

1 **Satellitome comparison of two oedipodine grasshoppers highlights the**
2 **contingent nature of satellite DNA evolution**

3

4 **Juan Pedro M. Camacho¹, Josefa Cabrero¹, María Dolores López-León¹, María**
5 **Martín-Peciña¹, Francisco Perfectti^{1,2}, Manuel A. Garrido-Ramos¹, Francisco J.**
6 **Ruiz-Ruano^{3,4,*}**

7

8 ¹Departamento de Genética, Universidad de Granada, 18071, Granada, Spain.

9 ²Research Unit Modeling Nature, Universidad de Granada, Granada, Spain.

10 ³Department of Organismal Biology – Systematic Biology, Evolutionary Biology
11 Centre, Uppsala University, SE-752 36, Uppsala, Sweden.

12 ⁴School of Biological Sciences, University of East Anglia, Norwich Research Park,
13 Norwich, NR4 7TU, UK.

14

15 *Corresponding author: Francisco J. Ruiz-Ruano (francisco.ruiz-ruano@ebc.uu.se,
16 f.ruiz-ruano-campana@uea.ac.uk)

17

18

19

20

21 **Abstract**

22 **Background:** The full catalogue of satellite DNA (satDNA) within a same genome
23 constitutes the satellitome. The Library Hypothesis predicts that satDNA in relative
24 species reflects that in their common ancestor, but the evolutionary mechanisms and
25 pathways of satDNA evolution have never been analyzed for full satellitomes. We
26 compare here the satellitomes of two Oedipodine grasshoppers (*Locusta migratoria* and
27 *Oedaleus decorus*) which shared their most recent common ancestor about 22.8 Ma ago.

28 **Results:** We found that about one-third of their satDNA families (near 60 in every
29 species) showed sequence homology and were grouped into 12 orthologous
30 superfamilies. The turnover rate of consensus sequences was extremely variable among
31 the 20 orthologous family pairs analyzed in both species. The satDNAs shared by both
32 species showed poor association with sequence signatures and motives frequently
33 argued as functional, except for short inverted repeats allowing short dyad symmetries
34 and non-B DNA conformations. Orthologous satDNAs frequently showed different
35 FISH patterns at both intra- and interspecific levels. We defined indices of
36 homogenization and degeneration and quantified the level of incomplete library sorting
37 between species.

38 **Conclusions:** Our analyses revealed that satDNA degenerates through point mutation
39 and homogenizes through partial turnovers caused by massive tandem duplications (the
40 so-called satDNA amplification). Remarkably, satDNA amplification increases
41 homogenization, at intragenomic level, and diversification between species, thus
42 constituting the basis for concerted evolution. We suggest a model of satDNA evolution
43 by means of recursive cycles of amplification and degeneration, leading to mostly
44 contingent evolutionary pathways where concerted evolution emerges promptly after
45 lineages split.

46

47 **Keywords:** Satellite DNA, Library Hypothesis, Satellitome Evolution, Cytogenomics.

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71 **Background**

72 Satellite DNA (satDNA) was first described by Kit [1] in mouse and guinea-pig DNA
73 with its repetitive nature demonstrated by Waring and Britten [2]. The first model for
74 satDNA evolution was devised by Smith [3], who demonstrated that DNA sequences
75 that are not maintained by natural selection evolve a tandem repeat structure due to
76 unequal crossing-over. Later, theoretical analyses assumed that satDNA evolution
77 usually depends on mutation, unequal crossing-over, and random drift, with purifying
78 selection controlling for excessive copy number [4,5,6,7,8,9,10,11].

79 Changes in satDNA amount are mainly due to unequal crossing-over, although
80 other mechanisms have been proposed to explain both amplification and spread of
81 satDNA repeats (for review, see Garrido-Ramos [12]). Walsh [13] proposed the
82 replication of extrachromosomal circles of tandem repeats by the rolling-circle
83 mechanism and reinsertion of replicated arrays as a powerful satDNA amplification
84 process, a mechanism for which Cohen et al. [14,15] have found some support.
85 Additionally, transposition may operate in satDNA emergence and amplification
86 [16,17,18,19]. Ultimately, replication-slippage might be an amplification process
87 [10,13], mainly involved in lengthening satellite monomers from basic shorter ones [20].

88 To explain the conservation of satellite sequences over long evolutionary periods,
89 Fry and Salser [21] suggested the Library Hypothesis. According to this hypothesis, a
90 group of related species should share a common library of satDNA sequences that
91 mostly show quantitative differences among species due to differential amplification.
92 Therefore, a given member of the library may appear as an abundant satDNA, while
93 others remain at low amounts and technically undetectable. Now we know that the
94 former can be visualized by FISH and the latter discovered by next-generation
95 sequencing [22]. Fry and Salser [21] suggested that an essential step in the evolution of

96 some satDNA families may be the acquisition of a biological function, in which case
97 natural selection would conserve its sequence for long evolutionary periods [23,24,25].

98 There are some examples of satDNA persisting for long, i.e., more than 40 Ma
99 (see Arnason et al. [26]; Garrido-Ramos et al. [27,28]; de la Herrán et al. [29,30];
100 Mravinac et al. [31,32]; Robles et al. [33]; Cafasso and Chinali [34]; Chaves et al. [35]).
101 Whereas the conservation of functional satDNA repeats is explained by purifying
102 selection (see references above), the persistence over time of other satDNA arrays
103 lacking apparent function might be simply due to chance events [8,9,13,37]. Therefore,
104 whether satDNA conservation in two or more species is just chance or due to selective
105 events remains unanswered.

106 Dover [37,38] suggested unequal crossing-over, gene conversion, and
107 transposition as molecular drive mechanisms for the concerted fixation of paralogous
108 variants, which operate independently of natural selection and drift. Recently, this
109 evolutionary pattern has been replaced by the birth-and-death model in the case of
110 coding multigene families [39,40]. Concerted evolution implies that paralogous copies
111 are more homogenized than orthologous ones when two species are compared. SatDNA
112 families comprise thousands or millions of copies of non-coding paralogous repeat units,
113 frequently arranged in many short arrays spread at different genomic locations
114 [17,22,41,42,43,44,45], so that fixation is improbable in these conditions. In fact,
115 although concerted evolution is the predominant pattern for satDNA evolution, non-
116 concerted evolution has also been reported and explained through various factors such
117 as life-history, population, location, organization, number of repeat-copies, or functional
118 constraints (for review, see Garrido-Ramos [12,44]). However, the ultimate causes for
119 concerted or non-concerted patterns are still unknown.

120 In this paper, we compare the full catalogue of satDNA families (i.e., the
121 satellitome) between two grasshopper species belonging to the subfamily Oedipodinae,
122 *Locusta migratoria* (Lmi) and *Oedaleus decorus* (Ode), which diverged 22.81 Ma [45].
123 We show the presence of about one-third of orthologous satDNA families whose
124 sequence comparison pointed to mutation and drift as the main drivers of satDNA
125 evolution. We also got estimates of nucleotide turnover rate at the level of consensus
126 sequences (consensus turnover rate, CTR), using 20 orthologous pairs present in both
127 species, and found that they were highly variable and depended on the history of
128 satDNA amplifications. We also analyzed repeat landscapes and developed indices for
129 satDNA homogenization and degeneration and an index for concerted evolution, which
130 may be useful for future research. Also, we propose a general model for satDNA
131 evolution and suggest that the evolution of these sequences constitute a good example
132 of contingent evolution (see Blount et al. [46]).

