover-slipped and sealed with rubber cement. Hybridization reactions were allowed to
316 proceed overnight in a humidified chamber at 37°C. After gently removing the rubber cement
317 and soaking off coverslips, slides were washed 2 times in 50% formamide/2x SSC followed by
318 an additional 2 washes in 2x SSC for 5 min at room temperature. Slides were counterstained in
319 0.05ng/mL DAPI (Thermo Fisher Scientific/Invitrogen) for 10 min and air dried at room
320 temperature. Lastly, slides were mounted with ProLong Gold AntiFade (Thermo Fisher
321 Scientific/Invitrogen) and stored at -20C until imaging.

322

323 Image capture and fluorescence intensity quantification

324 FISH reactions were imaged at 63x magnification on a Leica DM6B upright fluorescent
325 microscope equipped with fluorescent filters and a cooled monochrome 2.8 megapixel digital
326 camera. Images were collected at a plane with maximal intensity using consistent exposure
327 settings across slides (DAPI at 20ms, TxRed at 50ms and FITC at 200ms). The mean intensities
328 of FISH signals at each centromere were calculated in areas drawn around centromeres based on
329 thresholding with background subtraction. Signals were quantified from all centromeres within a
330 cell (n = 40). FISH fluorescent intensity signals were collected from two independent cell lines
331 (biological replicates) from each strain and two independent experiments were conducted for
332 each cell line with fluorophores swapped for each sequence (technical replicates). We collected
333 images from 8-10 cells per replicate, amounting to >320 individual centromere measurements
334 per replicate (40 centromeres signal/cell x 8 cells ≥ 320). Differences in fluorescent intensity
335 between strains were assessed by ANOVA (baseR). Fluorescent intensity is represented in
336 arbitrary units (AU).

337

338 Evaluating signals of meiotic drive in the Diversity Outbred mapping population

339 We utilized genotype probability data from five Diversity Outbred (DO) mapping studies
340 conducted on mouse cohorts from outbreeding generations 6 to 22 (Bult MegaMUGA, Svenson-
341 183 MegaMUGA, Churchill-181 MegaMUGA, Attie-232 GigaMUGA, and Chesler-192
342 MegaMUGA; all data from [https://www.jax.org/research-and-faculty/genetic-diversity-
343 initiative/tools-data/diversity-outbred-database](https://www.jax.org/research-and-faculty/genetic-diversity-initiative/tools-data/diversity-outbred-database)). All DO mice were genotyped at a common set
344 of loci (Churchill, Gatti, Munger, & Svenson, 2012). For mice in each outbreeding generation, we
345 first determined the frequency of each parental haplotype at every genotyped marker. We then
346 looked for linked clusters of markers that exhibit a consistent departure from the expected
347 haplotype frequency (0.125) and that displayed a monotonic increase in the frequency of one or
348 more haplotypes over successive outbreeding generations. Such regions may harbor loci subject
349 to meiotic drive loci.

350

351 Mouse Phenotype Data

352 Spearman correlation tests were used to test for relationships between chromosome instability
353 phenotypes and estimated centromere satellite copy number across inbred lab strains.

354 Chromosome instability phenotypes were obtained from the Mills1 dataset deposited in the
355 Mouse Phenome Database (Bogue et al., 2020).

356 **Results**

357

358 ***k*-mer analysis reveals striking differences in the abundance of centromere satellite repeats**
359 **across *Mus***

360

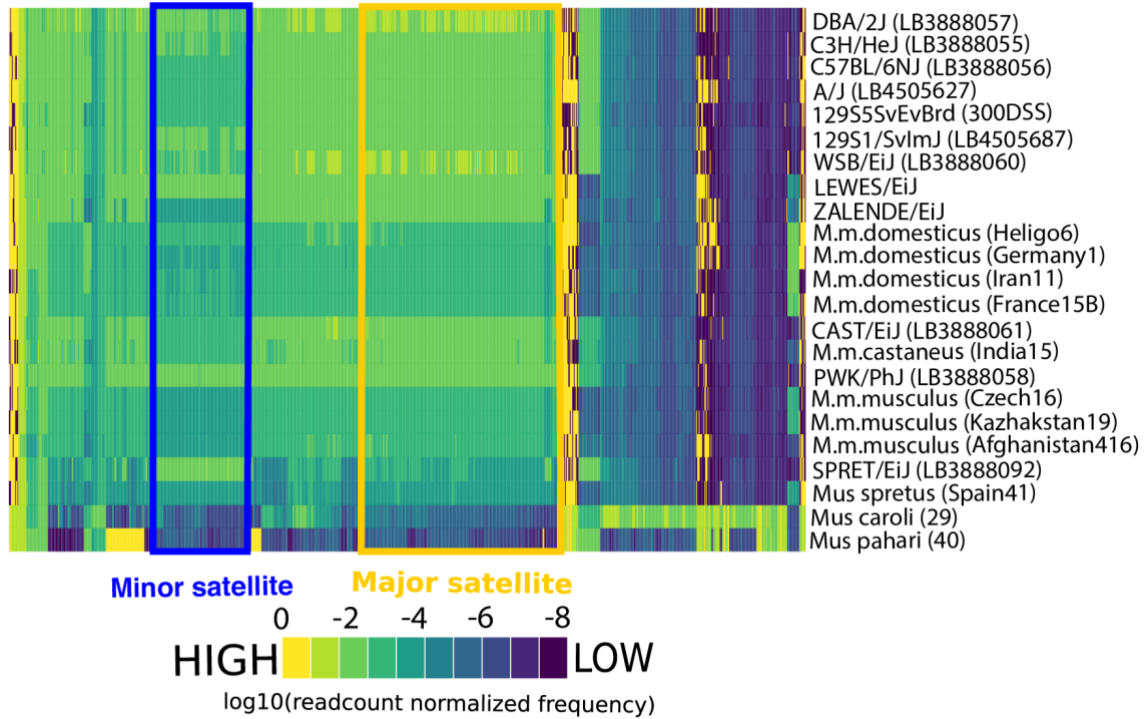
361 Standard approaches for surveying sequence diversity are not readily extendable to the
362 centromere due to its repeat-rich architecture and gapped status on the current mouse reference
363 assembly. To circumvent these challenges, we employed a *k*-mer based approach to quantify the
364 diversity of satellite DNA in mouse genomes. Our *k*-mer strategy is predicated on the insight that
365 the relative frequency of a given nucleotide word, or *k*-mer, in a shot-gun sequencing library is
366 proportional to its frequency in the parent sample genome. Thus, the observed frequency of a
367 particular *k*-mer within a pool of sequenced reads can be used as a proxy for its relative
368 abundance in a genome. We normalized *k*-mer counts to adjust for potential GC-biases
369 introduced during library preparation, and confirmed through rigorous comparisons of replicate
370 libraries for individual samples that our corrected *k*-mer counts provided a reliable readout of the
371 relative frequency of nucleotide motifs in diverse mouse genomes (Supplementary Figure 1; See
372 Methods).

373

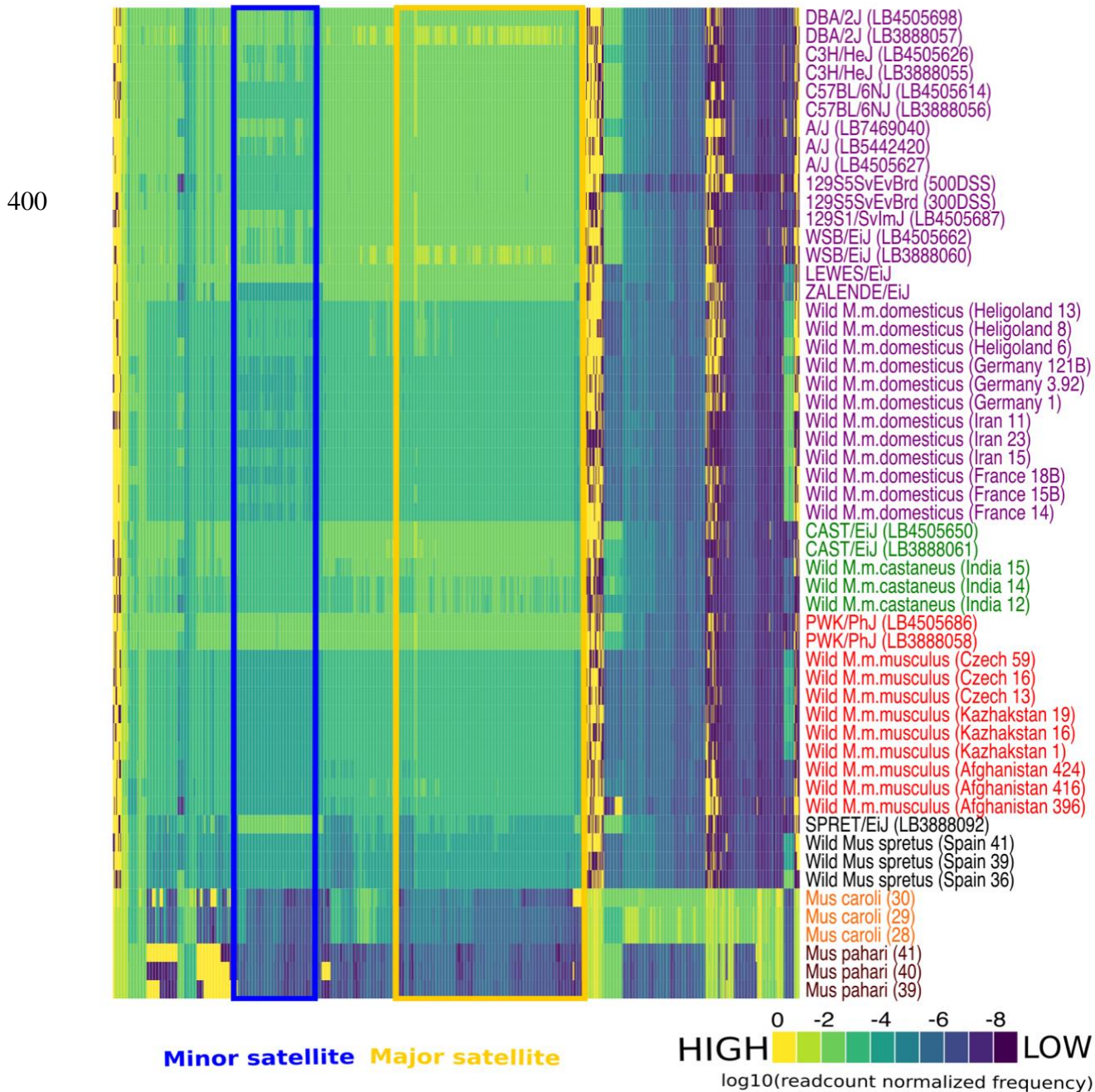
374 We first identified the most abundant 15-mers across a subset of 54 diverse mouse genomes.
375 These genomes included common inbred mouse strains, wild-caught mice from multiple
376 populations from each of the three principle house mouse subspecies (*M. m. domesticus*, *M. m.*
377 *castaneus*, and *M. m. musculus*), and three divergent *Mus* taxa (*M. spretus*, *M. caroli*, and *M.*
378 *pahari*). Consistent with prior reports (Komissarov, Gavrilova, Demin, Ishov, & Podgornaya,
379 2011), *Mus musculus* minor and major centromere satellite 15-mers were among the most
380 abundant sequences in mouse genomes (top 0.01% of all 15-mers). Interestingly, centromere 15-
381 mers were also among the most differentially abundant 15-mers across diverse *Mus musculus*
382 genomes (Figure 1), hinting at extensive centromere satellite copy number variation.

