1  **Selection constrains high rates of satellite DNA mutation in *Daphnia pulex***

2

3  Jullien M. Flynn[*], Ian Caldas[†], Melania E. Cristescu[‡], Andrew G. Clark[*†]

4

5  [*] Department of Molecular Biology and Genetics, Cornell University, Ithaca, USA

6  [†] Department of Biological Statistics and Computational Biology, Cornell University, Ithaca,

7  USA

8  [‡] Department of Biology, McGill University, Montreal, Canada

9

10

11  Raw data for this manuscript are available at NCBI Sequence Read Archive under BioProjectID:

12  PRJNA341529 BioSamples SAMN05725816 and SAMN05725817.

13    **RUNNING TITLE: High mutation rates in satellite DNA**

14    **KEY WORDS: tandem repeats, mutation accumulation (MA), stabilizing selection, satellite**

15    **evolution**

16

17    **Corresponding author:**

18    Jullien Flynn

19    Biotechnology-Rm 223

20    526 Campus Road, Ithaca, NY

21    14850

22    jmf422@cornell.edu

23

24   **Abstract**

25   A long-standing evolutionary puzzle is that all eukaryotic genomes contain large amounts of

26   tandemly-repeated satellite DNA whose composition varies greatly among even closely related

27   species. To elucidate the evolutionary forces governing satellite dynamics, quantification of the

28   rates and patterns of mutations in satellite DNA copy number and tests of its selective neutrality

29   are necessary. Here we used whole-genome sequences of 28 mutation accumulation (MA) lines

30   of *Daphnia pulex* in addition to six isolates from a non-MA population originating from the same

31   progenitor to both estimate mutation rates of abundances of satellite sequences and evaluate the

32   selective regime acting upon them. We found that mutation rates of individual satellite sequence

33   "kmers" were both high and highly variable, ranging from additions/deletions of 0.29 – 105

34   copies per generation (reflecting changes of 0.12 - 0.80 percent per generation). Our results also

35   provide evidence that new kmer sequences are often formed from existing ones. The non-MA

36   population isolates showed a signal of either purifying or stabilizing selection, with 33 % lower

37   variation in kmer abundance on average than the MA lines, although the level of selective

38   constraint was not evenly distributed across all kmers. The changes between many pairs of kmers

39   were correlated, and the pattern of correlations was significantly different between the MA lines

40   and the non-MA population. Our study demonstrates that kmer sequences can experience

41   extremely rapid evolution in abundance, which can lead to high levels of divergence in genome-

42   wide satellite DNA composition between closely related species.

43

44   **INTRODUCTION**

45      Up to half of the genome of higher eukaryotes is composed of tandem arrays of simple

46   repetitive motifs that can span megabases, called satellite DNA (Platero et al. 1998, Padeken et

3

47    al. 2015). Satellite DNA had initially been thought to carry no useful function, and because it

48    posed a replication burden, it became known as "junk DNA" (Orgel and Crick 1980). It also has

49    the potential to be harmful because it can cause deleterious genomic rearrangements facilitated

50    by recombination between similar motifs on different chromosomes (Bzymek and Lovett 2001).

51    Thus, evolutionary biologists have long wondered why repeated sequences accumulate and are

52    maintained at such abundance in the genome. Satellite DNA is typically found in

53    heterochromatic regions where expression is silenced and recombination is low, in Y

54    chromosomes, and in centromeric and telomeric regions (Charlesworth and Charlesworth 2000;

55    Henikoff et al. 2001). One hypothesis is that being situated in heterochromatin, where

56    recombination is suppressed, could both minimize the cost of repeated DNA and reduce the

57    efficacy of selection against it, allowing its accumulation (Charlesworth et al. 1994). Moreover,

58    the observation that, genome-wide, satellite repeats undergo rapid turnover, with nearly complete

59    replacement of these repeat sequences even between closely related species poses challenging

60    questions about the evolutionary forces that shape their evolution (Subirana et al. 2015; Wei et

61    al. 2014). The study of satellite DNA evolution on the timescales of diverged populations and

62    species have revealed interesting patterns, but understanding the forces that generated these

63    patterns requires direct observation on a shorter timescale.

64          Early evolutionary models did not consider selection to be influencing satellite sequence

65    evolution, except perhaps selection on overall genome size (MacGregor and Sessions 1986;

66    Charlesworth et al. 1994; Stephan and Cho 1994). Stephan and Cho (1994) proposed a model

67    where tandemly-repeated sequences could be generated and expand from mutation,

68    recombination, and replication slippage, combined with selection to maintain the genome size.

69    Another model, the "library hypothesis," attempts to explain the differences in satellite

70  composition by posing that the common ancestor of related species contained a library of many

71  satellites, which are then differentially amplified within each lineage as they diverge (Fry and

72  Salser 1977). However, the rate at which satellite sequences expand and contract is difficult to

73  quantify. Most models pose evolutionary neutrality of satellite repeats, but we now know that

74  satellite sequences also carry vital cellular functions. Since much of satellite DNA is located in

75  centromeric regions, there is evidence that selection acts on the sequence motifs themselves, as

76  they serve to bind centromere and histone proteins for centromere maintenance and chromosome

77  separation (Henikoff et al. 2001; Malik 2009). These functions have implications for genome

78  stability, and changes in satellite DNA sequences have been shown in some cases to be drivers of

79  speciation by inducing chromosome rearrangements leading to karyotype divergence (Paco et al.

80  2015). Moreover, some satellite sequences are transcribed and may be involved in the regulation

81  of heterochromatin formation (Palomeque and Lorite 2008; Plohl et al. 2008). Even more

82  intriguing is the observation that perturbation of satellite content can alter gene expression

83  genome-wide (Lemos et al. 2008). At present, it is not known whether these examples of

84  functional importance of satellite DNA are the exceptions or the rule and to what extent selection

85  governs the evolution of satellite DNA.

86      Here, we refer to the units or words of tandemly repeated satellite sequences as "kmers".

87  Several potential and non-mutually exclusive selection regimes could be operating on kmer

88  arrays: negligible selection where satellite DNA content is primarily governed by mutation and

89  drift, stabilizing selection to maintain a particular "optimal" composition, negative selection to

90  purge satellite DNA, and positive selection for rapidly generating new kmer motifs. Studying

91  satellite DNA changes over an evolutionary timescale showed clear differences across

92  geographic subpopulations of a single species and almost complete turnover between species,

93    confirming that they evolve rapidly (Wei et al. 2014; Subirana et al. 2015). However, quantifying

94    the relative contributions of mutation, genetic drift, and selection is difficult since all these forces

95    are at play in influencing genetic variation in natural populations. Additionally, it has been a

96    challenge to quantify the genome-wide satellite composition since their repetitive nature makes

97    them problematic to sequence and assemble (Hoskins et al. 2007). Early studies describing

98    satellite composition in different species have focused on single satellites or a small family of

99    satellite sequences (Lohe and Brutlag 1987). Obtaining a genome-wide view of the rate of

100    mutation in satellite sequences and how selection shapes their evolution would facilitate the

101    understanding the longstanding puzzles of satellite DNA evolution.

102        Mutation accumulation (MA) experiments reduce selection to a minimum by enforcing

103    bottlenecks every generation and reducing the effective populations size so that both neutral and

104    deleterious mutations can accumulate almost neutrally (Simmons and Crowe 1977). An MA

105    experiment combined with observations from a population with the same genetic background

106    where selection is not removed allows the study of the effects of mutation alone versus the

107    effects of mutation combined with selection (Flynn et al. 2017). *Daphnia pulex* is an ideal

108    organism in which to study satellite DNA mutation, firstly because MA studies can be conducted

109    effectively under asexual reproduction and single-progeny bottlenecks. Short microsatellite

110    arrays have been studied in *D. pulex* (Seyfert et al. 2008; Sung et al. 2010), but the sequence

111    motifs and abundances of satellite DNA have not been characterized, although 25% of the

112    genome is estimated to be heterochromatic (Colbourne et al. 2011). *D. pulex* has a potentially

113    dynamic genome, with evidence of high rates of deletion and duplication (Keith et al. 2016). The

114    lineage we use reproduces exclusively asexually via apomixis, meaning changes in satellite

115    content represent spontaneous mutations in the germline without the opportunity for meiotic

116    drive to play a role (Malik 2009; Wei et al. 2017).