133

134 **Results**

135 **One-third of satDNA families showed sequence homology between species**

136 The range of variation for repeat unit length (RUL) was 8-400 bp for the 60 satDNA
137 families found in *L. migratoria* and 12-469 bp for the 58 families found in *O. decorus*.
138 For subsequent analyses we included only those satDNA families showing more than
139 100 copies, which excluded the four least abundant satDNAs in *L. migratoria*
140 (Additional file 1: Table S1). After comparing the consensus sequences of all satDNA
141 families present in both species, we found that 21 families in *O. decorus* showed
142 homology with 20 in *L. migratoria* (Additional file 1: Table S2). We assume that these
143 sets of satDNAs showing some sequence identity were already present in the most
144 recent common ancestor of these two species (dated about 22.81 Ma) and thus belonged

145 to the ancestor satDNA library. Therefore, these homologous sets constituted 12
146 orthologous superfamilies (OSFs) including 31 and 44 subfamilies in *O. decorus* and *L.*
147 *migratoria*, respectively (Additional file 1: Table S2). On the other hand, the non-shared
148 satDNA families (37 in *O. decorus* and 36 in *L. migratoria*) could have arisen *de novo*
149 after both lineages split, or else they were lost in one of the species.

150 Between species comparison of basic satellitome features (Table 1) revealed that
151 shared satDNAs did not show significant differences between species for RUL, A+T
152 content, and abundance, but divergence was lower in *L. migratoria*. However, the non-
153 shared satDNAs showed higher RUL and abundance in *O. decorus*. Within species
154 comparisons between shared and non-shared satDNAs failed to show differences in *O.*
155 *decorus*. In *L. migratoria*, however, the shared satDNA families showed higher RUL,
156 A+T content and abundance, and lower divergence, than the non-shared ones (Table 1).
157 Taken together, these results revealed the presence of many satDNA families showing
158 short monomers among the non-shared ones in *L. migratoria* which also showed lower
159 A+T content and abundance, but higher divergence than those shared with *O. decorus*.

160

161 **Tandem structure and association with other repetitive elements**

162 The quantification of homogeneous and heterogeneous read pairs allowed estimating the
163 degree of tandem structure (TSI) for each satDNA family (Additional file 1: Table S1).

164 The annotation of the heterogeneous read pairs allowed identifying other genomic
165 elements adjacent to satDNA (Additional file 1: Table S3). This revealed that

166 LmiSat03-195 (TSI= 99.7%) was associated with LINES in 57 out of the 100

167 heterogeneous read pairs observed. However, only 2% of the 1,356 heterogeneous read

168 pairs showed association with LINES for its orthologous OdeSat02-204 (TSI= 95.9%),

169 suggesting that association with LINES occurred only in *L. migratoria*. Likewise,

170 OdeSat17-176 and LmiSat02-176 showed association with Helitron TEs in 93% and
171 76% of the 2,379 and 1,356 heterogeneous read pairs observed, respectively. Bearing in
172 mind that the sequence of the LmiSat02-176 repeat unit shows homology with Helitron
173 TEs (Ruiz-Ruano et al. 2016), the high frequency of association with Helitron observed
174 for OdeSat17-176 and the low TSI (11.1%) suggest that most units detected for this
175 satDNA were part of the TE itself and are not in tandem (i.e., 1-TSI= 88.9%). However,
176 LmiSat02-176 showed high TSI (94.7%) and lower association with the TE (76%),
177 suggesting that this satDNA arose from this TE, but it also constitutes an independent
178 entity which has reached quite long arrays in *L. migratoria* (longer than 20 kb in the
179 MinION reads). The FISH pattern of both satDNAs (see below) reinforced this
180 conclusion, as OdeSat17-176 yielded no hybridization signals (Table 2), whereas
181 LmiSat02-176 showed pericentromeric bands on six chromosome pairs (see Ruiz-
182 Ruano et al. [22] and Additional file 1: Table S1).

183

184 **A same orthologous satDNA may show different FISH patterns at intra- and**
185 **interspecific levels**

186 FISH analysis for 14 OdeSat families, which showed homology with 20 LmiSat ones,
187 revealed that six OdeSats displayed conspicuous bands on chromosomes (B-pattern
188 from hereafter). In contrast, the eight remainders failed to show FISH signal (NS-
189 pattern from hereafter), of which seven showed the B-pattern in *L. migratoria* (Table 2).
190 This revealed that a same OSF may show FISH signals in one species but not in a close
191 relative.

192 To search for molecular differences between satDNAs showing the B- and NS-
193 patterns, we analyzed MinION long reads in *L. migratoria* to score the maximum array
194 length (MAL) for each LmisatDNA (Table 2). Even though coverage was very low

195 (0.02x), we found that none of the seven NS families analyzed showed arrays higher
196 than 2,500 bp, whereas almost half of those showing the B pattern did (Gardner-Altman
197 unpaired mean difference= 2930, 95.0%CI: 1540, 4790), and the three orders of
198 magnitude of the difference indicated that satDNAs with the B-pattern have been
199 submitted to more (and extensive) amplification events than those showing the NS-
200 pattern. This difference justifies using the presence of FISH signals as an indication of
201 the degree of satDNA amplification. The fact that 18 out of 20 orthologous satDNA
202 families in *L. migratoria* showed the B-pattern, whereas only six out of the 14
203 orthologous families analyzed in *O. decorus* showed it, represent the first indication for
204 a higher incidence of satDNA amplifications in *L. migratoria* (RxC contingency test,
205 with 50,000 replicates: P= 0.00562, SE= 0.00077). This result was reinforced by the
206 fact that the 14 OdeSat families included 24 subfamilies whereas the 20 LmiSat ones
207 included 44 subfamilies (Table 2) (Wilcoxon matched-pairs test: $z= 2.11$, $N=12$, $P=$
208 0.035). As subfamilies represent different amplification events, the former results
209 demonstrate that a same orthologous satDNA may show different amplification
210 trajectories during their independent evolution in different species.

211 Careful examination of orthologous satDNAs revealed a unique case of no
212 satDNA amplification in both species during the 22.8 Ma of separate evolution, as the
213 LmiSat27-57 and OdeSat41-75 OSF showed the same NS-pattern. Consistently with
214 their low degree of amplification, these two satDNAs showed very low values for
215 tandem structure (TSI: 9% in *O. decorus* and 32% in *L. migratoria*) and
216 homogenization (RPS: 29% and 32%) indices (see next section), indicating poor tandem
217 structure and homogenization (see Table 2 and Additional file 1: Table S4). The
218 remaining OSFs, however, showed amplification in at least one species. One of the
219 most dramatic differences was found for the orthologues OdeSat59-185 and LmiSat01-

220 185, which were the scarcest and the most abundant satDNAs in *O. decorus* and *L.*
221 *migratoria*, respectively, with the latter showing pericentromeric FISH bands on all
222 chromosomes [22] and OdeSat59-185 showing the NS-pattern (Fig. 2 and Table 2). In
223 fact, seven orthologous satDNA families with the NS-pattern in *O. decorus* showed the
224 B-pattern in *L. migratoria* (Table 2 and Fig. 3).