383

384 *Mus spretus* shares an identical minor satellite consensus sequence with *Mus musculus*
385 (Narayanswami et al., 1992), and exhibited a high abundance of minor satellite centromere 15-
386 mers (Figure 1). In contrast, *Mus caroli* harbors divergent centromere satellite sequences from
387 those in *M. musculus* (Kipling et al., 1995). Expectedly, we found very weak enrichment for *M.*
388 *musculus* major and minor consensus centromere 15-mer sequences in the *M. caroli* genome.
389 Similarly, we found no enrichment for *M. musculus* major and minor centromere 15-mers in *M.*
390 *pahari*, suggesting that *M. pahari* centromeres are also defined by a unique and divergent
391 satellite.



392
393 **Figure 1: Consensus centromere 15-mers were the most abundant and variable 15-mers in**
394 **diverse *Mus* genomes.** Heatmap displaying the observed frequencies of the 1000 most variable
395 15-mers (columns) across a sample of diverse *Mus* genomes (rows). Supplementary Figure 4
396 profiles a larger set of samples (n = 54). The color scale represents the frequency of all
397 centromere satellite 15-mers, normalized by the number of sequenced reads. 15-mers present in
398 the *Mus musculus* minor and major satellite consensus sequences are noted by the blue and
399 yellow boxes, respectively.

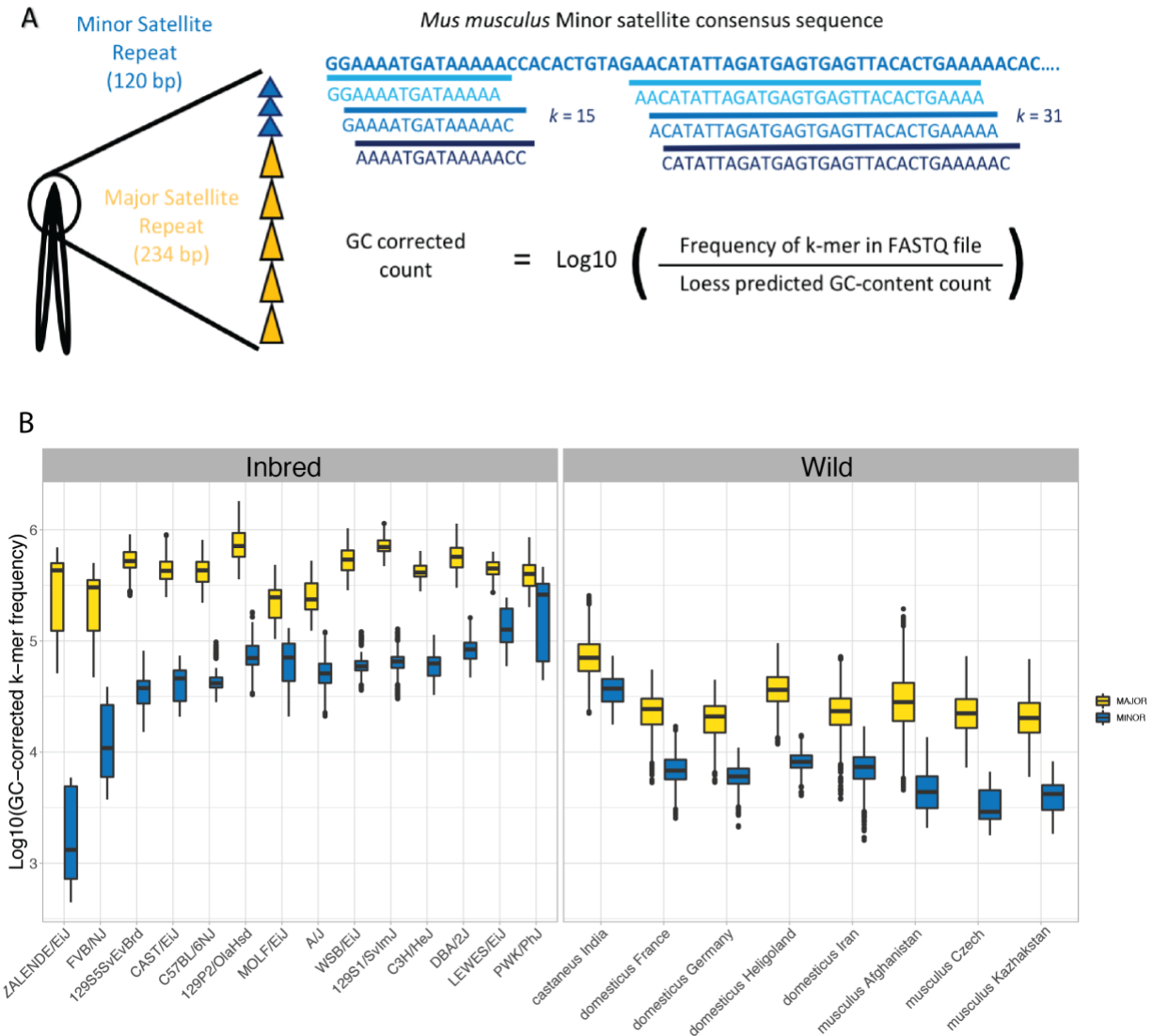


401
402 **Supplementary Figure 4: Consensus centromere 15-mers were the most abundant and**
403 **variable 15-mers in diverse *Mus* genomes.** Heatmap displaying the observed frequencies of the
404 1000 most variable 15-mers (columns) across 54 diverse samples (rows). The color scale
405 represents the frequency of all centromere satellite 15-mers, normalized by the number of
406 sequenced reads. 15-mers present in the *Mus musculus* minor and major satellite consensus
407 sequences are noted by the blue and yellow boxes, respectively.

408
409
410

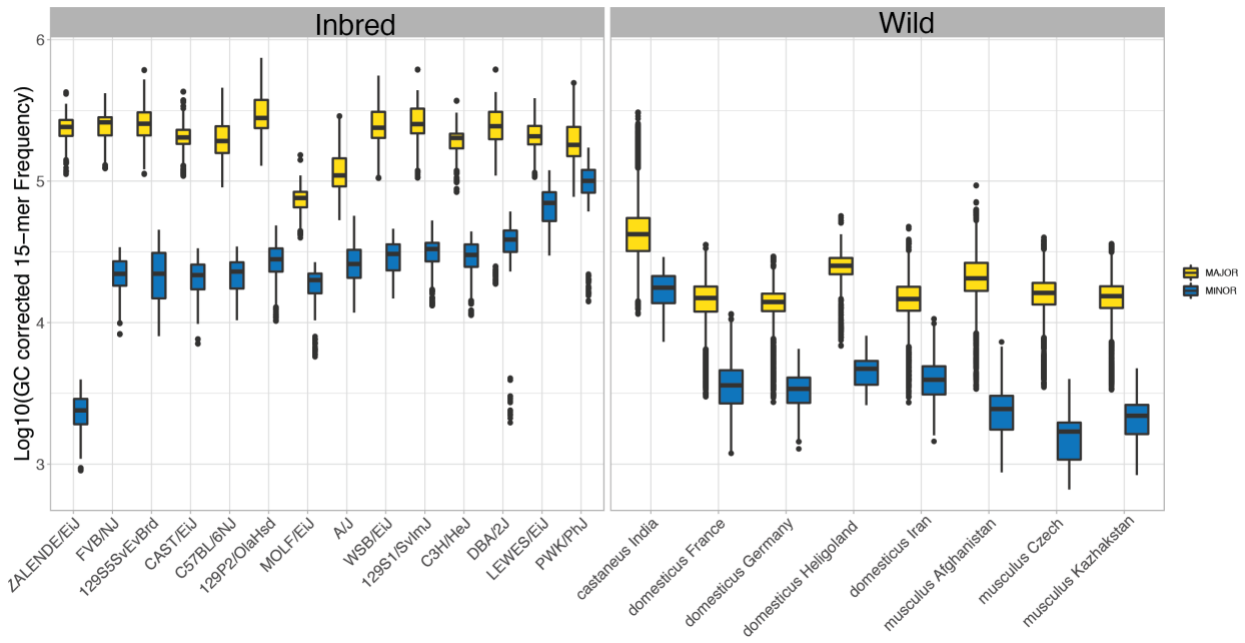
411 **Strain and population-level variation in the abundance of the *Mus musculus* consensus**
412 **centromere satellites**

413
414 Owing to the high prevalence and striking variability in the abundance of centromere satellite 15-
415 mers among *Mus* genomes, we sought to further define strain, subspecies, and species variation
416 in both major and minor satellite copy number. We compared the GC-corrected frequencies of
417 minor and major satellite 31-mers across *Mus musculus* genomes (Figure 2A) and uncovered
418 several noteworthy trends. For the following observations we present data for 31-mers but
419 observed qualitatively identical trends for 15-mers (Supplementary Figure 5). First, we
420 uncovered a greater abundance of major satellite 31-mers compared to minor satellite 31-mers
421 (Figure 2B; Kruskal-Wallis one-way ANOVA, $P < 2.2 \times 10^{-16}$). This difference is consistent with
422 the known size differences between the major and minor satellite array in *Mus musculus* (Cazaux
423 et al., 2013; Komissarov et al., 2011; Wong & Rattner, 1988). Second, there was greater strain-
424 to-strain variation in the abundance of minor satellite k -mers as compared to the major satellite k -
425 mers across the inbred strains (Kruskal-Wallis one-way ANOVA d.f. = 13, minor satellite F
426 value = 907.8, major satellite F value = 309.3, $P < 2.2 \times 10^{-16}$). We converted our normalized k -
427 mer counts into absolute satellite copy number estimates to approximate strain and subspecies
428 differences in consensus centromere size (see Methods). We estimate between 3100 – 125,000
429 minor satellite copies and 19,000 – 630,000 major satellite copies in the genomes of 14 inbred
430 *Mus musculus* strains. These estimates include only exact matches to the consensus satellite
431 sequences and ignore the potential presence of other sequence elements that modify centromere
432 size differences between samples. Nonetheless, the 40- (33-) fold range in absolute minor
433 (major) consensus satellite copy number suggests remarkable differences in centromere size
434 between closely related inbred *Mus musculus* strains. Third, inbred strains harbored higher
435 satellite 31-mer frequencies than wild-caught mice (Figure 2; Student's t -test = 212.76; $P <$
436 2.2×10^{-16}). Indeed, PCA analyses of minor and major 31-mer frequencies across diverse *Mus*
437 *musculus* samples identified inbred versus wild (i.e., outbred) as the major axes of differentiation
438 (Supplementary Figure 6). This outcome is not an artifact of systematically undercounting
439 centromeric k -mers with sequence mismatches to the consensus, as we also observe a reduced
440 fraction of reads mapping to the centromeric consensus in wild mice compared to inbred strains
441 (Supplementary Figure 2). We speculate that inbreeding may lead to the expansion of
442 centromeric repeats in house mice, similar to observations and an earlier proposal for maize
443 (Schneider, Xie, Wolfgruber, & Presting, 2016).