117         In this study, we compare the kmer composition of MA lines (without selection) to

118    individuals raised in a competitive non-MA population (with selection) of the identical *D. pulex*

119    genetic background. We were interested in exploring the mutational dynamics of satellite DNA

120    that can often not be reliably mapped to the reference genome. We quantify the kmer

121    composition using a mapping-independent method to capture all tandem kmers up to 20 bp

122    detectable from short-read data, without introducing biases from the reduced representation of

123    repetitive sequences in the reference genome. Our approach allows us to (1) estimate the kmer

124    mutation rates, including expansions, contractions, complete loss of kmers, and generation of

125    new kmers; and (2) evaluate the type of selection (if any) acting on the satellite DNA sequences.

126    Given the observed natural variation in satellite DNA content among populations and species, we

127    expect expansions and contractions to occur at high rates in the MA lines. If satellite DNA is

128    under selective constraint, we expect there to be less variation in kmer abundance in the

129    population evolving under selection than the neutrally-evolving MA lines.

130

131    **MATERIALS AND METHODS**

132

133    *Daphnia line setup and DNA sequencing*

134    The current study includes a total of 34 genomes: 24 MA line genomes and 6 population isolate

135    genomes that were sequenced in Flynn et al. (2017), and four additional previously unpublished

136    MA genomes. The asexually reproducing MA lines and a non-MA population were initiated

137    from a single progenitor (see Flynn et al. 2017 for a detailed description of the MA experiment).

138    MA lines were propagated for an average of 82 generations before whole-genome sequencing

139    with single-individual bottlenecks between generations prior to DNA isolation and whole-

140    genome sequencing. In contrast, the non-MA population was maintained without inducing

141    bottlenecks for 46 months in a 15 L tank, before six individuals were isolated for sequencing.

142    The approximate census population size was estimated to be 100-250. Overlapping generations

143    occurred so the exact number of generations that the population isolates progressed could not be

144    recorded. Thus, we used a life history experiment and estimated the slowest and fastest moving

145    lineages to calculate that the population underwent at least 62 generations (Supplementary

146    Material File S2). Prior to sequencing, *Daphnia* individuals were subjected to a brief antibiotic

147    treatment and fed with sterile Sephadex beads to reduce contaminants before sequencing (Fields

148    et al. 2015).

149        Sequencing libraries were prepared with Illumina Nextera procedures in two batches and

150    all genomes were sequenced to approximately the same coverage (10x). To ensure

151    reproducibility between multiple library preparations, we performed technical replication on two

152    MA lines, such that two independent library preparations and Illumina sequencing runs were

153    done for the MA lines C01 and C35 (in a completely separate, third sequencing batch). We

154    analyzed independently all technical replicates to ensure that variation produced by different

155    library preparations is smaller than the variation due to biological expansion/contraction

156    mutations. Libraries were sequenced on an Illumina HiSeq 2000 instrument at Genome Quebec

157    of McGill University with 100 bp paired-end reads.

158

159    *Satellite quantification*

160

8

161 To remove redundant sequences, adapters were trimmed and overlapping reads were merged

162 with SeqPrep (https://github.com/jstjohn/SeqPrep). To identify and quantify satellites, the

163 resulting unmapped read files were used as input for the program k-Seek (Wei et al. 2014,

164 https://github.com/weikevinhc/k-seek). The current version of k-Seek detects words or "kmers"

165 of length 1-20 bp that are repeated tandemly to cover at least 50 bp within the same read. The

166 program allows for one single nucleotide mismatch per repetition of each kmer (see Wei et al.

167 2014 for details). All offsets and reverse complements of each kmer are compiled and the total of

168 all individual unit copies are summed across all reads. Kmer sequences are presented as the

169 strand/offset that is alphabetically ordered (i.e. A's will be at the start of the sequences if

170 possible). This method of tandem kmer detection and quantification has been found to be

171 reproducible across different library preparations (Wei et al. 2014).

172 To obtain a quantitative comparison between samples we normalized the kmer counts by

173 both individual library sequencing depth and GC content, as PCR-based library preparations are

174 known to show a bias in the GC content of the fragments amplified and sequenced (Benjamini

175 and Speed 2012). This is especially problematic for satellite analysis since, by nature, many of

176 the repeated sequences may have extreme GC contents. First, reads were mapped against the

177 *Daphnia pulex* reference genome (version 1, Colbourne et al. 2007) using BWA (Li and Durbin

178 2009) v0.7.10 with default settings. Output BAM files spanning the whole genome were given as

179 input to a custom shell script (https://github.com/jmf422/Daphnia-MA-

180 lines/tree/master/GC_correction) to calculate correction factors following Benjamini and Speed

181 (2012). The correction factors produced for high and low GC contents were extreme and

182 variable, which was likely due to low read counts giving unstable estimates as well as mapping

183 biases to the reference genome (Flynn et al. 2017). Thus, to smooth the correction, we used a

9

184  Python script to bin values of GC content together, employing wider bins for GC content < 0.25

185  and > 0.60 (https://github.com/jmf422/Daphnia-MA-lines/tree/master/GC_correction). The

186  correction factor was then applied to kmer counts according to which GC bin their content

187  belonged.

188

189  *Mutation and interspersion metrics*

190  We define mutations in kmer repeats as any change in the number of copies of each specific

191  kmer. The rate of mutation is then the observed change in the number of each kmer per

192  generation. Each kmer may be in tandem arrays found in one or in several locations across the

193  genome, and our method sums the counts of each kmer across all genomic regions. A mutation

194  rate for a given kmer reflects the sum of changes in copy number across these potential multiple

195  loci. To calculate the mutation rates, we used the mean abundance of the kmer in the non-MA

196  population isolates as a proxy for the ancestral abundance using the following equation:

197

198  $$u_{i,m} = \frac{m_i - \overline{P_i}}{G_m}$$

199
200  Where $u_{i,m}$ is the mutation rate of kmer $i$ in MA line $m$, $m_i$ is the abundance of kmer $i$ in MA line

201  $m$, $\overline{P_i}$ is the mean abundance in the population for kmer $i$, and $G_m$ is the number of generations

202  propagated for MA line $m$. $u$ could be negative (for contractions) or positive (for expansions).

203  We used the same equation to calculate the absolute rate of change of each kmer, except we took

204  the absolute value of the numerator.

205       As mutation rates could be correlated with copy number, we also calculated the absolute

206  mutation rate normalized by initial copy number using the following equation:

207

208 $\qquad u_{i,m} \; = \; \dfrac{|\, m_i - \overline{P_i}\, |}{G_m \times \overline{P_i}}$

209

210     We used the same equations for calculating the "realized" kmer mutation rates (mutations that

211     made it through selection in the population isolates), and we used the conservative estimate of 62

212     for the number of generations.

213         We also searched for new kmers generated *de novo* in the MA experiment. In order to be

214     considered as a new kmer, putative new kmers had to have at least 3 copies in at least one MA

215     line and be completely absent from all other lines. New shared kmers had to have at least 3

216     copies in one MA line and at least 2 copies in a second line. Since k-Seek has a detection limit

217     for kmer arrays that are at least 50 bp long, we checked if the putative new kmers were present in

218     shorter arrays (~25 bp) in the other MA lines to determine if the kmer sequence was generated

219     truly de novo or if it expanded from an already present repeated motif. Similarly, we searched for

220     kmers that had been lost in individual and pairs of MA lines throughout the course of the

221     experiment. Our criteria for identifying lost kmers were that they had to have 0 copies in the

222     affected line or pair of lines, and at least 2 copies in all other lines. We checked if putative lost

223     kmers were completely undetectable or if they were still present in short arrays (~25 bp), below

224     the detection threshold of k-Seek.

225         Since we have paired-end reads, we have the potential to detect the level to which kmers

226     are present on both reads of the pair (on the same genomic fragment). To quantify this

227     interspersion level, we used the metric *I*, as in Wei et al (2014):

228     $I \; = \; n_{ij} / \sqrt{n_i n_j}$

229     where $n_{ij}$ is the number of read pairs containing kmer *i* and kmer *j*, and $n_i$ and $n_j$ are the number

230     of read pairs containing kmer *i* and kmer *j* respectively.

231

11

232 *Data availability*

233 Raw sequence files have been deposited in the NCBI Sequence Read Archive BioProject ID:

234 PRJNA341529 BioSamples SAMN05725816 and SAMN05725817.

235 All analyses were performed using RStudio (V 0.99.903), Python V2 and Perl V5 scripts.

236 Html files showing the code and output to all the analyses are available in the Supplementary

237 Material, and all scripts as well as the required processed input files are available at

238 https://github.com/jmf422/Daphnia-MA-lines. File S1 contains descriptions of all supplementary

239 files, as well as the supplementary table and figure legends.