225 An interesting case was OSF7, where one of the five *L. migratoria* families
226 showed the NS-pattern (LmiSat24-266) whereas the four remaining (LmiSat28-263,
227 LmiSat43-231, LmiSat45-274 and LmiSat54-272) showed the B- pattern (Table 2).
228 Likewise, one of the two *O. decorus* families (OdeSat28-276) showed the B-pattern
229 whereas the other (OdeSat58-265) showed the NS one. This shows that homologous
230 satDNAs can display the NS or B patterns at intra- and interspecific levels. Finally,
231 even those satDNAs with FISH bands in both species showed remarkable differences
232 regarding chromosome location (proximal, interstitial, or distal; see Additional file 1:
233 Table S1). Taken together, these results show that orthologous satDNAs can display
234 disparate chromosome distribution in separate species due to their independent
235 evolution, a fact previously reported in the literature [47,48,49,50]. These differences
236 can range from short arrays being undetectable by FISH, which may eventually serve as
237 seeds for species-specific amplification (as suggested by Ruiz-Ruano et al. [22]), up to
238 long arrays yielding conspicuous FISH bands.

239

240 **SatDNA homogenization and degeneration**

241 SatDNA homogenization and degeneration are considered important drivers of satDNA
242 evolution, but their relative importance has been debated. It would thus be desirable to
243 find satDNA parameters being good indices for these two alternative states. To search
244 for a homogenization index, we hypothesized that it should show a high negative

245 correlation with intraspecific divergence. Spearman rank correlation analysis showed
246 that, in both species, RPS (relative peak size, see methods and Fig. 1) showed a very
247 high negative correlation with divergence (measured as K2P) ($r_s = -0.9$ in both species)
248 (Table 3), which revealed RPS as a good homogenization index. On the contrary, a
249 degeneration index should be negatively correlated with homogenization, and Spearman
250 rank correlations revealed that DIVPEAK (i.e. the divergence value showing the
251 maximum abundance in a repeat landscape, see Fig. 1) showed the highest negative
252 correlation index with RPS in both species (Table 3). This means that the relative size
253 of amplification peaks decreases as satDNA sequences accumulate divergence through
254 mutational decay since the last satDNA amplification (see repeat landscapes in Fig. 2,
255 Additional file 2: Fig. S1 and Additional file 3: Dataset 1).

256 To ascertain whether satDNA degeneration, measured by DIVPEAK, is
257 associated with any of the satDNA parameters analyzed (RUL, A+T, no. subfam and
258 TSI), we performed Spearman rank correlation analyses, which revealed that RUL was
259 the only satDNA property showing significant correlation with DIVPEAK (Table 3) and
260 it was negative and of similar magnitude as that between DIVPEAK and RPS. This
261 suggests that RUL is an important determinant of satDNA degeneration, with shorter
262 satDNAs degenerating faster. A possible explanation is that short monomers degenerate
263 faster through mutational decay because every point mutation implies a higher
264 proportion of degeneration for short than for long monomers, as if the Muller's ratchet
265 would have fewer teeth for short than long repeat units and the same number of new
266 mutations would imply a higher number of ratchet's turns for short repeating units than
267 for long ones.

268 The analysis of the statistical properties of RPS and DIVPEAK indicated that, in
269 both species, RPS fitted a normal distribution (ODE: $\chi^2 = 4.45$, $df = 3$, $P = 0.215$; LMI:

270 $\chi^2= 4.78$, $df= 3$, $P= 0.189$ whereas DIVPEAK fitted an exponential distribution (ODE:
271 $\chi^2= 4.55$, $df= 2$, $P= 0.103$; LMI: $\chi^2=4.93$, $df= 3$, $P= 0.177$). Their scales ranged between
272 0 and 1 for RPS and between 0 and 27% (within the 0-40% scale of divergence
273 measured here) for DIVPEAK.

274 To apply these indices to satDNA evolution, we consider that satDNA families
275 follow evolutionary pathways that include recursive cycles of homogenization (through
276 amplification by tandem duplication) and degeneration (through random mutation).
277 After an amplification event, homogenization (measured by RPS) will increase, and
278 degeneration (measured by DIVPEAK) will decrease. As time goes by, with no other
279 amplification events, RPS will decrease and DIVPEAK will move towards higher
280 values. An expected outcome of mutation accumulation is reducing the kurtosis of the
281 repeat landscape (RL) distribution (i.e., curve flattening, Fig. 1 for examples). In fact,
282 kurtosis was correlated negatively with DIVPEAK (Ode: $N=58$, $r_s= -0.80$, $t= 9.89$,
283 $P<0.000001$; Lmi: $N=56$, $r_s= -0.76$, $t= 8.58$, $P<0.000001$) and positively with RPS (Ode:
284 $N=58$, $r_s= 0.80$, $t= 9.68$, $P<0.000001$; Lmi: $N=56$, $r_s= 0.83$, $t= 10.98$, $P<0.000001$).
285 Kurtosis is thus proportional to RPS, so that highly homogenized satDNAs show
286 leptokurtic RLs whereas highly degenerated ones show platikurtic RLs. Therefore,
287 kurtosis and RPS are expected to be high for recently amplified satDNAs and low for
288 satDNAs that have not been amplified for a long time (see some examples in Fig. 2 and
289 Additional file 2: Fig. S1). Although these parameters do not constitute absolute
290 measures of time, however, they can be useful as measures of "time since the last
291 satDNA amplification". As satDNA can undergo successive amplifications across
292 evolutionary time, we can also consider RPS and kurtosis as homogenization indices
293 indicating how far is a satDNA from degeneration.

294 To analyze whether conservation of the orthologous satDNA families in both
295 species was associated with homogenization and degeneration indices, we compared
296 them between the shared and non-shared satDNA families found in each species. In *O.*
297 *decorus*, the effect size (unpaired mean difference) found between non-shared and
298 shared satDNAs by means of Gardner-Altman estimation plots, revealed no mean
299 differences for RPS (unpaired mean difference= -0.0682, 95.0%CI: -0.159, 0.0348),
300 kurtosis (unpaired mean difference= 0.678, 95.0%CI: -1.62, 5.78) and DIVPEAK
301 (unpaired mean difference= 1.13, 95.0%CI: -0.954, 5.61), indicating similar levels of
302 homogenization and degeneration in both groups. In *L. migratoria*, however, the three
303 indices showed differences between shared and non-shared satDNA families, indicating
304 higher homogenization and lower degeneration for the shared ones (Fig. 4).

305

306 **Amplification explains the concerted evolution of satDNA**

307 *O. decorus* and *L. migratoria* shared their most recent common ancestor 22.81 Ma, on
308 which basis we could perform estimations of interspecific rates of turnover in the
309 consensus sequences (CTR). For this purpose, we compared the consensus DNA
310 sequences of 20 pairs of orthologous satDNA, representing half of the 40 estimations
311 that could be done at family level (see Additional file 1: Table S2). The values obtained
312 for CTR in the 20 orthologous pairs ranged from 0.013% (between LmiSat02-176 and
313 OdeSat17-176) to 2.86% (between LmiSat03-195 and OdeSat02-204) nucleotidic
314 changes in their consensus sequences per million year (mean= 1.11%, see Table 2), with
315 two orders of magnitude between the extreme values.