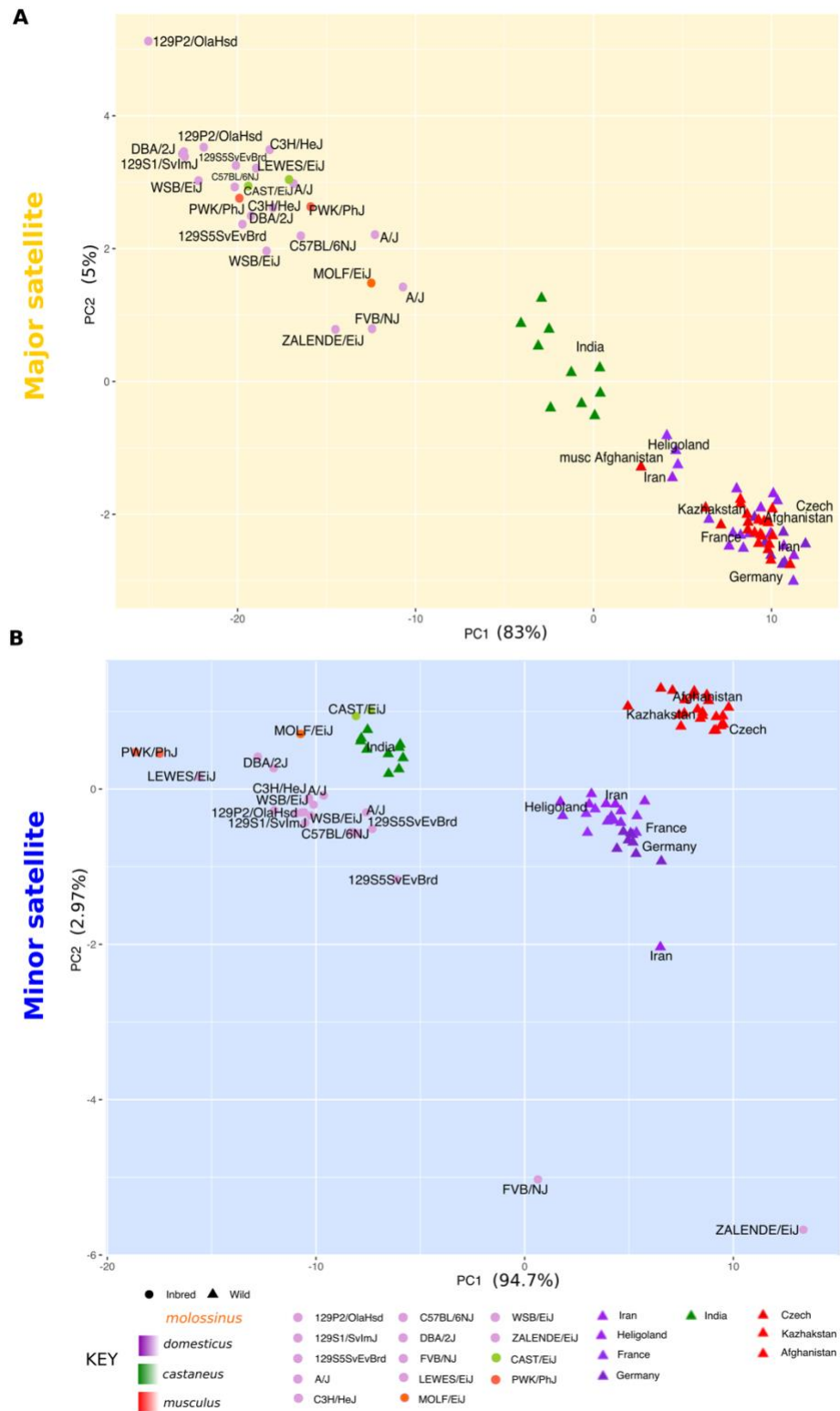
444
445
446
447
448

Figure 2: **Significant differences in consensus centromere satellite copy number across *Mus musculus*.** (A) Schematic overview of the approach used to quantify the frequencies of *k*-mers in centromeric satellite repeats. (B) Boxplots showing the distribution of major (yellow) and minor (blue) satellite 31-mer frequencies across inbred strains and wild-caught mouse populations.



449
450
451
452
453

Supplementary Figure 5: **Variation in centromere consensus 15-mer frequencies across diverse *Mus musculus* genomes.** Boxplots reveal the distribution of major (yellow) and minor (blue) satellite 15-mer frequencies across inbred strains and wild-caught mouse populations.



454
455
456
457

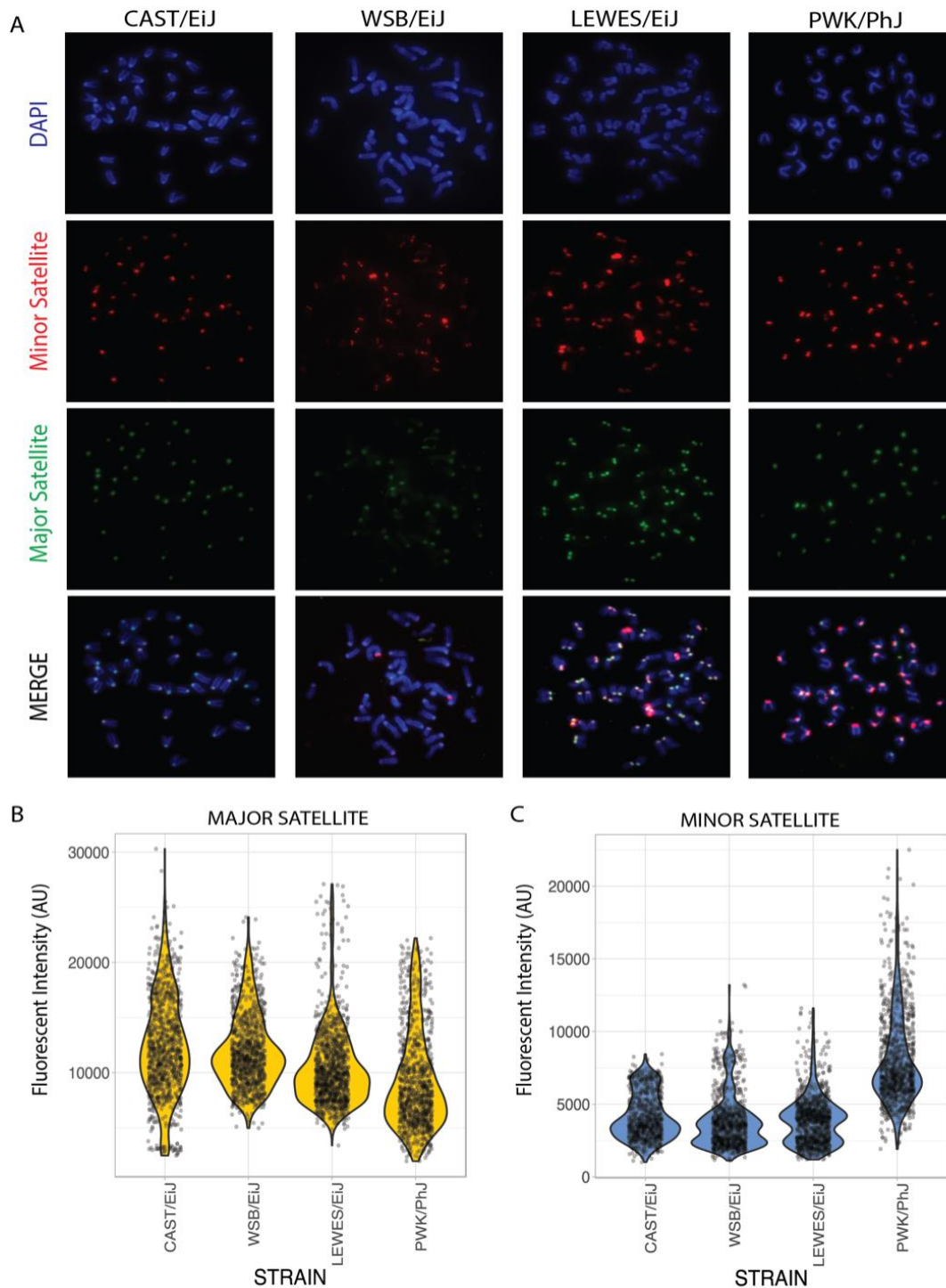
Supplementary Figure 6: **Inbred strains and wild-caught mice exhibit distinct centromere k -mer frequencies.** Principal component analysis of (A) major and (B) minor satellite centromere 31-mer frequencies in inbred strains and wild-caught *M. musculus* samples.

458 **Cytogenetic validation of strain differences in consensus centromere satellite abundance**

459

460 We used quantitative FISH (qFISH) to validate our *k*-mer based estimates of strain variation in
461 consensus centromere satellite abundance. We focused on a subset of strains that encompass a
462 range of estimated minor satellite copy numbers and span three principle house mouse
463 subspecies: CAST/EiJ (*M. m. castaneus*), WSB/EiJ (*M. m. domesticus*), LEWES/EiJ (*M. m.*
464 *domesticus*), and PWK/PhJ (*M. m. musculus*) (Figure 2B). We observed strong overall
465 concordance between relative copy number and qFISH signals at both the minor and major
466 centromere satellites (Figure 3B). Notably, both methods yielded a similar rank order of strains
467 with respect to the minor satellite abundance (median fluorescent intensity ranking in arbitrary
468 units WSB/EiJ = 3440 < CAST/EiJ = 3780 < LEWES/EiJ = 3820 < PWK/PhJ = 7055).

469 Interestingly, in WSB/EiJ and LEWES/EiJ, several chromosomes consistently showed a higher
470 minor satellite signal intensity relative to other chromosomes (Figure 3A). This observation
471 contrasted more uniform intensity of the minor satellite signal across all chromosomes in
472 CAST/EiJ and PWK/PhJ (Figure 3A). These findings point to chromosome-specific minor
473 satellite accumulation and/or loss in some *Mus musculus*, highlighting an additional dimension
474 of centromere diversity.