240

241 **RESULTS**

242

243 **Library and GC correction**

244 Over 713 million unmapped reads across all samples were scanned for tandemly-repeated

245 sequences with k-Seek. Of these, over 5.85 million reads (0.82 %) were found to be composed of

246 kmers of unit length 1-20 bp and encompassed at least 50 bp of the read. Although the second

247 library batch had overall more reads, a similar proportion of the total reads contained kmers in

248 both libraries (Supplementary Figure S1). Unless otherwise noted, the abundance of each kmer is

249 presented in copies per 1x read depth after GC normalization.

250 We found that our qualitative conclusions were robust whether or not we applied a

251 correction factor to normalize the kmer counts based on each kmer's GC content. Our results

252 were also robust across a wide range of parameters used for the GC correction. As library

253 preparations involving PCR have been shown to result in a biased representation of fragments

12

254    based on GC content (Benjamini and Speed 2012), we present results after GC correction as

255    described in the Materials and Methods.

256         Since the libraries were prepared in two separate batches, we were concerned that this

257    could confound our results. A Principal Components Analysis (PCA) was able to separate

258    samples based on library batch along PC2, but PC1 appeared to separate samples based on

259    biological mutation patterns (Supplementary Figure S2). There were no consistent differences

260    across all kmers or GC contents of kmers that could be consistently corrected for (Supplementary

261    Material File S4 ). In order to ensure that library batch would not be a confounding factor, we

262    performed a technical replicate of two MA lines and showed that the abundances of each kmer

263    between the library prep batches were highly similar (Supplementary Material File S4).

264    Although there was minor variation between technical replicates in some kmers, their overall

265    mutational patterns were highly similar, shown by their high clustering on PC1 (Figure S2).

266    Additionally, we tested for differences between the mean abundances of the kmers between the

267    MA lines prepared in the first batch and the MA lines prepared in the second batch. The mean

268    abundances of most of the kmers (28 out of 39) in the second-batch MA lines were within the

269    99% confidence intervals from 1000 subsampling replicates of four first-batch MA lines. 30 of

270    the kmers had mean abundances in the second batch within the range of observed abundances of

271    the first-batch MA lines (Supplementary Material File S4). For these reasons, we did not

272    consider library preparation batch to be a significant confounding factor.

273

274    **Description of satellite DNA content in *Daphnia pulex***

275

276    We first sought to describe and quantify the genome-wide short repeat content in *Daphnia pulex*

13

277 and make inferences about kmer origins. There were 162 kmers that had at least 2 copies per 1x

278 coverage in all our *D. pulex* lines. There were 39 kmers that had an average abundance of at least

279 100 copies, and 12 that had an average abundance of at least 1000 copies after normalization

280 (Table S1). We chose to focus most of our analysis on the 39 kmers that had at least 100 copies

281 after normalization (Figure 1). The most abundant kmer, the poly-A repeat, was present at an

282 average of 79,528 copies in each genome. The second most abundant kmer, present at an average

283 of 62,281 copies, was AACCT. This 5-mer is known to be the ancestral telomere repeat in

284 Arthropods (Sahara et al. 2009) and was previously found in *Daphnia* (Colbourne et al. 2011).

285 Runs of poly-C were also found to be in the "top 39 kmers", having an average of 6379 copies.

286 Two of the four possible 2-mers were included, AG and AC, having an average of 408 and 244

287 copies, respectively. The most abundant kmer sizes were 5-mers (there were six 5-mers of the

288 most abundant 39 kmers), 10-mers (7/39), and 20-mers (7/39). No 15-mers were abundant (>100

289 copies) in the genome, but there were 13 among all 162 kmers. Overall, 20-mers were the most

290 represented with 50 kmers of the 162. Adding up the total short kmer content per genome in base

291 pairs, we found the median to be 1.20 Mb per 1x coverage, which represents 0.6 % of the

292 estimated 200 Mb genome.

293 We also compared the satellite composition in *D. pulex* to that in *Drosophila*

294 *melanogaster,* the only other arthropod that has had its genome-wide satellite content

295 characterized to the same level of detail (Wei et al. 2014). This comparison might indicate the

296 extent of satellite diversity across arthropods and functional conservation of some satellites.

297 Comparing the most abundant kmers (at least average 100 copies normalized) between these two

298 species, we found that 10 short kmers (mostly derivations of AG repeats) were present in both

299 species: A, C, AC, AG, AAC, AAG, ACT, AAAAG, AAGAG, and AAGGAG. We found

14

300    several short (3-5 bp) motifs to be highly represented in the most abundant kmers, and among the

301    162 kmers (Table 1). Most of these motifs were completely absent or rare within the top 108 *D.*

302    *melanogaster* kmer sequences, except for AG-type motifs (Table 1).

303        In order to understand the origins and dynamics of kmer sequences in *Daphnia*, we

304    inspected the kmers for sequence similarities. We found that many kmers belonged to families of

305    related sequences separated by a small number of potential mutational steps. Figure 2 shows a

306    network diagram highlighting the potential mutational relationships between 29 of the 39 of the

307    top kmers. The most striking example of a family of related kmers is the 9 kmers of length 10-20

308    bp shown in the top left corner of Figure 2. Ten kmers were not included in the network, either

309    because they were mono-, di-, or tri- nucleotides whose sequence similarities were not

310    meaningful (6 kmers: A, C, AC, AG, AGG, ACT), or because they did not have a clear

311    mutational relationship with any of the other top kmers (4 kmers: AGCCTG, ACAGC, ATCC,

312    AACGGTACGG).

313

314    **Extremely high rate of mutation in satellite repeats**

315    Because we could not call single nucleotide changes reliably without mapping to the reference

316    genome, we were not able to estimate mutation rates for single base changes in the kmers, and

317    restrict our attention to changes in copy number. To estimate the copy number mutation rates, we

318    used the mean abundance of each kmer in the population isolates as a proxy for the ancestral

319    abundance of the kmer in the MA progenitor. This should be a reasonable estimation assuming

320    that individuals in the population roughly maintained the ancestral kmer content (see below). In

321    fact, all kmers were similar in abundance between the MA lines and the population isolates: all

322    kmer abundances in the population were within 2 standard deviations of the abundance in the

323  MA lines, and 69% were within 1 standard deviation. This is consistent with the MA lines both

324  gaining and losing kmer copies. The absolute mutation rates (summing expansions and

325  contractions) per kmer sequence ranged from 0.29 to 105 copies/generation, with a median of

326  1.26 copies/generation. This is not including the telomere repeat, which had the highest rate of

327  change at 199 copies/generation. It is likely that the changes affecting the telomere repeats are

328  from a different, non-mutational process (i.e. caused by a mechanism involving telomerase

329  which is not transgenerationally inherited). The suppression of meiosis in these *Daphnia* lines

330  may have resulted in relaxed selection on telomere maintenance, but we have no direct evidence

331  in support of this.

332      We found that the mutation rates were correlated with the kmer's initial copy number

333  abundance, proxied by the mean abundance of the population isolates, and that the best fit of this

334  relationship is linear ($r^2 = 0.88$, $p < 2.2 \times 10^{-16}$). Thus, to compare mutation rates of individual

335  kmers, we normalized the mutation rates by the initial unit copy number to give mutation rates in

336  copies per generation per original copy. Even after this normalization, we still found mutation

337  rates to be exceptionally high – on the order of $10^{-3}$ per copy per generation – compared to

338  previous estimates of other types of mutation such as microsatellite mutations. Normalized

339  mutation rates ranged from $1.23 \times 10^{-3}$ to $7.97 \times 10^{-3}$, with a mean of $2.74 \times 10^{-3}$ copies per

340  generation per original copy. Expansions were more frequent than contractions, with more MA

341  lines experiencing expansions in most kmers, and the overall rate of change was in a positive

342  direction for 124/162 kmers (76%) (Figure 3a and 3b). Expansion and contraction rates were not

343  correlated with either GC content or kmer size (Supplementary Material Figure S3).

344      A multi-generational MA experiment provides the potential opportunity to observe

345  completely new kmer sequences arising. We searched for kmers that were generated *de novo*

16

346    during the experiment, in single lines and shared across pairs of MA lines. We found five kmers

347    gained in single MA lines and three gained in three independent pairs of MA lines (Table 2).