316 To search for possible causes for such an extreme variation in the observed rates,
317 we performed forward stepwise multiple regression of CTR (dependent) on four factors
318 related to satDNA amplification: for each species, the number of subfamilies per

319 satDNA family (subfam), the absolute number of copies included in the 5% divergence
320 peak (peak-copies), RPS, and TSI. The results revealed that only three out of the eight
321 factors entered a model that explained 85% of the total variance in CTR, with
322 Ode_subfam explaining 56.4%, Ode_peak_copies explaining 25.7%, and TSI_Ode
323 explaining only a nonsignificant 2.8% (Table 4). Variance inflation factors of this
324 regression analysis ranged between 1.07 and 3.01 indicating the absence of
325 multicollinearity. Likewise, the standardized residuals of this regression fitted a normal
326 distribution (Shapiro-Wilks test: $W= 0.97$, $P= 0.82$). Finally, partial correlations were
327 0.85 for Ode_subfam, 0.76 for Ode_peak_copies, and 0.40 for TSI_Ode, whereas they
328 were much lower for the five factors failing to enter in the model (from -0.25 to -0.02).

329 As we defined satDNA subfamilies by sharing 95% or higher sequence identity,
330 i.e., up to 5% divergence, which was exactly the same figure used to define RPS and
331 DIVPEAK on RLS, we consider that the number of subfamilies actually represents the
332 number of independent amplification events being apparent within each family, as it
333 also coincides with the number of different consensus sequences per family. As peak-
334 copies represents the total number of repeat units in the amplification peak, we can infer
335 that the rate of nucleotide change estimated from consensus sequences (CTR), which is
336 positively correlated with the two former parameters, roughly represents the rate of
337 nucleotide changes driven by satDNA amplification to be part of the consensus
338 sequence. It was remarkable that only *O. decorus* variables entered in the stepwise
339 multiple regression model, as it is the species showing the lowest number of subfamilies
340 (31 versus 44 in the 12 OSFs, as a whole, and 24 and 44 in the 14 orthologous pairs
341 analyzed) and thus showed fewer amplification events, suggesting that CTR value is
342 limited by the species showing fewer amplification events. We thus conclude that the
343 same molecular mechanism, i.e., satDNA amplification, causes intraspecific

344 homogenization and interspecific diversification, thus explaining the concerted
345 evolution pattern of satDNA.

346 **Most satDNA families showed concerted evolution in both species**

347 Concerted evolution predicts that $CEI > 0$, and this was met for 16 orthologous pairs, the
348 four exceptions being the OdeSat17-LmiSat02 pair and three satDNA families in *O.*
349 *decorus* (OdeSat41, OdeSat57, and OdeSat59) where $CEI < 0$ thus showing signs of non-
350 concerted evolution (Table 2). Remarkably, these four OdeSats failed to display FISH
351 bands, suggesting that poor amplification might be related with non-concerted evolution.
352 In both species, CEI was positively correlated with RUL (Ode: $r_s = 0.70$, $N = 14$, $t = 3.4$,
353 $P = 0.0051$; Lmi: $r_s = 0.56$, $N = 20$, $t = 2.83$, $P = 0.011$) and RPS (Ode: $r_s = 0.73$, $N = 14$, $t =$
354 3.67 , $P = 0.0032$; Lmi: $r_s = 0.68$, $N = 20$, $t = 3.88$, $P = 0.0011$) but not with A+T content
355 ($P > 0.05$ in both species). In addition, CEI was positively correlated with TSI in *O.*
356 *decorus* ($r_s = 0.78$, $N = 14$, $t = 4.26$, $P = 0.0011$) but not in *L. migratoria* ($r_s = 0.43$, $N = 20$,
357 $t = 2.04$, $P = 0.056$). Finally, in *O. decorus*, CEI was higher in the six satDNAs showing
358 the FISH B-pattern than in the eight showing the NS-pattern (unpaired mean
359 difference = 2.63; 95% CI: 0.883, 5.36).

360 These results indicate that satDNAs displaying longer monomers, higher levels
361 of homogenization and the FISH B-pattern show higher indices of concerted evolution.
362 Exceptional non-concerted patterns were observed for satDNA families showing a low
363 number of amplifications since all showed a single subfamily in both species.

364

365 **The persistency of satDNA in these two species was not associated with functional** 366 **constraints**

367 Several sequence features have hitherto been associated with a variety of putative
368 satDNA biological roles, the most relevant being centromere function. We searched for

369 short internal repeats within each satDNA family's consensus sequences since these
370 repeats have been associated with sequence function. We found no direct repeats within
371 the sequence span of any satDNA sequence. On the contrary, it was common to find
372 short inverted repeats in all satDNA families that might facilitate non-B DNA
373 conformations such as stem-loops and cruciform structures, but they were found in both
374 shared and non-shared satDNA families.

375 To ascertain whether Gibbs free energy (dG) of satDNA sequence depends on
376 some satDNA properties, we performed forward stepwise regression, in each species,
377 with dG as dependent variable and RUL, A+T, sharing status and degeneration status
378 (DIVPEAK) as independent factors. In *Ode*, the regression model explained 67% of the
379 variance in dG (59% by RUL, 5% by A+T, and 3% by DIVPEAK). The correlation was
380 negative with RUL and positive with the two other factors. In *L. migratoria*, the result
381 was highly similar, except that DIVPEAK did not enter in the model, but the dG
382 variance explained was higher, reaching 83% (79% by RUL and 4% by A+T). As
383 higher free energy values correspond to lower dG values, the former results indicate that
384 free energy of satDNA sequence depends positively on RUL, as it determines the
385 likelihood of autopairing, and, at lower extent, also depends on two other sequence
386 properties influencing the number of hydrogen bonds in the double helix, as higher A+T
387 content implies more A-T pairs and fewer hydrogen bonds, thus lower free energy,
388 whereas higher DIVPEAK indicates higher mutational decay that might difficult
389 autopairing thus decreasing the number of hydrogen bonds. The fact that DIVPEAK of
390 the shared satDNAs was higher in *O. decorus* than *L. migratoria* (paired mean
391 difference= 2.6, 95.0%CI: 0.55, 6.8) is consistent with their higher degeneration in *O.*
392 *decorus*.

393 We found that most of the shared satDNA families failed to show a propensity to
394 acquire stable curvatures (Additional file 1: Table S1), even though the curvature-
395 propensity plots contained a peculiar maximum in some of them. However, the
396 magnitude of these peaks (11 to 13 degrees/10.5 bp helical turn) was far from the values
397 calculated for other highly curved motifs [51,52]. Most intriguingly, these peaks were
398 similar for satDNAs showing the NS or B FISH patterns or, in the latter case, whether
399 they were located on pericentromeric regions or not. In total, only 11 (7 in *L. migratoria*
400 and 4 in *O. decorus*) out of the 34 shared satDNA families showed curvature propensity,
401 all showing $RUL \geq 185$ bp. They belonged to five different OSFs, three of which showed
402 curvature propensity in both species, whereas the two remaining showed it in only one
403 species, suggesting that this property does not depend only on RUL, which was highly
404 similar in both species for these satDNA families.