475
476 **Figure 3: Quantitative FISH reveals consensus centromere satellite copy number variation**
477 **across inbred mouse strains.** (A) Representative FISH images for four genetically diverse
478 inbred strains: CAST/EiJ, WSB/EiJ, LEWES/EiJ, and PWK/PhJ (B and C) Quantification of
479 fluorescent intensity using DNA probes derived from the (B) major and (C) minor centromeric
480 satellite repeats across inbred strains. Points correspond to fluorescent intensity measurements
481 for a single chromosome. A minimum of 40 centromeres from 36 cells were examined per strain.
482 Fluorescent intensity is represented in arbitrary units (AU).

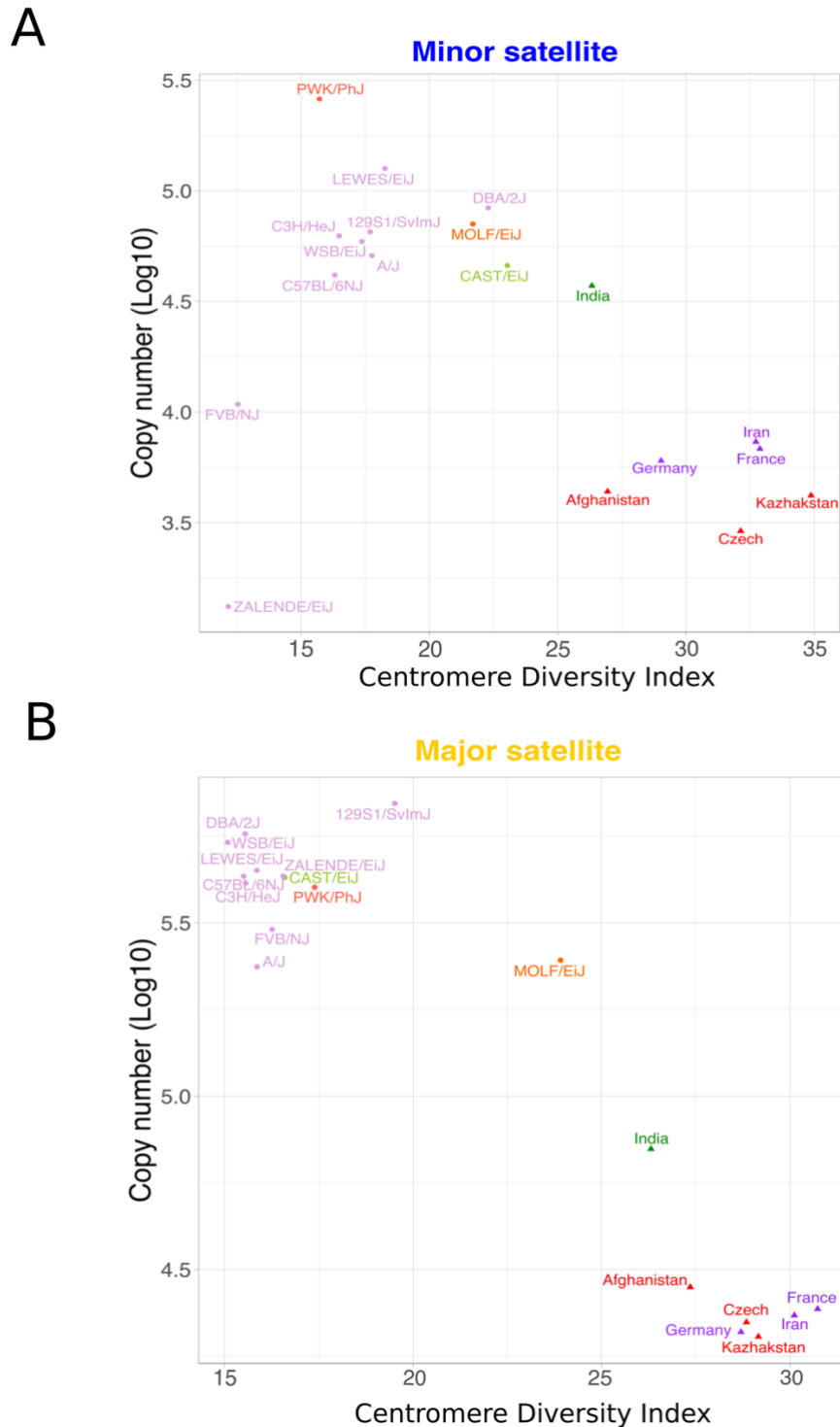
483 **Centromere satellite repeat heterogeneity at the strain and population level**

484

485 Using k -mers with exact matches to the consensus sequences, we have uncovered significant
486 variation in consensus centromere satellite copy number across diverse *Mus musculus* samples.
487 However, focusing on k -mers with exact matches to the consensus limits our ability to discover
488 and analyze centromere repeat diversity within a genome. To study this important class of
489 centromere variation, we calculated the average number of pairwise sequence differences
490 between centromere satellite repeats in each *Mus musculus* sample (see Methods). We refer to
491 this metric as the *centromere diversity index* (CDI). The minor satellite CDI across inbred *Mus*
492 *musculus* strains is lower than in wild-caught *Mus musculus* mice (inbred range: 12.1-23.0, wild
493 range: 24.8-35.6), revealing greater homogeneity of minor satellite repeats in inbred strains
494 compared to wild-caught mice. For both inbred and wild mice, the average minor satellite CDI
495 (inbred = 17.6, wild = 30.7) is slightly higher than the major satellite CDI (inbred = 16.9, wild =
496 28.7), despite the increased length and greater genomic abundance of the major satellite. Based
497 on these findings, we conclude that the mouse minor satellite harbors appreciably higher
498 sequence diversity than the major satellite.

499

500 We next assessed the relationship between centromere diversity and consensus satellite copy
501 number. There is an overall negative correlation between satellite copy number and CDI for both
502 the minor and major satellite repeats (Figure 4A and B; minor satellite: Spearman's $\rho = -0.40$,
503 $P = 0.08$; major satellite: Spearman's $\rho = -0.7$, $P = 0.001$). Samples with high satellite copy
504 number had more homogenous repeats, whereas samples with lower copy numbers had higher
505 repeat heterogeneity. This relationship is largely driven by the striking distinction between wild-
506 caught mice and inbred strains. Relative to inbred strains, wild-caught mice harbored smaller and
507 more diverse centromere arrays. The similarity in minor satellite size and diversity in inbred
508 strain CAST/EiJ and wild-caught *M. m. castaneus* represents one possible exception to this
509 pattern (Figure 4A). Together, these trends suggest that phenomena specific to inbred strain
510 genomes (or, potentially, the very process of inbreeding itself) may have influenced centromere
511 architecture in pronounced ways.



512
 513 **Figure 4: Negative correlation between centromere satellite copy number and sequence**
 514 **diversity in *Mus musculus*.** Estimated centromere satellite copy number and centromere
 515 diversity index for the (A) minor or (B) major satellite sequence. Copy number was estimated
 516 from the median frequency of consensus centromere 31-mers in each sample. The three primary
 517 house mouse subspecies are denoted by different colors: red - *M. m. musculus*, purple - *M. m.*
 518 *domesticus*, and green - *M. m. castaneus*, orange – *M. m. molossinus*. Shapes distinguish inbred
 519 strains (circles) from wild-caught mice (triangles).

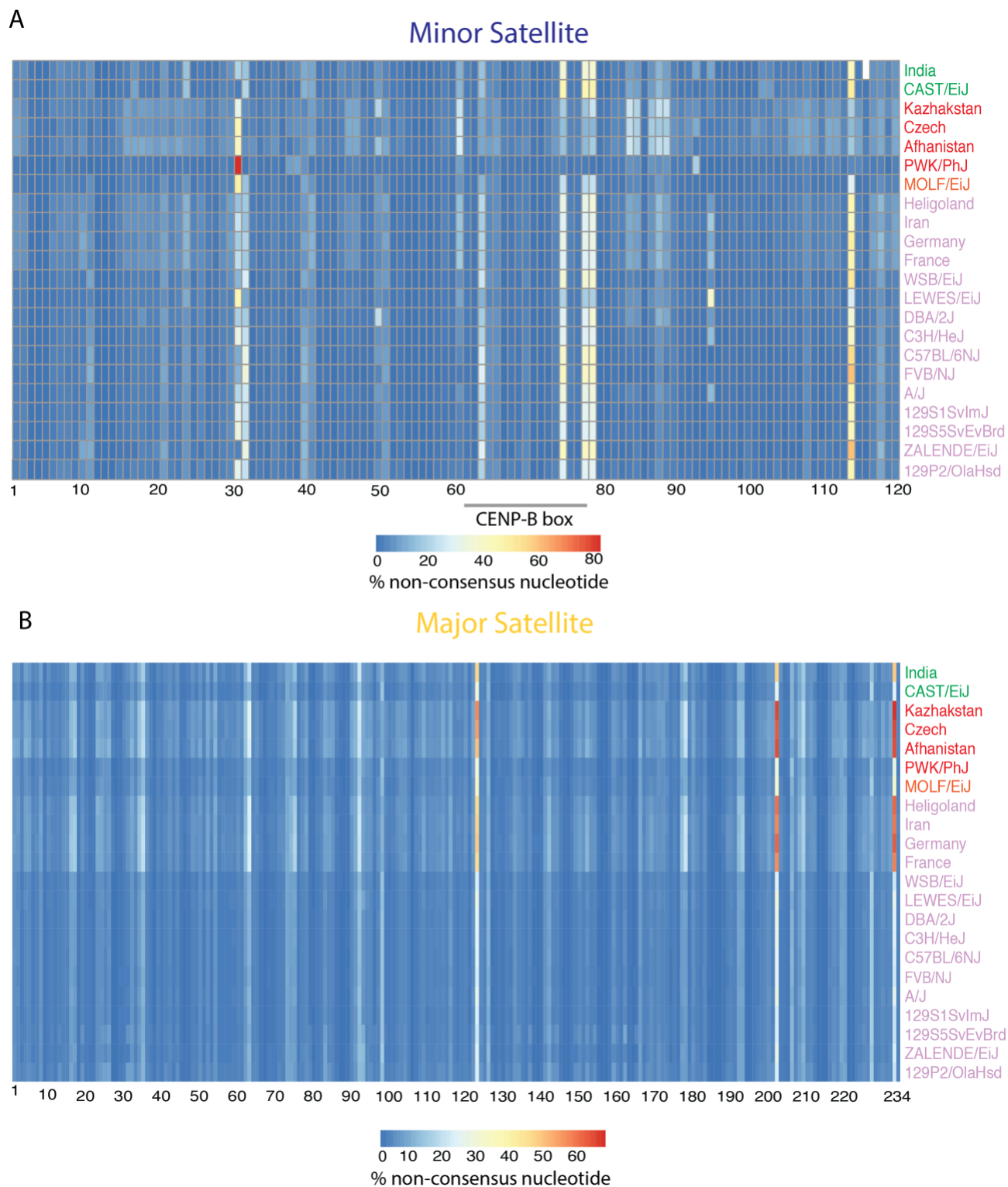
520 **Sequence landscape of *Mus musculus* satellite diversity**

521
522 Our CDI measure captures overall centromere satellite diversity within single genomes, but does
523 not pinpoint specific satellite sequence positions that subject to high variability. To investigate
524 the landscape of sequence polymorphisms along the major and minor centromere satellite
525 repeats, we relaxed the criterion for perfect k -mer matching by considering all 15- and 31-mers
526 with ≤ 2 and ≤ 5 mismatches, respectively, from the *Mus musculus* centromere satellite
527 consensus sequences. These relaxed edit-distance k -mers can be unambiguously assigned to
528 positions in the minor and major satellite consensus sequences, allowing us to quantify the
529 proportion of k -mers harboring nucleotide mismatches at each position. Using the percentage of
530 non-consensus nucleotides at each position, we then identified sites with variable nucleotide
531 usage across samples.