348    These MA lines did not share single nucleotide mutations within uniquely mapping regions of

349    the genome (Flynn et al. 2017). Five of the eight new kmers seemed to originate from a single

350    base substitution of an existing kmer present in high abundance (Table 2). Two of the other new

351    kmer sequences contained motifs that were abundant in existing kmers. Only 3 out of 8 putative

352    new kmers were present in shorter arrays (~25 bp, half of k-Seek's requirement) in the ancestral

353    lineage (Table 2). We also searched for kmers that were lost from individual and pairs of MA

354    lines. There were seven kmers lost uniquely and two lost in a pair of lines (Table 3). Three of the

355    seven lost kmers were lost from MA line C40. C40 is the MA line that was found to have ~100

356    kb of homozygous deletions of non-repetitive sequence mapped to the reference genome on one

357    chromosome (Flynn et al. 2017), so it is possible that these lost kmers were present exclusively

358    in the deleted regions. Of the lost kmers, 5 out of 7 were still detectable in shorter arrays (~25

359    bp) in the affected MA line (Table 3).

360

361    **MA lines versus population with selection**

362

363    Comparing the changes in kmer composition between the MA lines and the population isolates,

364    which experienced selection, reveals the potential influences of selection on satellite DNA. The

365    population provides a valid comparison to the MA lines because it was previously shown to

366    show signs of purifying selection in both single nucleotide mutations (Flynn et al. 2017) and

367    copy number variants (CNVs) (Chain et al. 2017, in prep). The population did not experience a

368    recent bottleneck and in fact the most recent common ancestor of the isolates analyzed here was

17

369    the progenitor of the experiment (Flynn et al. 2017). We also made a conservative estimate of the

370    number of generations the population underwent in order to compare mutation rates

371    (Supplementary File S2). First, we compared the total genome-wide amount of tandem repeats

372    by multiplying the length of the kmer with its number of copies in each genome. The MA lines

373    diverged by a factor of 1.67 in overall kmer abundance, from 1.00 Mb to 1.67 Mb total, with a

374    median of 1.23 Mb. The population isolates had both a lower genomic amount of kmers and a

375    narrower range among isolates, ranging from 0.97 to 1.22 Mb, with a median of 1.05 Mb (Figure

376    4a). A Levene's test did not detect a significant difference in the variances between these two

377    groups (p = 0.12). However, the mean total abundance (but not the variance) of kmers in the

378    population was below the 5% confidence interval produced from 1000 replicates of 6 randomly-

379    sampled MA lines (Supplementary Material File S3).

380        The MA lines underwent expansions across many kmers, with some MA lines

381    experiencing a considerable increase in satellite DNA content in their genome over the course of

382    the MA experiment compared to the non-MA population isolates (Figure 4b). Most notably was

383    MA line C20, which experienced expansions in 34 of the 39 top kmers, with an overall increase

384    in 580.2 kb in satellite content over 81 generations. On the other hand, population isolates

385    deviated less from an overall balance of expansions and contractions (Figure 4b).

386        To perform a contrast in the variation of individual kmers between the MA lines and the

387    population, we calculated the coefficient of variation (CV, the standard deviation divided by the

388    mean) for each kmer for both the MA lines and the population. We found that the CV was higher

389    in the MA lines than the population for 37 of the top 39 kmers (Figure 4c, $P < 2.84 \times 10^{-9}$, sign

390    test). We considered the possibility that the CV of the MA lines was artificially inflated if the

391    mean abundances of some kmers in the second library preparation MA lines were different that

18

392   of the first library preparation MA lines. We also considered that the CV may be greater in the

393   MA lines partially because the MA lines may have undergone more generations than the

394   population. In order to ensure these factors were not founding our results, we also calculated the

395   CV of the MA lines comparing only the ones that were prepared in the same library and

396   normalized the CV of the MA lines by the difference in the generation number between the MA

397   lines and the population. After this conservative calculation, the CV was higher than the

398   population isolates across 25 of the 39 top kmers (Supplementary Material File S3).

399   Next, we compared the realized mutation rate – the mutations that made it through the "filter" of

400   selection in the non-MA population – to the mutation rate in the MA lines. We used the

401   conservative calculation of 62 generations propagated in the population (Supplementary Material

402   File S2). The realized mutation rates in the population were lower than the MA lines for 32 of

403   the top 38 kmers not including the telomere repeat (Figure 4d, $P < 2.43 \times 10^{-5}$, sign test). Since

404   the sample size of the non-MA population isolates is lower than the sample size of the MA lines

405   (6 versus 28), we performed 1000 subsampling replicates for each kmer by sampling 6 random

406   MA lines and calculating the mean mutation rate. We found that 22 out of 38 kmers had a lower

407   realized mutation rate in the population than the 5% quantile of the 1000 MA line subsample

408   replicates (Supplementary Material File S3). Moreover, some kmer sequences seemed to be more

409   constrained than others. We roughly quantified constraint by the difference between the mutation

410   rate of the MA lines to the realized mutation rate of the population for each kmer. Five kmers

411   experienced high levels of constraint by this measure, with the realized mutation rate being at

412   least 75% lower in the population than the MA lines: AGG, C, AGCCTG,

413   AAGCCAGTGCAGC, and AATCTGGAATGGAATGG. All of these highly constrained kmers

414   had statistically significant lower realized mutation rates by the subsampling test above. On the

415    other hand, 9 kmers seem to have experienced little constraint in their expansions and

416    contractions, with their realized mutation rates in the population being either slightly higher in

417    the population or very close to that of the MA lines (Figure 4d, Supplementary Material File S3).

418          Wei et al. (2014) noted strong patterns of correlation of kmer abundances across different

419    populations of *Drosophila melanogaster*, suggesting some sort of evolutionary non-

420    independence. These correlation patterns may suggest constraints on the mutational process, or

421    they may be driven by selection. To investigate correlations in the expansions/contractions

422    between kmers, we computed a correlation matrix (using the Pearson correlation) between each

423    kmer for the MA lines and the population. This would indicate correlations between changes in

424    kmer content inherent to mutation as well as correlations caused by similar selective regimes

425    among kmers. We found that some (mainly positive) correlations existed in the MA lines, but the

426    correlations between kmer changes in the population were more common and stronger for both

427    positive and negative correlations (Figure 5a and 5b). A Mantel test indicated that these

428    correlation matrices were significantly different ($P < 0.001$), indicating a role of selection in

429    interactions between satellite sequences.

430          For both the MA lines and the population, strong correlations between kmers that had GC

431    contents between 0.6-0.7 were apparent. Although it is possible that the correlations reflect

432    technical bias, we think this is unlikely because these strong correlations were between the nine

433    kmers closely related in sequence (Figure 2). These nine 10-20 bp kmers have a GC content

434    between 0.6 and 0.7, and their deviations were in fact strongly positively correlated (Figure 5c).

435    Aggregations of kmers that are related in sequence might easily explain some of the patterns of

436    correlation among kmer sequences that we identified. To follow up on this, we analyzed the

437    paired-end reads to measure the level of interspersion between kmer sequences. All kmers were

20

438    highly interspersed with themselves, indicating the that kmer arrays we are studying span at least

439    the length of the sequenced fragments (~250 bp). Several groups of kmers were interspersed with

440    each other (Supplementary Material, File S5). The telomere repeat, unsurprisingly, was not

441    interspersed with any other of the top 39 kmers. The group of nine related kmers mentioned

442    above were the most mutually interspersed as a group. Eight of the nine were interspersed with at

443    least one other in the group, and the more abundant kmers of the group were interspersed with up

444    to six other kmers within the group.

445

446    **DISCUSSION**

447

448    To our knowledge this is the first study that assays the mutation rates of genome-wide tandemly

449    repeated satellite DNA using MA lines. The inclusion of data from a population that experienced

450    selection also enables us to begin to detect whether selective forces act on satellite DNA. We

451    found that mutation rates in satellite DNA are high, but they appear to be constrained by

452    selection. We assumed the population isolates maintained the ancestral satellite composition

453    when calculating mutation rates, which was reasonable given that the means of the population

454    and MA lines were similar for most satellites, and the population isolates showed narrow

455    changes in kmer abundance. The high variance in the abundance of kmers across the MA lines

456    resulted in high estimated mutation rates. We are confident that these results are biological and

457    not due to technical biases because we normalized kmer counts by the GC-normalized read

458    coverage (see Materials and Methods), and after normalization, there was no longer a correlation

459    between sequencing depth and kmer abundance (Supplementary Material Figure S4). All MA

460    lines and population isolates contained a similar proportion of reads with kmers to total reads

461    (Supplementary Material Figure S1). There was also a high level of concordance in the kmer

462    content between technical replicates. Another line of support that these data are reliable for this

463    type of study is that the analysis of copy number variation of mapped reads revealed mutation

464    rate estimates close to those of previous studies (Chain et al. 2017 in prep). Going forward, using

465    PCR-free library preparation methods can reduce batch-based biases in sequencing results (Wei

466    et al. 2017, in prep).