405 We also analyzed curvature propensity for the non-shared satDNAs, and none
406 of them showed it to a large degree. Notwithstanding, as observed for shared satDNAs, a
407 few families (one in *L. migratoria* and five in *O. decorus*) showed a conspicuous peak
408 of magnitudes between 11 to 14 degrees/10.5 bp helical turn. It has been suggested that
409 DNA curvature may be involved in the recognition of DNA-binding protein
410 components of the heterochromatin [53]. Our results show that curvature propensity is
411 not differentially frequent or relevant in the 34 shared satDNAs analyzed in both species,
412 compared with the non-shared ones. Therefore, we believe that curvature propensity is
413 not a relevant feature of satDNA or the cause for satDNA conservation in these two
414 species.

415 Finally, we searched for the presence of short sequence motifs common to the
416 shared satDNA families in both species. We isolated individual monomers from each
417 satDNA family and calculated nucleotide diversity (π) per position (not shown). We did

418 not find conserved motifs in these satDNAs, irrespectively of their FISH pattern or
419 chromosomal location.

420 Taken together, these results show that, in these two species, there is no
421 sequence conservation for pericentromeric satDNAs, which also lack significant
422 sequence signatures other than A+T richness and repeat length. On the other hand, all
423 putative functional signatures analyzed here were not more frequent in the shared
424 satDNAs than in the non-shared ones. We interpret this as evidence that satDNA
425 conservation is mostly a contingent event. This conclusion is logically conditioned by
426 data and methodology limitations, such as testing based just on sequence data and
427 genomic location, and using a long time scale.

428

429 **Incomplete sorting of the satDNA library**

430 The satellitomes of relative species show sequence homology for a fraction of their
431 satDNA families, which is the best support for the satDNA library hypothesis [21]. Joint
432 analysis of RLs and MSTs revealed interesting properties of the satDNA library (Fig. 2
433 and Additional file 2: Fig. S1): i) OdeSat02A and LmiSat03A were the two OSF02
434 subfamilies showing the highest amplification peaks in the RLs (Fig. 5a, plot on the
435 left), and they also showed the highest CTR observed among all those analyzed here
436 (2.86% per Ma). Remarkably, the MST plot for all subfamilies and families comprising
437 OSF02 revealed complete sorting per species for this component of the library (Fig. 5a,
438 right). ii) On the other hand, OSF12 included two families in *L. migratoria* (LmiSat01
439 and LmiSat13) which were fully sorted in the MST (Fig. 5b, right), whereas the single
440 *O. decorus* family (OdeSat59) was remarkably similar to LmiSat01A, with only two
441 nucleotidic differences in their sequence, which is lower than those shown by the four
442 other *L. migratoria* subfamilies with LmiSat01A. This illustrates an extreme case of

443 incomplete library sorting (ILibS) and the second lowest CTR value (0.26% per Ma).
444 Other OSFs showed intermediate situations. For instance, OSF04 showed CTR values
445 between 1.16 and 1.60 and their MST revealed the existence of ILibS, with OdeSat32A
446 being connected with three different LmiSats (37A, 26A and 51A), the latter being
447 placed between OdeSat32A and OdeSat21A (see Additional file 2: Fig. S1a). On the
448 contrary, OSF5 (Additional file 2: Fig. S1b) showed high CTR values (>2% per Ma)
449 and complete library sorting, with the satDNAs properly separated between species.
450 Finally, OSF07 showed CTRs between 0.56 and 1.43 and apparent ILibS, with high
451 level of intermixing between the satDNAs of both species (Additional file 2: Fig. S1c).
452 Taken together, these observations suggest that CTR values are inversely associated
453 with the level of ILibS. On this basis, we used the maximum CTR value (maxCTR=
454 2.86) as reference to estimate the degree of ILibS as one minus the quotient between
455 CTR_i and maxCTR (see Table 2). This indicated that the satDNA library of *O. decorus*
456 and *L. migratoria* shows, on average, 61% of incomplete sorting after 23 Ma. Finally,
457 the fact that the four OdeSats showing the non-concerted pattern were those showing
458 the highest ILibS figures (0.88-1), whereas ILibS values up to 0.84 corresponded with
459 patterns of concerted evolution (see OSF8 in Table 2), suggested the possible existence
460 of a threshold for ILibS (between 0.84 and 0.88) below which satDNA evolution is
461 concerted.

462

463 **Discussion**

464 **SatDNA evolution is mostly contingent**

465 Comparative analysis of the satellitome in the grasshoppers *O. decorus* and *L.*
466 *migratoria*, two species belonging to the Oedipodinae subfamily, which shared their
467 most recent common ancestor about 23 Ma, gave us a chance to take a look into

468 satDNA library evolution during this period. We assume that the 41 satDNA families
469 (20 in *L. migratoria* and 21 in *O. decorus*) that showed sequence homology between
470 species belong to 12 orthologue groups already present in the ancestor library, which
471 have been conserved up today. However, the remaining 84 families (36 in *L. migratoria*
472 and 37 in *O. decorus*) could represent either remnant satDNAs conserved in only one
473 species or satDNAs arisen *de novo* during the separate evolution of these species. To
474 distinguish between these two possibilities, it is necessary to analyze other oedipodine
475 species. The occurrence of a species-specific profile of satDNAs resulting from
476 differential amplifications and/or contractions from a pool of sequences shared by
477 related genomes is a prediction of the library hypothesis of satDNA evolution with the
478 subsequent replacement of one satDNA family for another in different species [21]. By
479 analogy with incomplete lineage sorting (ILS) in phylogenetic studies, satDNA
480 amplifications and/or contractions between close relative species may yield a pattern of
481 incomplete library sorting (ILibS). We have detected here this phenomenon using
482 consensus sequences, but the use of physical sequences would yield even higher rates of
483 ILibS.

484 The library hypothesis predicts the residual retention of low-copy counterparts
485 of the dominant satDNA of one species in the other [21]. For instance, OdeSat02A-204
486 and LmiSat03A-195 have been independently amplified in both species, reaching
487 among the highest genomic abundances in both species, and showed the highest CTR
488 and extensive diversification, with four subfamilies in *O. decorus* and six in *L.*
489 *migratoria* (see Fig. 5a). In addition, a joint MST for OSF02 (to which both satDNA
490 families belong) revealed the absence of ILibS as all satDNA families and subfamilies
491 appeared well separated between species in the MST (see Fig. 5a). Conversely, the
492 consensus sequences of LmiSat01A-185 and OdeSat59-185 only differed in two

493 positions, thus showing higher interspecific similarity than that found, at intraspecific
494 level, between the five *L. migratoria* subfamilies (see Fig. 5b), thus constituting an
495 extreme example of ILibS. The high similarity in the consensus sequences of
496 OdeSat59A and LmiSat01A cannot be explained by functional conservation because
497 only the latter shows FISH bands on centromeric regions of all chromosomes thus
498 probably playing a centromeric function in *L. migratoria*, whereas OdeSat59A is the
499 most scarce satDNA found in *O. decorus* thus being only a relic. Likewise, while
500 OdeSat01-287 is the most abundant satDNA in *O. decorus*, its orthologous (LmiSat09-
501 181) is a relict in *L. migratoria*. We thus believe that the observed sequence similarity
502 between OdeSat59A and LmiSat01A might be due to chance convergence, as the
503 likelihood of nucleotide coincidence in each position of the consensus sequence is a
504 function of the relative frequency of the four possible nucleotides in each species, thus
505 being a probabilistic issue.