532
533 Overall, sequence diversity is not uniformly distributed across the minor and major satellite
534 sequences, but instead restricted to a limited number of sites that are variable between genomes
535 (Figure 5). Despite its smaller size, the minor satellite harbors more sites with at least 20% non-
536 consensus nucleotide usage than the major satellite (107 versus 79; Figures 5A and 5B).
537 Although divergence from the satellite consensus is concentrated at a minority of sites, different
538 samples vary in the frequency of non-consensus nucleotides present at a given position. For
539 example, LEWES/EiJ, WSB/EiJ and 129S1/SvImJ have similar centromere diversity indices
540 (CDI=17-18), but their minor satellite sequence landscapes are distinct from each other (Figure
541 5A).

542
543 Intriguingly, three positions within the CENP-B binding motif of the minor satellite show high
544 levels of nucleotide variability among *Mus musculus*. CENP-B binding is important but
545 dispensable for kinetochore assembly and chromosome segregation (Hudson et al., 1998).
546 Whether observed variants in the CENP-B box lead to differences in the binding efficiency of
547 CENP-B across *Mus musculus* remains unknown. Similarly, the potential functional significance
548 of nucleotide variation at other satellite positions will require future investigation.

549
550 We also uncover clear differences in the satellite sequence landscape between wild-caught mice
551 and inbred strains. On average, inbred strains have lower rates of non-consensus nucleotide
552 usage (minor satellite 2.9-4.9%; major satellite 3.6-4.4%) compared to wild-caught mice (minor
553 satellite 4.8-6.8%; major satellite 5.7-6.5%). This finding aligns with the higher CDI observed in
554 wild-caught compared to inbred mice, suggesting that wild-caught *M. musculus* have more
555 diverse and heterogenous centromere satellites than inbred strains.



556
557 **Figure 5: Landscape of nucleotide variation across centromeric satellite repeats.** Heatmap of
558 non-consensus nucleotide usage for positions in the (A) minor satellite consensus sequence and
559 (B) major satellite consensus sequence. Each row corresponds to a single sample with sample
560 names color-coded by subspecies origin: green – *M. m. castaneus*, red – *M. m. musculus*, purple
561 – *M. m. domesticus*, orange – *M. m. molossinus*.

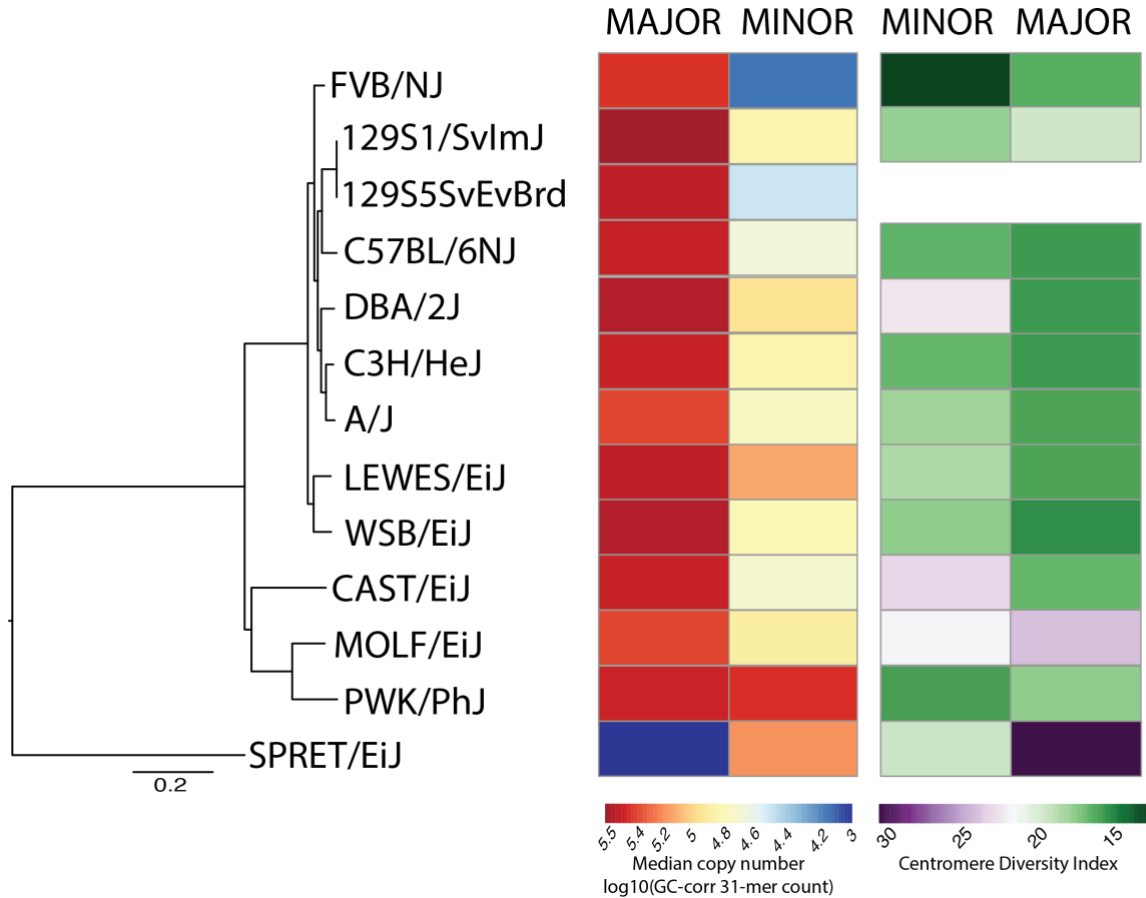
562 **Rapid evolution of minor satellite copy number and repeat heterogeneity**

563

564 To investigate how centromere architecture evolves in house mice, we analyzed the distribution
565 of centromere diversity metrics in a phylogenetic framework (Lynch, 1991). Using phylogenetic
566 comparative methods, we quantified the proportion of strain-to-strain variation in major and
567 minor satellite copy number and CDI that is explained by the underlying strain phylogeny
568 (Figure 6). We limited this analysis to inbred strains owing to the large contrast in centromere
569 architecture between inbred and wild-caught mice. Inclusion of both types of animals would
570 likely mask any legitimate phylogenetic signals within either group. We further excluded the
571 inbred strain ZALENDE/EiJ from this analysis, as it harbors multiple Robertsonian
572 chromosomal rearrangements associated with the loss of centromere satellite sequence that could
573 similarly confound the detection of phylogenetic patterns.

574

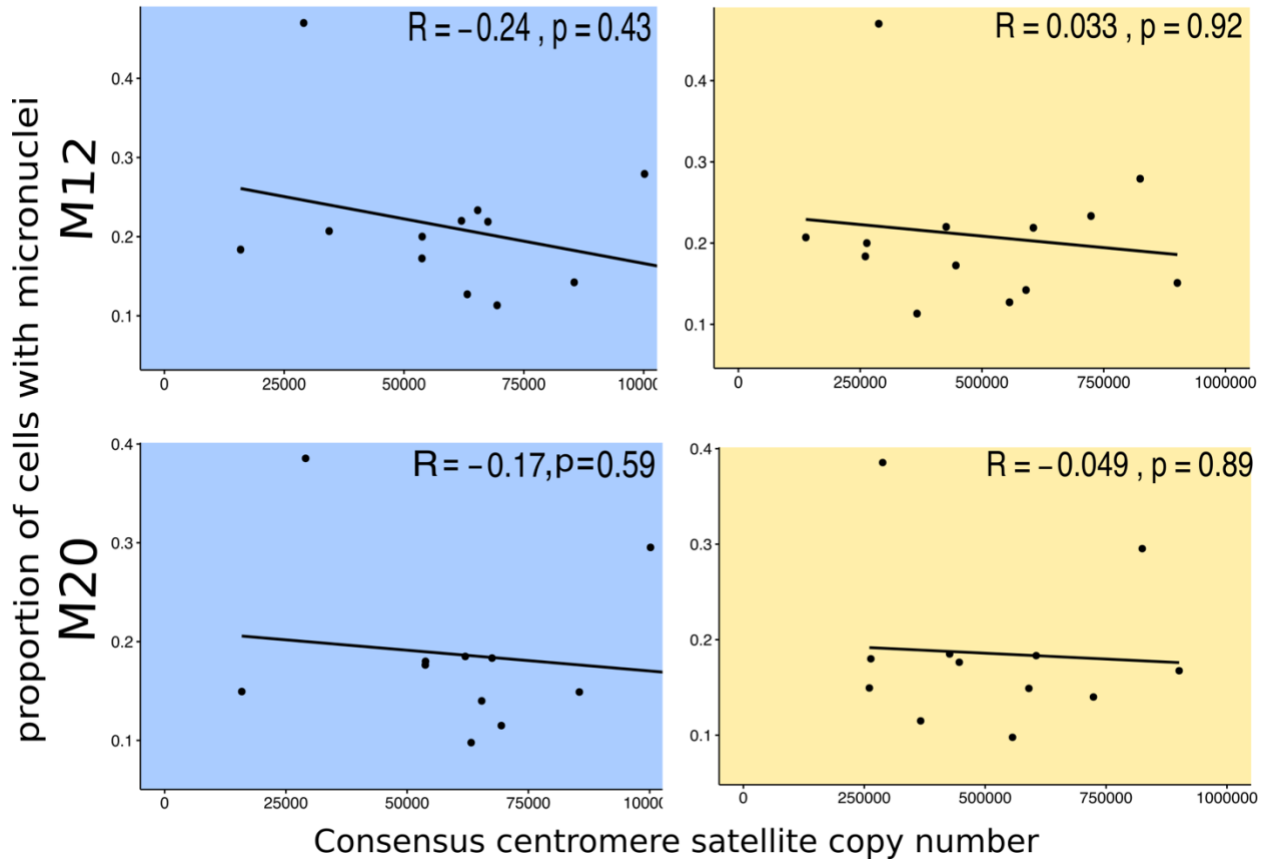
575 If variation in satellite copy number or satellite heterogeneity (*i.e.* CDI) is well-predicted from
576 the evolutionary relationships among inbred strains and *M. musculus* subspecies, these metrics
577 should exhibit a high phylogenetic heritability, H_p^2 . In contrast, if centromere satellite copy
578 number or levels of satellite heterogeneity evolve at exceptionally high rates, these measures of
579 centromere variation should exhibit a weak phylogenetic signal (*i.e.*, low H_p^2). Consistent with
580 this latter prediction, the phylogenetic heritability of both minor satellite copy number ($H_p^2 =$
581 0.45 ; $P=0.18$) and CDI ($H_p^2 = 0.15$; $P=0.21$) was low, and not significantly different from zero.
582 Evidently, both measures of minor satellite variation evolve sufficiently rapidly to outpace
583 signals of strain relatedness. In contrast, variation in both major satellite copy number ($H_p^2 =$
584 0.98 ; $P=0.24$) and CDI ($H_p^2 = 0.99$; $P=0.07$) exhibited a high, albeit non-significant, phylogenetic
585 heritability. Although modest sample sizes limit the power of this analysis, the absolute
586 differences in the H_p^2 estimates between the minor and major satellites suggests that these two
587 centromere satellites are evolving via distinct regimes, potentially mediated by differences in
588 selective pressures or mutational mechanisms.