467

468    **High mutation rates can explain the rapid turnover of satellites**

469

470    The rates of mutation in satellite sequences in our *D. pulex* MA lines were extremely high,

471    ranging from 0.29 to 105 copy changes per generation for a given kmer sequence. Our mutation

472    rates include the sum of changes across all loci of a particular kmer, thus can include both

473    loss/gain of entire repeat arrays (i.e. from recombination), and changes in the number of repeats

474    within individual loci (i.e. from replication slippage). This makes our study distinct from

475    previous studies of microsatellite mutation rates that have focused on single loci (Seyfert et al.

476    2008). Here, mutation rates estimated on a per-copy basis for satellite repeats we examined were

477    on the order of $10^{-3}$ copies/generation/copy. Seyfert et al. (2008) estimated the per locus mutation

478    rate of *D. pulex* dinucleotide microsatellites 13-47 repeats long to be on the order of $10^{-5}$ to $10^{-4}$

479    copies/locus/generation. This indicates not only that many satellite sequence arrays are present as

480    multiple loci and long arrays, but that they mutate at much higher rates than microsatellites

481    typically studied. One contributing factor could be that the short microsatellite arrays are more

482    constrained since they are often found in introns or even coding regions (Ellegren 2004). Past

483    studies of microsatellites have found that rates of mutation are positively correlated with the

22

484     number of repeats, but limited data points prevented the determination of a linear relationship

485     (Wierdl et al. 1997; Brinkmann et al. 1998). We had mutation rate estimates for 39 kmer

486     sequences ranging in abundance from 100 - 50,000 copies and found a significant positive linear

487     correlation between copy number and mutation rate. This suggests a simple relationship between

488     the number of copy units present and the potential for copy units to be gained or lost, consistent

489     with replication slippage being the dominant mechanism of mutation (Schlötterer 2000). We also

490     found that expansions are more common than contractions, indicating that satellite DNA in

491     *Daphnia* has the intrinsic tendency to expand, which is in agreement with the consensus finding

492     in microsatellites across taxa (reviewed in Ellegren 2004).

493           This high mutation rate demonstrates that there is a high potential for rapid turnover in

494     satellite DNA sequences between species, and even substantial differences in abundances among

495     diverged populations of the same species. Wei et al. (2014) used k-Seek to characterize the

496     satellite sequences in eight different populations of *D. melanogaster,* distributed globally. They

497     found that the overall satellite content varied by a factor of 2.5, accounting for about a 4 Mb

498     difference. This initially surprising finding makes sense in light of our results that the MA lines

499     that were only 82 generations diverged differed in their overall satellite content by a factor of

500     1.67, accounting for about 0.67 Mb. Even the population isolates, which showed a reduction of

501     variation in satellite content compared to the MA lines, varied in their overall satellite content by

502     a factor of 1.27. This highlights the great potential for expansions and contractions in satellite

503     sequences, which can explain the high levels of differentiation in genome-wide satellite

504     compositions between populations and related species (Wei et al. 2014; Subirana et al. 2015;

505     Wei et al. 2017 in prep).

23

506    *Daphnia* and *Drosophila* are both arthropods, although they are estimated to have

507    diverged about 400 million years ago (Rehm et al. 2011). Given the high differentiation observed

508    between closely related species, the paucity of shared sequences and motifs between these taxa is

509    not surprising. It is worth noting, however, that the AGG and AAG motifs as well as a couple

510    AG-rich satellites were conserved between the two species. This conservation could be random

511    or it could indicate a conserved functional role for AG-rich satellites; for example, involvement

512    in binding conserved proteins. The GAGA factor is known to bind to AG-rich satellite sequences

513    in *Drosophila*, although it has been hypothesized to have evolved in the ancestor of Diptera and

514    Hymenoptera and is not thought to be present in *Daphnia* (Heger et al. 2013). However, it is

515    possible that other proteins could be playing a role in binding to AG-rich satellites in *Daphnia*

516    similarly to *Drosophila*. Another possibility is that AG-rich satellites result in a favourable DNA

517    curvature (Palomeque and Lorite 2008). *D. melanogaster* populations were found to be enriched

518    in kmers that were multiples of 5 bp long (i.e. 5, 10, 15, 20, -mers) (Wei et al. 2014).

519    Interestingly, we also found that *Daphnia* were enriched in kmers 5, 10, and 20 bp long.

520    Although there were no 15-mers at high abundance ($\geq$100 copies), 15-mers were enriched when

521    considering the dataset of 162 kmers. This pattern, conserved from *Drosophila* to *Daphnia*,

522    could indicate conservation in the periodicity of DNA wrapping around histones (Lohe and

523    Brutlag 1987).

524    Based on the new kmers generated *de novo* in the MA experiment and the relatedness of

525    existing kmers, we suggest that novel kmers typically result from a point mutation in an already

526    abundant kmer, followed by expansion of the new kmer. This is in agreement with the result

527    from Wei et al. (2017 in prep), who found new kmers originating in different *Drosophila* species

528    that also differed by only a single nucleotide from a pre-existing kmer. Additionally, satellite

529    kmers can be generated when they expand from a sequence motif that is already present in only a

530    few copies, as we suggest could have occurred for some of our new kmers (Table 2). Concordant

531    with these hypotheses, we found that two kmers were gained independently in two independent

532    MA lines, demonstrating the high potential for new satellite sequences to be generated from

533    existing ones or from an existing motif, and that parallel independent gains of the same satellite

534    are possible. We detected the newly-generated kmers at low abundance, but after many more

535    generations they would presumably have the potential to expand and achieve higher abundance

536    in the genome. Most of the lost kmers were not completely lost, but contracted to be part of a

537    short enough array to not be detected by k-Seek (Table 3). Some of the satellite losses we found

538    likely resulted from deletions of segments of chromosomes. Flynn et al. (2017) found that the

539    MA line C40 experiences homozygous deletions totalling ~100 Mb on chromosome 11. Here we

540    found that C40 lost 3 different satellite sequences, so we think it is likely that these satellites

541    were present on the deleted regions of chromosome 11 (although these tandem sequences were

542    not present in tandem in the genome assembly). In fact, C40 experienced a complicated

543    recombination event that resulted in these deletions and complete loss of heterozygosity across

544    the chromosome (Flynn et al. 2017). It is possible that the satellite DNA composition on this

545    chromosome causes it to be a hotspot for structural rearrangements. Nevertheless, if most of our

546    kmer losses resulted from deletions of large segments also containing functionally important

547    regions, these would likely be purged from natural populations and thus would not be detectable

548    in population data.

549

550    **Selection on satellite content**

551

25

552     In contrast to the MA lines, which accumulate mutations with minimal regard to their phenotypic

553     consequences, we found that most kmers experienced constraint in their evolution in the non-MA

554     population. The constraint was disproportionately distributed among kmer sequences, with some

555     being unconstrained in these conditions and others being highly constrained. Both the coefficient

556     of variation and mutation rates were lower for many kmers in the population. The population

557     reproduced asexually and it had a census population size of 100-250, indicating selection would

558     have to be substantially strong in order to produce the signal that we observed. Our results

559     suggest that satellite repeats have the intrinsic tendency to expand in the *Daphnia* genome, but

560     the magnitude of expansions are reduced under a selection regime. Although this study does not

561     investigate specific functional roles of satellites, we recognize the potential for some satellites to

562     have specific roles (e.g. Fry and Salser 1977; Palomeque and Lorite 2008; Plohl et al. 2008).

563     This may be true for AG-rich satellites, which are known to bind proteins and are conserved in

564     both *Drosophila* species and *D. pulex* (Raff et al. 1994). The selection we infer might stem from

565     conservation of a specific role for some satellites, and/or deleterious effects of changing the copy

566     number of satellites.

567         In calculating mutation rates, we assumed that the population maintained the ancestral

568     kmer content. Concordant with this assumption, there was little variation in the total summed

569     abundance of satellite sequences in the population. This would be consistent with the operation

570     of stabilizing selection on overall satellite content (implying that there is an "optimal" level).

571     However, we note that it is possible that negative selection (purging satellite DNA) could have

572     been acting in the population and the lower realized mutation rates in the population was a result

573     of mutation-selection balance (i.e. mutation pressure to increase satellite content, selection to

574     reduce it). Under stabilizing selection, the reduced realized kmer mutation rates in the

575    population could stem from selection on overall genome size or chromosome organization.