506 Our estimates of ILibS from CTR values indicated that the satDNA libraries of
507 *O. decorus* and *L. migratoria* still show 61% of incomplete sorting after 23 Ma of
508 independent evolution, i.e. about 39% of complete sorting (1.7% per Ma). This extreme
509 cohesiveness of the satDNA library is due to the highly paralogous nature of these
510 genomic elements, with thousand copies evolving at once, independently in both species,
511 through point mutation, amplification (tandem duplication) and drift (see below). This
512 39% expresses only part of library divergence, as the maximum divergence would be
513 reached when all homology signals between satDNAs in both species would have been
514 erased, as in the case of the non-shared ones, whereas the satDNAs belonging to OSF02
515 are still recognized as homologous between species even with 100% library sorting.
516 Anyway, the ILibS parameter of a given OSF (or orthologous pair of satDNAs)
517 inversely indicates its possible utility for phylogenetic analysis.

518 Another prediction of the library hypothesis is that the appearance of satDNA
519 families would usually represent amplification of one of the satellites already present at
520 a low level in the library, rather than actual *de novo* appearance. It is not easy to know if
521 any of the non-shared satDNA families actually arose *de novo*. However, in *L.*
522 *migratoria*, the lower RUL of non-shared satDNAs suggests that the satellitome of this
523 species might harbor some *de novo* arisen short satellites, in consistency with an
524 evolutionary trend towards increasing monomer length and complexity, suggested by
525 theoretical [54] and experimental [20,27,29,55] work.

526 Our estimates of CTR by the comparison of 20 orthologous pairs of satDNA
527 families indicated that it was 1.11% per Ma, which implies that two satellites can
528 diverge by more than 50% in about 50 Ma. This explains why *L. migratoria* and *O.*
529 *decorus*, belonging to the Acrididae family do not share a single satDNA family with
530 *Eumigus monticola* [56], a grasshopper belonging to the Pamphagidae family, as these
531 two orthopteran families shared their most recent common ancestor about 100 Ma [45].
532 Along with the stochastic nature of satDNA loss or gain during evolution, sequence
533 changes at the mentioned rate will make unrecognizable a satDNA family after 100 Ma
534 of separate evolution within the genomes of different species, which contrasts with the
535 case of some other satDNAs preserved for more than 60 Ma [28,30,31,34] or even more
536 than 100 Ma [29,33].

537 Our results suggest that the same OSF may be involved in the centromeric
538 function in a given species but not in a close relative species. According to Melters et al.
539 [57], the most abundant satDNAs in a genome are most likely involved in the
540 centromeric function. Another feature suggesting this fact is satDNA location on
541 pericentromeric regions of all chromosomes. Therefore, LmiSat01-185, OdeSat01-287
542 and/or OdeSat02-204 are the best candidate families in these species since all meet the

543 two conditions. However, all three satDNAs showed orthologous families in the other
544 species displaying much more limited chromosome distribution, suggesting that one or
545 both species have replaced the centromeric satDNA during the last 22.8 Ma. No
546 significant track of signatures such as conserved motifs or sequence mediated specific
547 stereo-spatial features were found for these or any other pericentromeric satDNAs found
548 in these species. We thus believe that, in the absence of other evidence, contingent facts
549 such as the opportunity to be in the right place when amplified might be responsible for
550 centromeric satDNA turnover. Zhang et al. [58] also revealed rapid divergence for
551 centromeric sequences among closely related *Solanum* species and suggested that
552 centromeric satellite repeats underwent boom-bust cycles before a favorable repeat
553 became predominant in a species. Indeed, there are species such as chicken [60],
554 common bean [60], or pea [61] that contain different satDNAs in different centromeres.

555 Whether a given satDNA is conserved for long due to functional reasons is an
556 open question. Fry and Salser [21] suggested that an essential step in the evolution of a
557 specific satDNA family may be acquiring a biological function. However, persistence
558 over time of a satDNA might also be explained in terms that do not depend on natural
559 selection [8,9,10,13,36]. Our results were consistent with this latter view. No conserved
560 functional motifs were found within the monomers of every grasshopper satDNA
561 analyzed as has been found in other satDNAs such as human centromeric satDNA
562 [62,63,64,65]. On the other hand, short dyad symmetries within satDNA repeats might
563 be associated with thermodynamically stable secondary structures and yield non-B-form
564 conformations, such as stem-loops or cruciforms. It has been claimed that these short
565 dyad symmetries may play an important role in satDNA repeats as targets for protein
566 binding and thus in satDNA function [12,44,53,66,67,68,69]. Kasinathan and Henikoff
567 [70] have proposed that that cruciform structures formed by dyad symmetries may

568 specify centromeres and that these non-B form DNA configurations in centromeric
569 repeats may facilitate centromere assembly [70,71]. In the two grasshopper species
570 analyzed here, short inverted repeats that might facilitate dyad symmetries and non-B
571 DNA conformations were frequent in both shared and non-shared satDNAs,
572 independently of their organization and chromosomal location. We believe that this
573 property is a simple outcome of stochastic processes of satDNA evolutionary dynamics.
574 Its ubiquity suggests that almost any satDNA can be recruited for functions being
575 dependent on the formation of non-B DNA conformations (see Kasinathan and
576 Henikoff [70]).

577 SatDNA evolution is a topic of high interest for the scientific community, but
578 the processes and mechanisms have sometimes been confused. Molecular drive was a
579 turnover mechanism suggested by Dover [37,38] as a directional force leading to repeat
580 fixation. It has been the prevalent hypothesis for satDNA evolution due to its apparent
581 explicative power as a mechanism for sequence change, turnover, and concerted
582 evolution. Nonetheless, when applied to satDNA, the presence of arrays on multiple
583 genomic sites makes it impossible, in practice, the fixation of a given repeat. The
584 dependence of CTR on the number and extent of satDNA amplifications in *O. decorus*
585 suggests that molecular drive mainly operates through satDNA amplification and is thus
586 a mutational force (e.g. tandem duplication by means of unequal crossing-over).
587 However, the reach of satDNA amplification is limited to changes in the relative
588 abundances of the pre-existing sequence variants for a given family, most frequently
589 leading to incomplete turnovers. A good way to visualize the role of molecular drive (or
590 amplification) in satDNA evolution is through repeat landscapes for families consisting
591 of several subfamilies showing platykurtic curves (i.e. with low abundance and high
592 divergence) and one or two subfamilies displaying leptokurtic distributions (i.e. with

593 high abundance and low divergence) (see Fig. 2 and Additional file 2: Fig. S1), the
594 latter being those sequences that acquire relevance through satDNA amplification. The
595 comparison of orthologous satDNA pairs between species thus reveal that satDNA
596 amplification implies molecular drive or drift at intra- and inter-specific levels,
597 respectively.