589
 590 **Figure 6: Phylogenetic distribution of centromere satellite copy number and satellite**
 591 **diversity in inbred mice.** Maximum likelihood phylogenetic tree for 12 inbred house mouse
 592 strains and the outgroup, SPRET/EiJ. For each strain and for both the major and minor satellites,
 593 estimated satellite copy number (median 31-mer frequency) and satellite heterogeneity (CDI) are
 594 indicated by boxes shaded according to the corresponding color scales.
 595

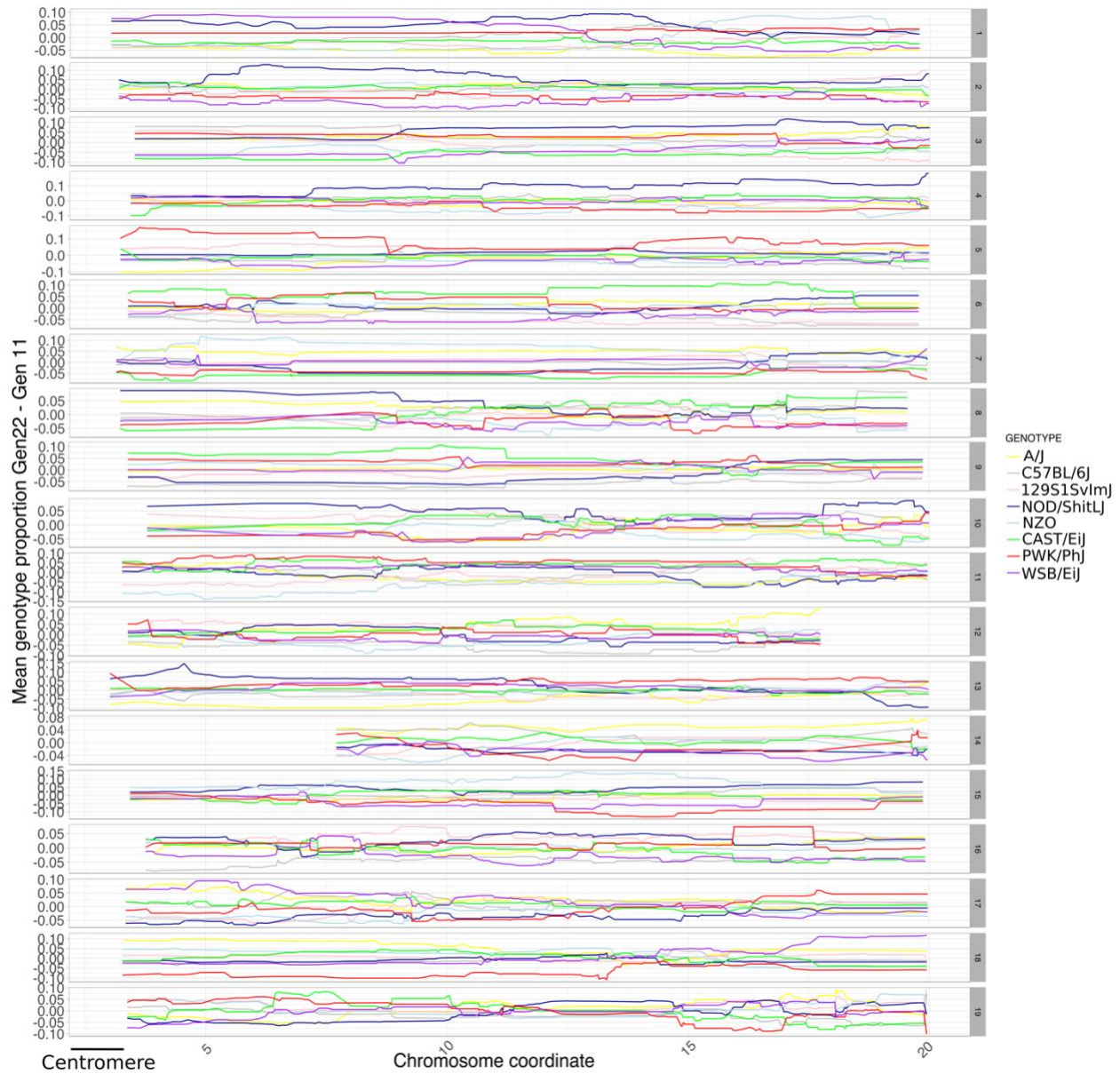
596 **Assessing the phenotypic consequences of centromere diversity in *Mus musculus***

597
 598 Centromere integrity is essential for genome stability and if not maintained can lead to cancer
 599 and infertility (Aldrup-MacDonald et al., 2016; Barra & Fachinetti, 2018; Hudson et al., 1998;
 600 Régnier et al., 2005; Zhang et al., 2016). We next asked whether observed centromere satellite
 601 diversity influences the stability of genome transmission. Using publicly available phenotype
 602 data from the Mouse Phenome Database (<https://phenome.jax.org/>) we searched for correlations
 603 between centromere satellite copy number and micronuclei formation, a hallmark of
 604 chromosome instability (Lee et al., 2013; Luzhna, Kathiria, & Kovalchuk, 2013). We found no
 605 significant correlation between this measure of genome stability and either major or minor
 606 satellite consensus copy number (Supplementary Figure 7). These results suggest no strong
 607 functional link between centromere size and this measure of chromosome instability. However,
 608 small sample sizes, uncertainty in our copy number estimates, and imprecision in the
 609 chromosomal instability phenotype may conceal true biological associations.



Supplementary Figure 7: **No correlation between micronuclei frequency and centromere satellite consensus copy number.** Spearman correlations between the proportion of peripheral blood cells (red blood cells and micronuclei) with micronuclei and median minor (left) or major (right) satellite 31-mer frequencies. The proportion of cells with micronuclei was determined for 12-month-old mice (top) and 20-month-old mice (bottom).

Centromeres are reservoirs for the accumulation of selfish drive elements (Chmátal et al., 2014; Kursel & Malik, 2018) that can hijack the inherent asymmetry of female meiosis to bias their own transmission into the oocyte (Kursel & Malik, 2018; Malik, 2009). We next asked whether centromere satellite copy number differences among inbred strains lead to systematic meiotic drive in diverse mouse populations. We profiled datasets from the Diversity Outbred (DO) mouse population (Churchill et al., 2012), a heterogenous stock population founded from 8 strains with distinct centromere satellite copy number states. We scanned genotypes of DO mice from 15 successive generations for evidence of the over-transmission of centromere-proximal alleles from one (or more) founder strain(s). We found no evidence for non-Mendelian transmission of centromere-proximal regions in the DO (Supplementary Figure 8). This result suggests (i) the absence of centromere-mediated meiotic drive in this complex population, (ii) the lack of power to detect weak drive signals, (iii) that drive is influenced by multiple genetic factors (Didion et al., 2016), or (iv) that aspects of centromere architecture other than minor satellite copy number may be critical for defining drive potential.



632
633
634
635
636
637
638
639

Supplementary Figure 8: Haplotype frequencies at centromere-proximal regions in the Diversity Outbred populations are not consistent with strong centromere drive.

Chromosome coordinates of genotyped markers in megabases (Mb) are provided on the x-axis. The difference in the frequency of each strain haplotype between generation 22 and generation 11 is shown on the y-axis. Line colors correspond to each of the 8 DO founder strains.

640 Discussion

641
642 Evolutionary theory predicts that genomic regions with key cellular roles should exhibit reduced
643 rates of evolution in order to preserve their biological function. Centromeres are paramount for
644 chromosome segregation and the maintenance of genome stability, but, paradoxically,
645 centromere satellite sequences are known to evolve rapidly between species (Alkan et al., 2011;
646 Feliciello, Akrap, Brajković, Zlatar, & Ugarković, 2014; Garrido-Ramos, 2017; Smith, 1976).
647 Despite this knowledge, comparatively little is known about the extent of centromere variation
648 over shorter evolutionary timescales, including at the population level. Here, we developed a
649 powerful *k*-mer based workflow for quantifying centromere satellite copy number and sequence
650 diversity from whole genome sequence data. We apply this analytical framework to 100
651 genomes from diverse inbred and outbred mice to characterize multiple dimensions of mouse
652 centromere variation.

653
654 Our analyses present several notable advances. First, whereas prior studies have used genomic
655 methods to survey the diversity of centromere satellite sequences across divergent taxa (Melters
656 et al., 2013), our study represents one of the first efforts to uncover the short-term evolutionary
657 dynamics of telocentric centromeres and the first in-depth study of centromere diversity in house
658 mice. Second, whereas earlier studies of centromere diversity in mice applied qualitative
659 approaches to small cohorts of inbred strains (Aker & Huang, 1996), our study provides
660 quantitative estimates of both centromere copy number and diversity across inbred strains and
661 wild-caught mice. Third, by situating our findings in a phylogenetic framework, we show that
662 centromere satellite copy number and heterogeneity are volatile and rapidly evolving properties
663 of centromere architecture. Lastly, by pursuing these investigations in a well-established model
664 system with rich phenotyping resources, our work presents an initial functional exploration of
665 observed centromere diversity and encourages further investigations of its functional effects on
666 the fidelity and dynamics of chromosome segregation.