576    Previous studies have shown the importance of satellites for maintaining chromatin organization

577    and that genetic variation in satellite content can affect genome stability and chromosome

578    function (Kim et al. 2009, Aldrup-MacDonald et al. 2016). If the genome has evolved to

579    accommodate satellite arrays, altering the length of an array could alter the chromatin

580    organization and disrupt nearby functional regions, causing deleterious effects. This could be

581    especially important for *Daphnia,* which has a compact genome with little intergenic space

582    (Colbourne et al. 2011). If selection to remove satellite DNA was the dominant form of selection,

583    this would imply that the satellite arrays were mostly "junk" and the genomes would be in a

584    constant state of purging satellite DNA. Moreover, it seems that the level of selection varies

585    across kmers. We found five kmers to be under very little constraint, shown by the lack of

586    difference between their CV and mutation rates between the MA lines and population (Figure 4 c

587    and d).

588

589    **Correlations in kmer changes**

590

591    Across the MA lines, we found mainly positive correlations between the kmers in their

592    expansions and contractions. These could have arisen randomly or been driven by physical

593    location/linkage. We found that kmers that were closely related to each other were strongly

594    correlated in their changes in abundance, which was especially evident for the family of 10-20 bp

595    kmers with the GCCAG motif (Figure 2, top left). Concordant with the hypothesis that new

596    satellite sequences come from point mutations in existing ones, we found that these closely

597    related kmers were highly interspersed within each other. Since these positive correlations were

598   present in both the MA lines and the population, it is likely to be a consequence of the mutational

599   process. Segmental duplications and deletions, which are common in *Daphnia* (Keith et al.

600   2016), could explain the observed positive correlations. When only the population isolates were

601   considered in the correlations, we identified more correlations of stronger effect, both positive

602   and negative. The positive correlations could arise from the same selective regime acting on

603   pairs of kmers, i.e. if expansions of both kmers of the pair in question are both selected against.

604   On the other hand, the negative correlations may arise if different satellite sequences are in

605   conflict (Wei et al. 2014). Specifically, there could be selection on a chromosome to maintain its

606   organization such that if one kmer happens to expand, another will have to contract. Kmers

607   localized on the same chromosomal region may be the ones negatively correlated with each

608   other, however, we were not able to test this thoroughly with our data because of the un-

609   mappable nature of repeats.

610

611   **Conclusions**

612   We were able to gain a picture of genome-wide satellite DNA expansions and contractions, and

613   provide evidence that satellite DNA is not always neutrally evolving but can experience strong

614   selection. We also show evidence that new satellite sequences are often generated from existing

615   ones, and there is a complex interaction structure between individual satellites that differs if

616   selection is at play. Future studies using different sequencing technologies, especially with

617   longer reads that can capture longer satellite sequences and span the boundaries of unique and

618   repeated sequences, will provide further insights into satellite DNA evolution.

619

620   **ACKNOWLEDGEMENTS**

28

629     **TABLES**
630

Table 1. Number of kmers that contain certain sequence motifs. The top 39 kmers are characterized by having an average of at least 100 copies per 1x coverage, and the top 162 kmers have at least 2 copies per 1x coverage in all samples. The *D. melanogaster* kmers are the 108 kmers that have at least 100 copies normalized (personal communication, K. Wei).

631

| Motif sequence | In top 39 kmers | In top 162 kmers | In *D. melanogaster* kmers |
|---|---|---|---|
| AAAA | 4 | 13 | 12 |
| AAC | 7 | 32 | 16 |
| AAG | 11 | 49 | 45 |
| ACGC | 8 | 52 | 0 |
| AGC | 12 | 69 | 4 |
| AGG | 11 | 41 | 14 |
| AGGAG | 3 | 15 | 5 |
| GCCAG | 7 | 48 | 0 |
| TAGG | 4 | 10 | 0 |
| TCCAG | 3 | 14 | 0 |

632

Table 2. New kmers generated *de novo* during the MA experiment. Many of them are potentially generated from single nucleotide mutations or rearrangements of existing kmers. It is indicated if the putative new kmer was present in shorter arrays (~25 bp) in the ancestral lineage.

633

| New kmer | Line(s) | Normalized abundance | Present in shorter arrays | Similar kmer, or common motifs (copy | Abundance of similar kmer (copy |
|---|---|---|---|---|---|

29

| | | | | number) | number) |
|---|---|---|---|---|---|
| AAATG | C18 | 4 | No | AATGG or AAAAG | 18853 or 471 |
| ACCAG | C02 | 11 | No | AACAG | 6846 |
| AAACCCTAGTGGGTTGT | C02 | 14 | No | AAC, TAGG, AGG | - |
| AAGCCACGC | C02 | 3 | No | AAG, AGC, ACGC | - |
| ACAT | C35 | 3 | Yes | - | - |
| AAACT | C13, C21 | 2, 3 | No | AACCT | 62281 |
| AATG | C12, C34 | 3, 5 | Yes | AAAG or AAGG | 708 or 11 |
| AGAT | C20, CC8 | 3, 2 | Yes | AGAA (AAAG) | 708 |

**Table 3.** Kmers that were lost during the course of the MA experiment. It is indicated if the putatively lost kmer is present in short arrays at low frequencies in the MA line in question.

634

| Lost kmer | Line(s) | Present in shorter arrays | Normalized mean abundance in other lines (copy number) |
|---|---|---|---|
| AAGGAATGG | C40 | No | 61.6 |
| AAAAAGAAGGAGATAG | C36 | No | 47.6 |
| AATAG | C25 | Yes | 16.3 |
| ACGCCCGAGC | C40 | Yes | 13.8 |
| AATCT | C20 | Yes | 7.2 |
| ACACCGACCACTACT | C40 | Yes | 5.6 |
| AAAGGCAACAACAGT | C35 | Yes | 4.2 |
| AATGT | C02, C34 | Yes | 12.4 |
| AAAAAAAC | C39, C44 | Yes | 7.5 |

**FIGURE CAPTIONS**

635

636    **Figure 1**. Abundance of the kmers that contain at least 100 counts after normalization. Kmers

637    are arranged in order of length, from 1-mers to 20-mers.

638

639    **Figure 2**. Network of relatedness of the 29 of the 39 most abundant kmers. Networks show

640    possible mutational steps between related kmers. Kmers are numbered based on their abundance

641    rank. The sizes of the circles also correspond to the kmer abundance.

642

643    **Figure 3**: Kmer mutation rates. (A) Boxplot of mutation rates in the MA lines (n=28). The top

644    20 kmers are shown, in descending order of abundance. Negative values indicate contractions

645    and positive values indicate expansions. (B) Plot of mean overall mutation rate across all 162

646    kmers at normalized abundance of at least 2 copies per 1x coverage (in descending order of

647    abundance). The red line indicates an overall mutation rate of 0.

648

649    **Figure 4.** Reduced variation in kmer content in the population compared to the MA lines. (A)

650    Total absolute amount of tandemly-repeated satellite sequences with length 1-20 bp in the MA

651    lines versus the population, including all 162 kmers with at least 2 copies per 1x coverage per

652    line. (B) Histogram of the cumulative copy changes across the top 39 kmers in the MA lines and

653    population. (C) The coefficient of variation across the top 39 kmers for the MA lines and the

654    population. Kmers are ordered by abundance. (D) The absolute normalized mutation rate (sum of

655    expansions and contractions), divided by the initial copy number of the kmer. Data from the MA

656    lines are shown in red and the population in blue.

657

31

658     **Figure 5.** Correlated changes in kmer abundance in (A) the MA lines, and (B) the population.

659     (C) shows the correlations between the family of 9 kmers that are closely related by sequence

660     and GC content, which are situated at the bottom corner of the matrix in B. Matrices were

661     computed from a pairwise Pearson correlation matrix between kmers in the deviation of their

662     abundance from the inferred ancestral abundance. Kmers are ordered by their GC content.

## LITERATURE CITED

Aldrup-MacDonald ME, Kuo ME, Sullivan LL, Chew K, Sullivan BA. 2016. Genomic variation within alpha satellite DNA influences centromere location on human chromosomes with metastable epialleles. Genome Res. 26:1301-1311.

Benjamini Y, Speed TP. 2012. Summarizing and correcting the GC content bias in high-throughput sequencing. Nucleic Acids Res. 40:e72.

Brinkmann B, Klintschar M, Neuhuber F, Hühne J, Rolf B. 1998. Mutation rate in human microsatellites: influence of the structure and length of the tandem repeat. American Journal of Human Genetics. 62:1408-1415.