598 The high or low degree of homogenization for a given satDNA is inversely
599 proportional to the time since the last amplification. It thus depends on i) the neutral
600 mutation rate introducing new sequence variants (increasing intra-specific divergence)
601 and ii) the rate of satDNA amplification, implying partial turnovers that promote
602 sequence variants that become new subfamilies. As satDNA amplification for
603 orthologous satDNA families is independent in relative species, it behaves as an inter-
604 specific drifting mechanism. This dual role of satDNA amplification as the major
605 homogenizing force at the intraspecific level and as the principal driver for interspecific
606 sequence divergence, forced by reproductive barriers, inevitably leads to the concerted
607 evolution pattern. In fact, 16 pairs of orthologous satDNAs met this pattern, with only
608 four showing a non-concerted one. Remarkably, these exceptions coincided with the
609 absence of major amplifications in *O. decorus* satDNAs that remain at low abundance.
610 This kind of variation can persist for long in the absence of (homogenizing)
611 amplification events [72]. Therefore, concerted evolution should be a reasonable
612 consequence of the stochastic nature of satDNA evolution, while exceptional non-
613 concerted patterns can result from differential amplifications among species. Other
614 exceptions can result from satDNA homology with TEs, as was the case for LmiSat02-
615 176, whose homology with Helitron might have biased the calculation of intraspecific
616 divergence. Other explanations have been raised as possible causes for non-concerted
617 evolution patterns, such as the effect of location, organization, and repeat-copy number

618 [55,72,73], population and evolutionary factors [29,33,75,76,77], biological factors
619 [68,77], or functional constraints [32].

620 We have shown here that concerted evolution is a pattern emerging from
621 satDNA amplification due to the resulting homogenization at intraspecific level and
622 diversification at interspecific level. To visualize this relationship, think about two
623 species recently emerged from a common ancestor. Their satDNA libraries are almost
624 identical at interspecific level but both retain the ancestral polymorphism at intraspecific
625 level. This situation would imply, for each OSF, ILibS values next to 1 and CEI<0 since
626 divergence would be higher at intra- than inter-specific level. As time goes by and
627 mutation and drift operate, ILibS will decrease and CEI will increase as new mutations
628 occur independently in both species. In absence of satDNA amplification, mutation and
629 drift would lead satDNA towards concerted evolution by increasing interspecific
630 divergence, although this process would be slow. However, the pathway to concerted
631 evolution would be paved away by satDNA amplification as the resulting
632 homogenization would reach CEI>0 values (by sharply decreasing intraspecific
633 divergence) when ILibS would decrease below a threshold which, in the case of *O.*
634 *decorus* and *L. migratoria*, lies between 0.84 and 0.88. The fact that this threshold is so
635 close to 1 reinforces the idea that concerted evolution is an unavoidable property fastly
636 emerging from satDNA amplification. In fact, the four satDNA families which in *O.*
637 *decorus* showed signs of non-concerted evolution showed low levels of homogenization
638 (RPS between 0.29 and 0.40) and high values of ILibS (0.88-1), presumably due to the
639 low level of amplification of these four satDNAs in this species. Taken together, our
640 results indicate that concerted evolution is a state of interspecific diversification of the
641 satDNA library, reached below a given ILibS threshold, which is fastly promoted by
642 satDNA amplification.

643

644 **A model for satDNA evolution**

645 Considering all findings derived from the quantitative analysis of 114 satDNAs in *O.*
646 *decorus* and *L. migratoria*, we suggest the following model for satDNA evolution (Fig.
647 6). Intragenomic changes are mainly stochastic, implying that satDNA families mainly
648 evolve under the domain of mutation and drift. SatDNA arises from any tandem
649 duplication yielding at least two monomers. Subsequent unequal crossover is the main
650 source for longer arrays with the consequent increase in tandem structure. This tandem
651 duplication is one of the two classes of mutation operating on satDNA. The other is
652 point mutation increasing divergence among the different monomers composing the
653 whole set of satDNA sequences belonging to a given family. When tandem duplication
654 occurs massively during a short time, it constitutes an **amplification** event that
655 decreases intra-specific divergence (i.e., increases homogenization as measured by RPS)
656 by adding a high number of repeats showing identical sequence. Next, intra-specific
657 divergence will grow across years by the incidence of point mutations, inevitably
658 leading to the **degeneration** of the satDNA sequence unless new amplifications occur.
659 This is characterized by a temporal decrease of RPS and kurtosis and an increase of
660 DIVPEAK as family sequences became more and more divergent. From time to time,
661 some monomers will lose their identity as members of a given satDNA family (reaching
662 identities lower than 80%) or even as members of the same superfamily (with no
663 recognizable homology). This process may shorten long arrays into pieces, thus
664 decreasing TSI and, finally, the satDNA may fade away across time.

665 Each new amplification event drives a satDNA family away from degeneration
666 (by promoting that a given subfamily shows the highest abundance and
667 homogenization), after which new point mutations will drive it towards

668 degeneration again, and even complete disappearance if new amplifications do
669 not take place. In summary, we suggest that satDNA undergoes recursive cycles
670 of amplification-degeneration that may keep them in the genome for a long time.
671 During this time, they can integrate into longer repeat units or higher-order
672 structures [79,80], or else disappear through sequence degeneration and/or
673 unequal crossover. The fact that short satDNAs degenerate faster than the longer
674 ones (see above) suggests that their cycle is usually shorter than that of long
675 satDNAs, partly explaining why many short satDNAs show high K2P
676 divergence and platykurtic distribution. For instance, LmiSat10-9 is made of
677 monomers of only 9 bp and is not found in Ode. Even if it would have been
678 present in the common ancestor, it is doubtful that it would have remained for
679 22.8 Ma in both species without losing identity in at least one of them. In fact,
680 there seems to be a minimum monomer length for homology conservation in
681 these two species, which was 57 bp (LmiSat27-57 and OdeSat41-75).
682 Alternatively, a satDNA formed by repeats of only 9 bp could have arisen *de*
683 *novo*, by chance, in the gigantic genome of *L. migratoria* [22].

684 In addition to all former intragenomic events, satDNA frequently undergoes
685 spread among chromosomes. Transposition and replication of extrachromosomal circles
686 of tandem repeats, by the rolling-circle mechanism, followed by reinsertion of
687 replicated arrays, have been postulated as the main mechanisms for the amplification
688 and spread of satDNA families and is supported by indirect [43,81] or direct [14,15]
689 evidence.

690 At intergenomic (population) level, the only conceivable way to spread an
691 amplification event (occurred in a single individual) is through differential reproduction,
692 as we believe that the molecular drive mechanism suggested by Dover [37,38] as a non-

693 selective fixing force even at the population level, is circumscribed at the intragenomic
694 level. Differential reproduction can occur at random, i.e., by genetic drift, or non-
695 random, i.e., through selection. The latter may be negative, setting up an upper limit to
696 the amount of satDNA tolerable by a genome. Purifying selection, mutation and drift
697 are the drivers in the mutational-hazard (MH) hypothesis [82,83], which suggests that
698 the efficacy of purifying selection is impaired by genetic drift in small populations. This
699 is especially applicable to satDNA, where CTR is highly variable among families
700 (intragenomically). The fact that all satDNA families within a genome have been
701 submitted to the same demographic changes at population level (excepting the
702 differences due to sex linkage) means that purifying selection appears to set few limits
703 to the variation in nucleotide substitution rate among satDNA families. Interestingly, 18
704 out of 20 shared satDNA families in *L. migratoria* showed amplification events giving
705 rise to FISH bands, whereas only six out of their 14 orthologous families in *O. decorus*
706 did it. This reveals that many of these OSFs have shown highly different evolutionary
707 paths in both species. Based on the MH hypothesis, we may speculate that the extreme
708 demographic changes associated with locust outbreaks in *L. migratoria* might have
709 helped to spread individual satDNA sequences at the population level during the
710 extreme bottlenecks that characterize the solitary phase and subsequent population
711 expansions during the gregarious one. This issue needs further research, including
712 quantitative population analyses of every satDNA family in this species.