667
668 We discovered key differences in the mode and rate of evolution of the *Mus musculus* major and
669 minor satellite sequences. Minor satellite arrays exhibited more extreme variation in copy
670 number and CDI in comparison to the major satellite arrays (Figure 2B and Figure 4). Using
671 phylogenetic comparative methods, we further showed that the rate of evolution in minor
672 satellite copy number and CDI is sufficiently rapid to erode signals of strain relatedness. In
673 contrast, major satellite copy number variation exhibits a stronger phylogenetic signal, although
674 we lack sufficient power to obtain statistically significant results (Figure 6). These differences
675 between the *Mus musculus* major and minor satellite repeat are presumably due to their distinct
676 biological functions. The major satellite repeat forms the pericentromeric heterochromatin and is
677 responsible for the establishment and maintenance of sister chromatid cohesion (McKinley &
678 Cheeseman, 2016). The minor satellite repeat binds to CENP-A, a specialized centromeric
679 histone variant responsible for kinetochore complex specification and assembly (McKinley &
680 Cheeseman, 2016). In many animal species, CENP-A is rapidly evolving, which imposes a
681 complementary selection pressure on the centromere satellite sequence to ensure protein-DNA
682 compatibility (Henikoff, Ahmad, & Malik, 2001; Malik & Henikoff, 2001; Talbert, Bryson, &
683 Henikoff, 2004). The CENP-A amino acid sequence is perfectly conserved among *M. musculus*
684 subspecies, but sequence diversity at the centromere satellite could influence the efficiency of
685 CENP-A binding, with potential downstream consequences for kinetochore assembly and

686 chromosome segregation (Iwata-Otsubo et al., 2017; Sullivan, Chew, & Sullivan, 2017). In
687 contrast, the *M. musculus* major satellite sequence does not serve as a sequence substrate for
688 kinetochore proteins. The co-evolutionary dynamics between the minor satellite DNA and
689 CENP-A have likely contributed to the accelerated evolution of the minor satellite relative to the
690 major satellite.

691
692 Our work also identified surprising differences in centromere satellite architecture between wild-
693 caught and inbred mice. Wild-caught mice exhibit lower major and minor centromere satellite
694 copy numbers and greater satellite heterogeneity than the inbred strains (Figure 5). Similar
695 observations have been previously reported for centromeres in inbred and outbred maize
696 (Schneider et al., 2016). Together, these findings suggest that the inbreeding process itself might
697 drive the homogenization of satellite arrays and facilitates the fixation of larger centromeres.
698 Indeed, prior studies have established that larger centromeres may recruit more kinetochore
699 proteins than smaller centromeres, enabling larger centromeres to selfishly bias their own
700 segregation into the oocyte during asymmetric female meiosis, a process known as centromere
701 drive (Akeru, Trimm, & Lampson, 2019; Chmátal et al., 2014; Iwata-Otsubo et al., 2017). In the
702 context of inbreeding, such “strong centromeres” should be rapidly fixed. Recurrent bouts of de
703 novo centromere expansion and fixation could lead to rapid, run-away amplification of
704 centromere satellites in inbred strains compared to wild-caught animals. Thus, centromere size
705 and repeat heterogeneity within inbred strains may not faithfully capture the native state of *M.*
706 *musculus* centromeres. Future investigations that chronical changes in centromere size from the
707 earliest stages of inbreeding onward could provide a real-time window into the mutational
708 processes that promote this architectural shift.

709
710 Our analyses define the extent of centromere copy number and sequence diversity in diverse
711 inbred strains, begging investigation into the phenotypic consequences of this variation. As an
712 initial attempt to address this outstanding challenge, we looked for correlations between satellite
713 copy number and a phenotype proxy for chromosome instability: the frequency of spontaneous
714 micronuclei formation in peripheral blood cells. We observed no significant relationship between
715 these variables, although the tested phenotype – spontaneous micronuclei formation – is likely an
716 imprecise measure of centromere-mediated genome instability (Luzhna et al., 2013).

717 Furthermore, our analysis was limited to a small number of inbred strains with available
718 published data, and is underpowered to find small to moderate strength genotype-phenotype
719 correlations. We also tested whether variation in minor satellite copy number leads to centromere
720 drive in an outbred mouse population from eight inbred strains with variable minor satellite copy
721 number (Churchill et al., 2012). We found no evidence for strong non-Mendelian transmission of
722 centromere-proximal variants although again, our analysis likely suffers from a lack of statistical
723 power to find weak to moderate drive signals. By providing the first quantitative estimates of
724 centromere satellite diversity in a panel of widely used inbred strains, our investigation critically
725 informs strain choice for future studies that aim to rigorously and explicitly test how centromere
726 diversity influences the fidelity of chromosome segregation and genome stability.

727
728 Ultra-long read sequencing technologies are now enabling sequence-level resolution of
729 mammalian centromeres (Jain et al., 2018; Longsdon et al., 2020; Miga, 2019; Miga et al.,
730 2020). However, the high cost of these methods and their labor-intensive analyses put their use
731 out of reach for most investigators and effectively limit the number of population samples that

732 can be analyzed. Our powerful *k*-mer-based workflow for assaying the architectural and
733 sequence diversity of centromeres circumvents these critical limitations and is readily extendable
734 to large numbers of genomes. Using only short read data in public repositories, our work has
735 provided key evolutionary insights into the scope of population and subspecies variation across
736 house mouse centromeres, establishing the needed foundation for functional tests of centromere
737 diversity in an important biomedical model system.

738

739 **Acknowledgements**

740 We thank Drs. Mary Ann Handel (JAX), Laura Reinholdt (JAX), and Christopher Baker (JAX)
741 for their comments while preparing the manuscript.

742

743 **Competing Interests**

744 None to declare.

745

746 **Funding**

747 This work was funded by NIGMS MIRA (R35 GM133415) awarded to BLD. UA is funded by
748 NICHD T32 (HD007065). RAL is supported by a JAX scholar Award.

749

750 **References**

- 751 Adams, D. J., Doran, A. G., Lilue, J., & Keane, T. M. (2015). The Mouse Genomes Project: a
752 repository of inbred laboratory mouse strain genomes. *Mamm Genome*, 26(9-10), 403-
753 412. doi:10.1007/s00335-015-9579-6
- 754 Aker, M., & Huang, H. V. (1996). Extreme heterogeneity of minor satellite repeat arrays in
755 inbred strains of mice. *Mamm Genome*, 7(1), 62-64.
- 756 Akera, T., Trimm, E., & Lampson, M. A. (2019). Molecular Strategies of Meiotic Cheating by
757 Selfish Centromeres. *Cell*, 178(5), 1132-1144.e1110. doi:10.1016/j.cell.2019.07.001
- 758 Aldrup-MacDonald, M. E., Kuo, M. E., Sullivan, L. L., Chew, K., & Sullivan, B. A. (2016).
759 Genomic variation within alpha satellite DNA influences centromere location on human
760 chromosomes with metastable epialleles. *Genome Res*, 26(10), 1301-1311.
761 doi:10.1101/gr.206706.116
- 762 Alexandrov, I., Kazakov, A., Tumeneva, I., Shepelev, V., & Yurov, Y. (2001). Alpha-satellite
763 DNA of primates: old and new families. *Chromosoma*, 110(4), 253-266.
764 doi:10.1007/s004120100146
- 765 Alkan, C., Cardone, M. F., Catacchio, C. R., Antonacci, F., O'Brien, S. J., Ryder, O. A., . . .
766 Ventura, M. (2011). Genome-wide characterization of centromeric satellites from
767 multiple mammalian genomes. *Genome Res*, 21(1), 137-145. doi:10.1101/gr.111278.110
- 768 Bakhoun, S. F., Thompson, S. L., Manning, A. L., & Compton, D. A. (2009). Genome stability
769 is ensured by temporal control of kinetochore-microtubule dynamics. *Nat Cell Biol*,
770 11(1), 27-35. doi:10.1038/ncb1809
- 771 Barra, V., & Fachinetti, D. (2018). The dark side of centromeres: types, causes and consequences
772 of structural abnormalities implicating centromeric DNA. *Nat Commun*, 9(1), 4340.
773 doi:10.1038/s41467-018-06545-y
- 774 Benjamini, Y., & Speed, T. P. (2012). Summarizing and correcting the GC content bias in high-
775 throughput sequencing. *Nucleic Acids Res*, 40(10), e72. doi:10.1093/nar/gks001
- 776 Bogue, M. A., Philip, V. M., Walton, D. O., Grubb, S. C., Dunn, M. H., Kolishovski, G., . . .
777 Chesler, E. J. (2020). Mouse Phenome Database: a data repository and analysis suite for
778 curated primary mouse phenotype data. *Nucleic Acids Res*, 48(D1), D716-D723.
779 doi:10.1093/nar/gkz1032
- 780 Cacheux, L., Ponger, L., Gerbault-Seureau, M., Richard, F. A., & Escudé, C. (2016). Diversity
781 and distribution of alpha satellite DNA in the genome of an Old World monkey:
782 *Cercopithecus solatus*. *BMC Genomics*, 17(1), 916. doi:10.1186/s12864-016-3246-5
- 783 Cazaux, B., Catalan, J., Justy, F., Escudé, C., Desmarais, E., & Britton-Davidian, J. (2013).
784 Evolution of the structure and composition of house mouse satellite DNA sequences in
785 the subgenus *Mus* (Rodentia: Muridea): a cytogenomic approach. *Chromosoma*, 122(3),
786 209-220. doi:10.1007/s00412-013-0402-4
- 787 Chmátal, L., Gabriel, S. I., Mitsainas, G. P., Martínez-Vargas, J., Ventura, J., Searle, J. B., . . .
788 Lampson, M. A. (2014). Centromere strength provides the cell biological basis for
789 meiotic drive and karyotype evolution in mice. *Curr Biol*, 24(19), 2295-2300.
790 doi:10.1016/j.cub.2014.08.017
- 791 Churchill, G. A., Gatti, D. M., Munger, S. C., & Svenson, K. L. (2012). The Diversity Outbred
792 mouse population. *Mamm Genome*, 23(9-10), 713-718. doi:10.1007/s00335-012-9414-2
- 793 Didion, J. P., Morgan, A. P., Yadgary, L., Bell, T. A., McMullan, R. C., Ortiz de Solorzano, L., .
794 . . Pardo-Manuel de Villena, F. (2016). R2d2 Drives Selfish Sweeps in the House Mouse.
795 *Mol Biol Evol*, 33(6), 1381-1395. doi:10.1093/molbev/msw036