Bzymek M, Lovett ST. 2001. Instability of repetitive DNA sequences: the role of replication in multiple mechanisms. Proc Natl Acad Sci. 98:8319-8325.

Charlesworth B, Sniegowski P, Stephan W. 1994. The evolutionary dynamics of repetitive DNA in eukaryotes. Nature 371, 215-220.

Charlesworth B, Charlesworth D. 2000. The degeneration of Y chromosomes. Philos Trans R Soc Lond B Biol Sci. 355:1563-1572.

Colbourne JK, Pfrender ME, Gilbert D, Thomas WK, Tucker A, Oakley TH, Tokishita S, Aerts A, Arnold GJ, Basu MK et al. 2011. The ecoresponsive genome of *Daphnia pulex*. Science. 331:555-561.

Ellegren H. 2004. Microsatellites: simple sequences with complex evolution. Nature Reviews Genetics. 5:435-445.

Fields PD, Reisser C, Dukić M, Haag CR, Ebert D. 2015. Genes mirror geography in Daphnia magna. Mol Ecol. 24:4521-4536.
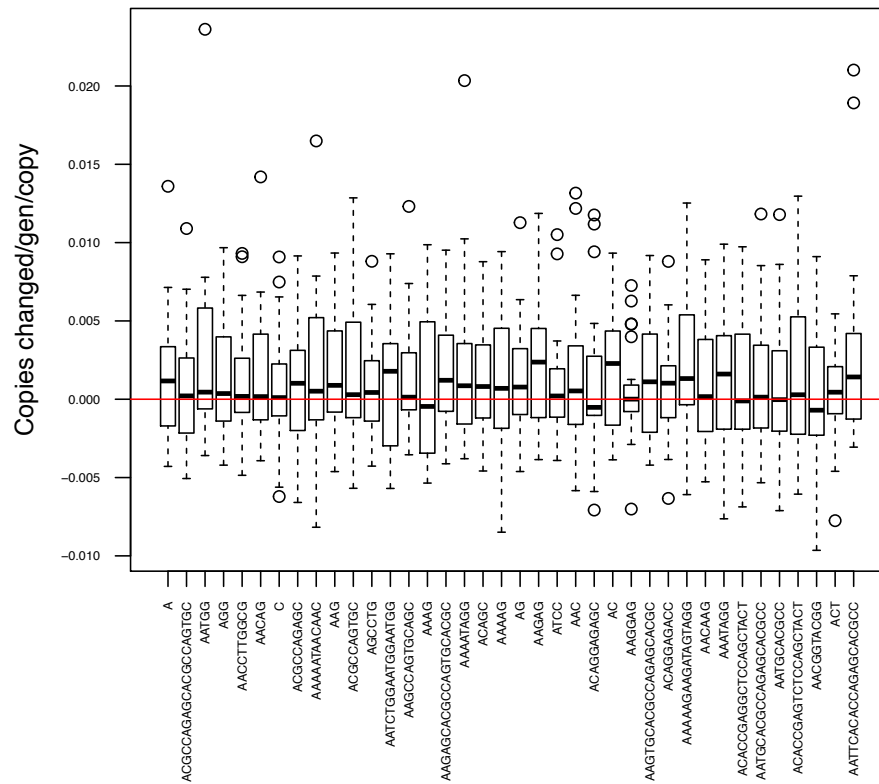
Flynn JM, Chain FJ, Schoen DJ, Cristescu ME. 2017. Spontaneous Mutation Accumulation in *Daphnia pulex* in Selection-Free vs. Competitive Environments. Mol Biol Evol. 34:160-173.

Fry K, Salser W. 1977. Nucleotide sequences of HS-α satellite DNA from kangaroo rat *Dipodomys ordii* and characterization of similar sequences in other rodents. Cell. 12:1069-1084.

Heger P, George R, Wiehe T. 2013. Successive gain of insulator proteins in arthropod evolution. Evolution. 67:2945-2956.

Henikoff S, Ahmad K, Malik HS. 2001. The centromere paradox: stable inheritance with rapidly evolving DNA. Science. 293:1098-1102.

Hoskins RA, Carlson JW, Kennedy C, Acevedo D, Evans-Holm M, Frise E, Wan KH, Park S, Mendez-Lago M, Rossi F et al. 2007. Sequence finishing and mapping of *Drosophila melanogaster* heterochromatin. Science. 316:1625-1628.

Keith N, Tucker AE, Jackson CE, Sung W, Lucas Lledo JI, Schrider DR, Schaack S, Dudycha JL, Ackerman M, Younge AJ et al. 2016. High mutational rates of large-scale duplication and deletion in *Daphnia pulex*. Genome Res. 26:60-69.

Kim JH, Ebersole T, Kouprina N, Noskov VN, Ohzeki J, Masumoto H, Mravinac B, Sullivan BA, Pavlicek A, Dovat S et al. 2009. Human gamma-satellite DNA maintains open chromatin structure and protects a transgene from epigenetic silencing. Genome Res. 19:533-544.

Lemos B, Araripe LO, Hartl DL. 2008. Polymorphic Y chromosomes harbor cryptic variation with manifold functional consequences. Science. 319:91-93.

Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics. 25:1754-1760.

Lohe AR, Brutlag DL. 1987. Identical satellite DNA sequences in sibling species of Drosophila. J Mol Biol. 194:161-170.

MacGregor HC, Sessions SK. 1986. The biological significance of variation in satellite DNA and heterochromatin in newts of the genus *Triturus*: an evolutionary perspective. Philos Trans R Soc Lond B Biol Sci. 312:243-259.

Malik HS. 2009. The centromere-drive hypothesis: a simple basis for centromere complexity. In: Anonymous Centromere. Springer. p. 33-52.

Orgel LE, Crick FH. 1980. Selfish DNA: the ultimate parasite. Nature. 284:604-607.

Paço A, Adega F, Meštrović N, Plohl M, Chaves R. 2015. The puzzling character of repetitive DNA in Phodopus genomes (Cricetidae, Rodentia). Chromosome Research. 23:427-440.

Padeken J, Zeller P, Gasser SM. 2015. Repeat DNA in genome organization and stability. Curr Opin Genet Dev. 31:12-19.

Palomeque T, Lorite P. 2008. Satellite DNA in insects: a review. Heredity. 100:564-573.

Platero JS, Csink AK, Quintanilla A, Henikoff S. 1998. Changes in chromosomal localization of heterochromatin-binding proteins during the cell cycle in Drosophila. J Cell Biol. 140:1297-1306.

Plohl M, Luchetti A, Meštrović N, Mantovani B. 2008. Satellite DNAs between selfishness and functionality: structure, genomics and evolution of tandem repeats in centromeric (hetero) chromatin. Gene. 409:72-82.

Raff JW, Kellum R, Alberts B. 1994. The Drosophila GAGA transcription factor is associated with specific regions of heterochromatin throughout the cell cycle. Embo j. 13:5977-5983.

Rehm P, Borner J, Meusemann K, von Reumont BM, Simon S, Hadrys H, Misof B, Burmester T. 2011. Dating the arthropod tree based on large-scale transcriptome data. Mol Phylogenet Evol. 61:880-887.

Schlötterer C. 2000. Evolutionary dynamics of microsatellite DNA. Chromosoma. 109:365-371.

Simmons MJ, Crow JF. 1977. Mutations affecting fitness in Drosophila populations. Annu Rev Genet. 11:49-78.

Stephan W, Cho S. 1994. Possible role of natural selection in the formation of tandem-repetitive noncoding DNA. Genetics. 136:333-341.

Subirana JA, Albà MM, Messeguer X. 2015. High evolutionary turnover of satellite families in Caenorhabditis. BMC Evolutionary Biology. 15:1.

Sung W, Tucker A, Bergeron RD, Lynch M, Thomas WK. 2010. Simple sequence repeat variation in the Daphnia pulex genome. BMC Genomics. 11:691.

Wei KH, Grenier JK, Barbash DA, Clark AG. 2014. Correlated variation and population differentiation in satellite DNA abundance among lines of *Drosophila melanogaster*. Proc Natl Acad Sci. 111:18793-18798.