- 796 Feliciello, I., Akrap, I., Brajković, J., Zlatar, I., & Ugarković, Đ. (2014). Satellite DNA as a
797 driver of population divergence in the red flour beetle *Tribolium castaneum*. *Genome*
798 *Biol Evol*, 7(1), 228-239. doi:10.1093/gbe/evu280
- 799 Fishman, L., & Kelly, J. K. (2015). Centromere-associated meiotic drive and female fitness
800 variation in *Mimulus*. *Evolution*, 69(5), 1208-1218. doi:10.1111/evo.12661
- 801 Fukagawa, T., & Earnshaw, W. C. (2014). The centromere: chromatin foundation for the
802 kinetochore machinery. *Dev Cell*, 30(5), 496-508. doi:10.1016/j.devcel.2014.08.016
- 803 Garrido-Ramos, M. A. (2017). Satellite DNA: An Evolving Topic. *Genes (Basel)*, 8(9).
804 doi:10.3390/genes8090230
- 805 Harr, B., Karakoc, E., Neme, R., Teschke, M., Pfeifle, C., Pezer, Ž., . . . Tautz, D. (2016).
806 Genomic resources for wild populations of the house mouse, *Mus musculus* and its close
807 relative *Mus spretus*. *Sci Data*, 3, 160075. doi:10.1038/sdata.2016.75
- 808 Henikoff, S., Ahmad, K., & Malik, H. S. (2001). The centromere paradox: stable inheritance
809 with rapidly evolving DNA. *Science*, 293(5532), 1098-1102.
810 doi:10.1126/science.1062939
- 811 Holland, A. J., & Cleveland, D. W. (2009). Boveri revisited: chromosomal instability,
812 aneuploidy and tumorigenesis. *Nat Rev Mol Cell Biol*, 10(7), 478-487.
813 doi:10.1038/nrm2718
- 814 Hudson, D. F., Fowler, K. J., Earle, E., Saffery, R., Kalitsis, P., Trowell, H., . . . Choo, K. H.
815 (1998). Centromere protein B null mice are mitotically and meiotically normal but have
816 lower body and testis weights. *J Cell Biol*, 141(2), 309-319.
- 817 Iwata-Otsubo, A., Dawicki-McKenna, J. M., Akera, T., Falk, S. J., Chmátal, L., Yang, K., . . .
818 Black, B. E. (2017). Expanded Satellite Repeats Amplify a Discrete CENP-A
819 Nucleosome Assembly Site on Chromosomes that Drive in Female Meiosis. *Curr Biol*,
820 27(15), 2365-2373.e2368. doi:10.1016/j.cub.2017.06.069
- 821 Jain, M., Olsen, H. E., Turner, D. J., Stoddart, D., Bulazel, K. V., Paten, B., . . . Miga, K. H.
822 (2018). Linear assembly of a human centromere on the Y chromosome. *Nat Biotechnol*,
823 36(4), 321-323. doi:10.1038/nbt.4109
- 824 Kalitsis, P., Griffiths, B., & Choo, K. H. (2006). Mouse telocentric sequences reveal a high rate
825 of homogenization and possible role in Robertsonian translocation. *Proc Natl Acad Sci U*
826 *S A*, 103(23), 8786-8791. doi:10.1073/pnas.0600250103
- 827 Kipling, D., Mitchell, A. R., Masumoto, H., Wilson, H. E., Nicol, L., & Cooke, H. J. (1995).
828 CENP-B binds a novel centromeric sequence in the Asian mouse *Mus caroli*. *Mol Cell*
829 *Biol*, 15(8), 4009-4020.
- 830 Kipling, D., Wilson, H. E., Mitchell, A. R., Taylor, B. A., & Cooke, H. J. (1994). Mouse
831 centromere mapping using oligonucleotide probes that detect variants of the minor
832 satellite. *Chromosoma*, 103(1), 46-55.
- 833 Komissarov, A. S., Gavrilova, E. V., Demin, S. J., Ishov, A. M., & Podgornaya, O. I. (2011).
834 Tandemly repeated DNA families in the mouse genome. *BMC Genomics*, 12, 531.
835 doi:10.1186/1471-2164-12-531
- 836 Kursel, L. E., & Malik, H. S. (2018). The cellular mechanisms and consequences of centromere
837 drive. *Curr Opin Cell Biol*, 52, 58-65. doi:10.1016/j.ceb.2018.01.011
- 838 Langley, S. A., Miga, K. H., Karpen, G. H., & Langley, C. H. (2019). Haplotypes spanning
839 centromeric regions reveal persistence of large blocks of archaic DNA. *Elife*, 8.
840 doi:10.7554/eLife.42989

- 841 Lee, H. S., Lee, N. C., Grimes, B. R., Samoshkin, A., Kononenko, A. V., Bansal, R., . . .
842 Larionov, V. (2013). A new assay for measuring chromosome instability (CIN) and
843 identification of drugs that elevate CIN in cancer cells. *BMC Cancer*, *13*, 252.
844 doi:10.1186/1471-2407-13-252
- 845 Li, H., & Durbin, R. (2010). Fast and accurate long-read alignment with Burrows-Wheeler
846 transform. *Bioinformatics*, *26*(5), 589-595. doi:10.1093/bioinformatics/btp698
- 847 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., . . . Subgroup, G. P. D. P.
848 (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, *25*(16),
849 2078-2079. doi:10.1093/bioinformatics/btp352
- 850 Longsdon, G. A., Voller, M. R., Hsieh, P., Mao, Y., Liskovych, M. A., Koren, S., . . . Eichler, E.
851 E. (2020). The structure, function, and evolution of a complete human chromosome 8.
- 852 Luzhna, L., Kathiria, P., & Kovalchuk, O. (2013). Micronuclei in genotoxicity assessment: from
853 genetics to epigenetics and beyond. *Front Genet*, *4*, 131. doi:10.3389/fgene.2013.00131
- 854 Lynch, M. (1991). METHODS FOR THE ANALYSIS OF COMPARATIVE DATA IN
855 EVOLUTIONARY BIOLOGY. *Evolution*, *45*(5), 1065-1080. doi:10.1111/j.1558-
856 5646.1991.tb04375.x
- 857 Malik, H. S. (2009). The centromere-drive hypothesis: a simple basis for centromere complexity.
858 *Prog Mol Subcell Biol*, *48*, 33-52. doi:10.1007/978-3-642-00182-6_2
- 859 Malik, H. S., & Henikoff, S. (2001). Adaptive evolution of Cid, a centromere-specific histone in
860 *Drosophila*. *Genetics*, *157*(3), 1293-1298.
- 861 Malik, H. S., & Henikoff, S. (2009). Major evolutionary transitions in centromere complexity.
862 *Cell*, *138*(6), 1067-1082. doi:10.1016/j.cell.2009.08.036
- 863 McKinley, K. L., & Cheeseman, I. M. (2016). The molecular basis for centromere identity and
864 function. *Nat Rev Mol Cell Biol*, *17*(1), 16-29. doi:10.1038/nrm.2015.5
- 865 Melters, D. P., Bradnam, K. R., Young, H. A., Telis, N., May, M. R., Ruby, J. G., . . . Chan, S.
866 W. (2013). Comparative analysis of tandem repeats from hundreds of species reveals
867 unique insights into centromere evolution. *Genome Biol*, *14*(1), R10. doi:10.1186/gb-
868 2013-14-1-r10
- 869 Miga, K. H. (2019). Centromeric Satellite DNAs: Hidden Sequence Variation in the Human
870 Population. *Genes (Basel)*, *10*(5). doi:10.3390/genes10050352
- 871 Miga, K. H., Koren, S., Rhie, A., Vollger, M. R., Gershman, A., Bzikadze, A., . . . Phillippy, A.
872 M. (2020). Telomere-to-telomere assembly of a complete human X chromosome. *Nature*.
873 doi:10.1038/s41586-020-2547-7
- 874 Musich, P. R., Brown, F. L., & Maio, J. J. (1980). Highly repetitive component alpha and related
875 alphoid DNAs in man and monkeys. *Chromosoma*, *80*(3), 331-348.
876 doi:10.1007/BF00292688
- 877 Narayanswami, S., Doggett, N. A., Clark, L. M., Hildebrand, C. E., Weier, H. U., & Hamkalo, B.
878 A. (1992). Cytological and molecular characterization of centromeres in *Mus domesticus*
879 and *Mus spretus*. *Mamm Genome*, *2*(3), 186-194.
- 880 Paradis, E., & Schliep, K. (2019). ape 5.0: an environment for modern phylogenetics and
881 evolutionary analyses in R. *Bioinformatics*, *35*(3), 526-528.
882 doi:10.1093/bioinformatics/bty633
- 883 Peters, A. H., O'Carroll, D., Scherthan, H., Mechtler, K., Sauer, S., Schöfer, C., . . . Jenuwein, T.
884 (2001). Loss of the Suv39h histone methyltransferases impairs mammalian
885 heterochromatin and genome stability. *Cell*, *107*(3), 323-337. doi:10.1016/s0092-
886 8674(01)00542-6

- 887 Rocchi, M., Archidiacono, N., Schempp, W., Capozzi, O., & Stanyon, R. (2012). Centromere
888 repositioning in mammals. *Heredity (Edinb)*, *108*(1), 59-67. doi:10.1038/hdy.2011.101
- 889 Régnier, V., Vagnarelli, P., Fukagawa, T., Zerjal, T., Burns, E., Trouche, D., . . . Brown, W.
890 (2005). CENP-A is required for accurate chromosome segregation and sustained
891 kinetochore association of BubR1. *Mol Cell Biol*, *25*(10), 3967-3981.
892 doi:10.1128/MCB.25.10.3967-3981.2005
- 893 Schalch, T., & Steiner, F. A. (2017). Structure of centromere chromatin: from nucleosome to
894 chromosomal architecture. *Chromosoma*, *126*(4), 443-455. doi:10.1007/s00412-016-
895 0620-7
- 896 Schneider, K. L., Xie, Z., Wolfgruber, T. K., & Presting, G. G. (2016). Inbreeding drives maize
897 centromere evolution. *Proc Natl Acad Sci U S A*, *113*(8), E987-996.
898 doi:10.1073/pnas.1522008113
- 899 Smith, G. P. (1976). Evolution of repeated DNA sequences by unequal crossover. *Science*,
900 *191*(4227), 528-535. doi:10.1126/science.1251186
- 901 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of
902 large phylogenies. *Bioinformatics*, *30*(9), 1312-1313. doi:10.1093/bioinformatics/btu033
- 903 Sullivan, L. L., Chew, K., & Sullivan, B. A. (2017). α satellite DNA variation and function of the
904 human centromere. *Nucleus*, *8*(4), 331-339. doi:10.1080/19491034.2017.1308989
- 905 Talbert, P. B., Bryson, T. D., & Henikoff, S. (2004). Adaptive evolution of centromere proteins
906 in plants and animals. *J Biol*, *3*(4), 18. doi:10.1186/jbiol11
- 907 Thybert, D., Roller, M., Navarro, F. C. P., Fiddes, I., Streeter, I., Feig, C., . . . Flicek, P. (2018).
908 Repeat associated mechanisms of genome evolution and function revealed by the *Mus*
909 *caroli* and *Mus pahari* Genomes. *Genome Res*, *28*(4), 448-459.
910 doi:10.1101/gr.234096.117
- 911 Ventura, M., Antonacci, F., Cardone, M. F., Stanyon, R., D'Addabbo, P., Cellamare, A., . . .
912 Rocchi, M. (2007). Evolutionary formation of new centromeres in macaque. *Science*,
913 *316*(5822), 243-246. doi:10.1126/science.1140615
- 914 Wong, A. K., & Rattner, J. B. (1988). Sequence organization and cytological localization of the
915 minor satellite of mouse. *Nucleic Acids Res*, *16*(24), 11645-11661.
- 916 Zhang, W., Mao, J. H., Zhu, W., Jain, A. K., Liu, K., Brown, J. B., & Karpen, G. H. (2016).
917 Centromere and kinetochore gene misexpression predicts cancer patient survival and
918 response to radiotherapy and chemotherapy. *Nat Commun*, *7*, 12619.
919 doi:10.1038/ncomms12619
920