Wei KH, Reddy HM, Rathnam C, Lee J, Lin D, Ji S, Mason JM, Clark AG, Barbash DA. 2017. A Pooled Sequencing Approach Identifies a Candidate Meiotic Driver in Drosophila. Genetics. 206:451-465.
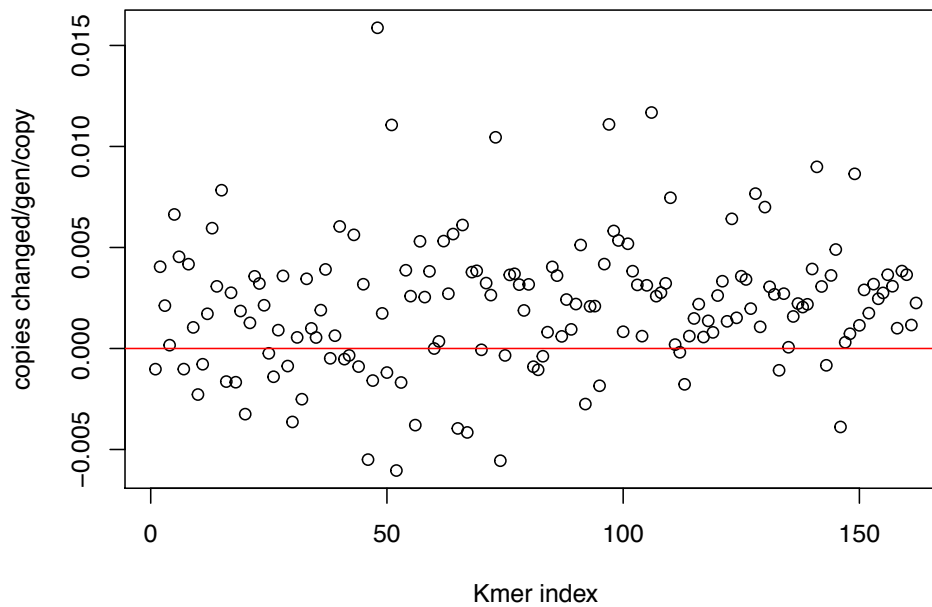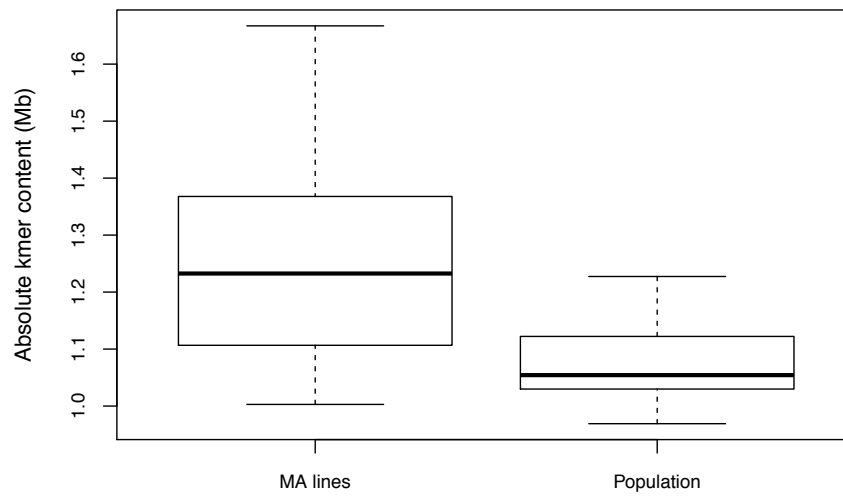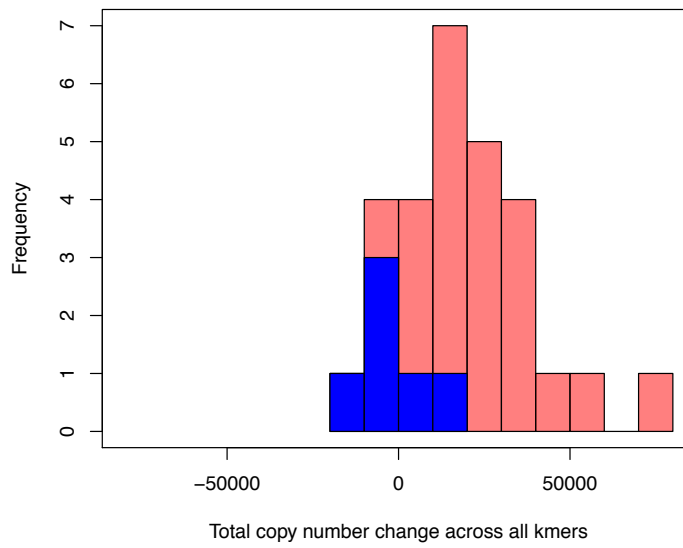
**Figure 1.**



**Figure 2.**
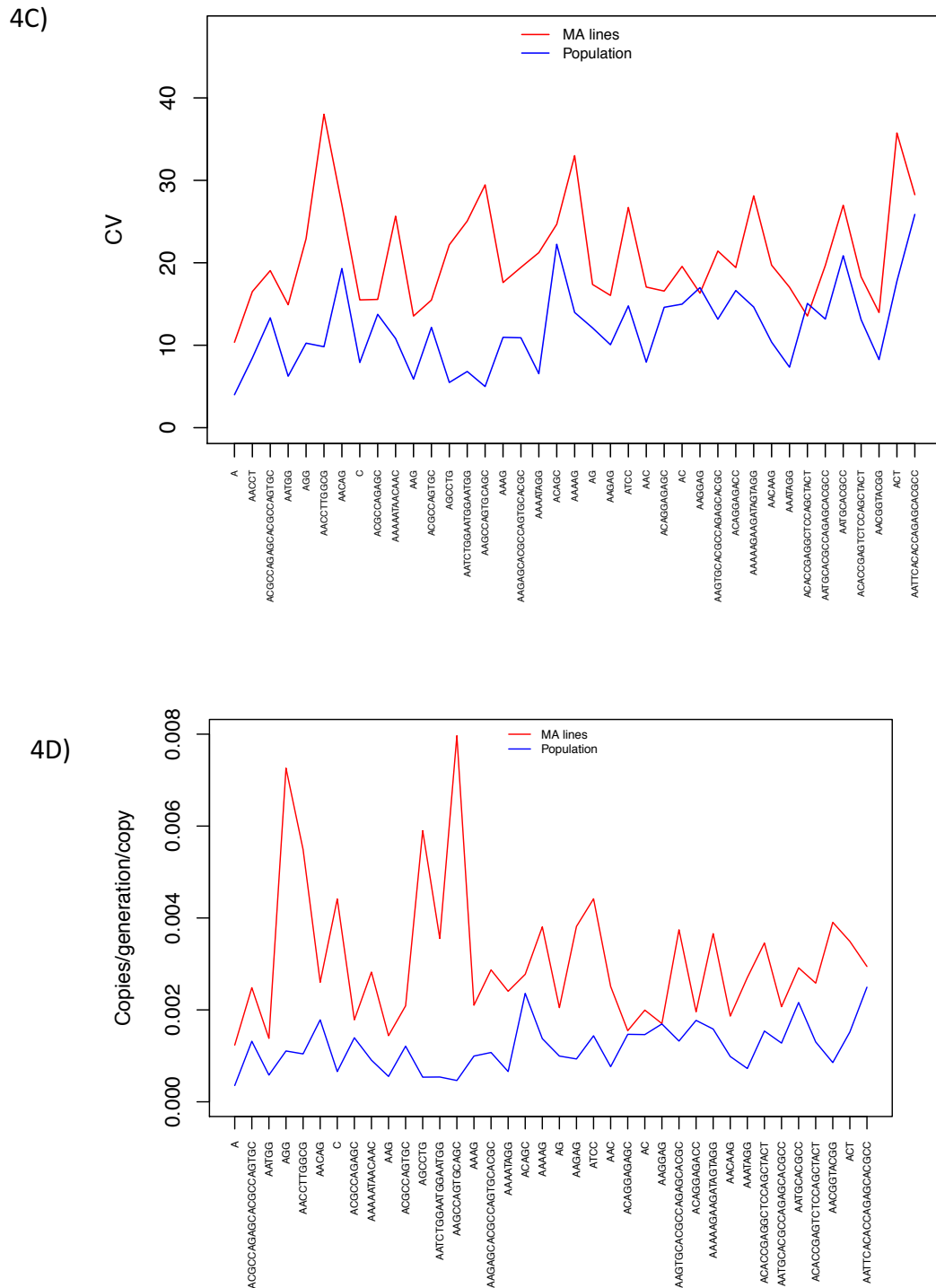
3A)



3B)



**Figure 3.**

4A)



4B)

4C)



4D)



**Figure 4.**

5A)



5B)

5C)



Figure 5.