713 In addition, selection can operate positively through non-phenotypic (i.e.,
714 meiotic drive) or phenotypic (functional recruitment) effects, as is the case for
715 centromeric and telomeric repeats. The latter is the extreme example of functional
716 recruitment since the repeat is actively homogenized by an RNA-protein complex
717 (telomerase) coded by the genome. Centromeric satDNA in primates resembles this

718 kind of recruitment as another gene (CENPB) is involved in the organization of
719 centromeric satDNA [62,63,64,65].

720 Our model is an extension of the models devised in the '70s and '80s
721 [4,5,6,7,8,9,10,11], with some more emphasis on the intragenomic level, and under the
722 light of the MH hypothesis [82,83]. Briefly, amplification is the homogenizing force of
723 satDNA whereas point mutation causes sequence degeneration, with both forces acting
724 recursively. We believe that our model brings about some essential term clarifications.
725 For instance, Escudeiro et al. [84] recently suggested a model of satDNA evolution in
726 bovids consisting of three stages, namely amplification, degeneration (deduced from
727 high satDNA similarity between some species and low between others) and
728 homogenization (high sequence identity among all species). These authors thus claimed
729 for degeneration and homogenization as if they were inter-specific processes. However,
730 in our model, both processes are intragenomic (i.e., intra-specific) resulting from
731 satDNA amplification and point mutation, respectively, whereas inter-specific
732 homogenization or degeneration is highly unlikely under contingent evolution. In fact,
733 homogenization to an identical sequence in several species could only be achieved by
734 functional (selective) recruit, as that occurred for the telomeric DNA repeat.

735 Finally, the paralogous nature of the satDNA library implies that its
736 diversification between species may show high levels of incomplete library sorting, and
737 this may be a problem for the use of satDNA for phylogenetical purposes beyond
738 satDNA evolution itself. However, the pathway followed by an ancestor satDNA library
739 after speciation can be monitored by satellitome comparison, as shown here for *O.*
740 *decorus* and *L. migratoria*. A new body of research is taking form recently about
741 contingency and determinism in evolution [46], trying to answer Gould's question on
742 whether evolutionary trajectories are repeatable [85]. In this respect, satellitome

743 evolution is a natural "parallel replay experiment" able to show many properties of
744 contingent evolution, as the initially identical libraries in the ancestor undergo
745 independent evolution after speciation reaching a high diversity of outcomes among
746 different OSFs. Within species, the environment (at both intragenomic and population
747 levels) is the same for all satDNA families (except for genomic location and
748 organization), but the pathway followed by each of them is highly variable: some
749 families show consensus sequences being highly similar to those in the other species,
750 thus showing high ILibS, whereas others are completely sorted between species, and
751 still others are unrecognizable between species because they have arisen *de novo* in one
752 species or else they have undergone so many sequence changes that have lost homology
753 between species. In analogy with Blount et al. [46] claiming at ecological level, the
754 evolutionary trajectory followed by each OSF in the satellitomes of two separate species
755 is mainly influenced by stochastic processes (i.e. mutation and drift), most likely
756 reaching different outcomes even when both species satellitomes started from the same
757 state in the ancestor and the different OSFs evolved under almost identical conditions at
758 intragenomic level. Therefore, the satellitome is a good example of contingent evolution
759 supporting that "disparate outcomes become more likely as the footprint of history
760 grows deeper" [46]. A rough estimate of the minimal degree of contingent evolution in
761 the *O. decorus* and *L. migratoria* satellitomes can be obtained from the 20 orthologous
762 satDNA pairs used here to estimate CTR. As Table 2 shows, only two of them showed
763 identity higher than 95%: OdeSat17-176/LmiSat02-176 showing a single nucleotide
764 difference in their consensus sequences, and OdeSat59-185/LmiSat01A-185 showing
765 two differences. The first pair showed homology with Helitron TEs which could have
766 biased identity calculations, and the second one appears to have little to do with
767 functional conservation (as explained above). Even assuming that these two cases are

768 adaptive convergences (which is unlikely), we can estimate that satDNA evolution in
769 these species was at least 90% contingent.

770 The comparison of the satellitomes in two grasshopper species belonging to the
771 subfamily Oedipodinae has allowed us to develop several indices that have proven to be
772 highly useful in the joint analysis of tens of different satDNA families. These were TSI
773 (tandem structure index), RPS (relative peak size) and kurtosis of the repeat landscape
774 distribution as homogenization indices, DIVPEAK as an index of degeneration, CEI as
775 an index of concerted evolution, CTR for consensus turnover rate, and IlibS for
776 incomplete library sorting. However, the main shortcoming of our present analysis was
777 the impossibility to ascertain whether those satDNA families showing no sequence
778 homology between these two species (i.e., non-shared satDNAs) arose *de novo* in one of
779 the species or else they had degenerated in one species but not in the other. To solve this
780 problem, it will be necessary to analyze many species belonging to the same
781 taxonomical group and thus sharing a given satDNA library. We are now sequencing
782 other oedipodine species to perform a multispecies satellitome comparison in the hope
783 that it will allow a better classification of the non-shared satDNA families into *de novo*
784 and partly extinct ones.

785

786 **Conclusions**

787 The analysis of the satellitomes of two species of grasshoppers separated by 22.8 Ma of
788 independent evolution has revealed that one-third of the nearly 60 satDNA families
789 found in each species showed sequence similarity to be considered orthologous and thus
790 descended from their last common ancestor. SatDNA turnover at the level of consensus
791 sequences (CTR) showed a range of variation up to two orders of magnitude among
792 orthologous superfamilies. The use of new satDNA parameters allowing to quantify

793 tandem structure (TSI), homogenization (RPS), degeneration (DIVPEAK), concerted
794 evolution (CEI) and incomplete library sorting (ILibS) showed that satDNA
795 amplification has a dual role by increasing homogenization at intra-specific level and
796 diversification at inter-specific level, thus being a molecular driver unavoidably leading
797 to concerted evolution. Most orthologous pairs of satDNAs analyzed in these species
798 showed the concerted pattern of evolution. The causes for the four non-concerted
799 evolution cases were identified as poor amplification in *O. decorus*. The highest levels
800 of concerted evolution were found for satDNAs displaying long repeat units, high levels
801 of homogenization and FISH bands. These results led us to put forward a general model
802 for satDNA evolution, which updates past models with new empirical data and new
803 statistical approaches to quantify key aspects of variation in satDNA dynamics. We also
804 provide a renewed view of the Library Hypothesis by which a satDNA library begins a
805 new divergence process with each cladogenetic event, during which some satDNA
806 families can disappear whereas other can form *de novo*. The contingent nature of
807 satDNA evolution will make unpredictable the precise set of satDNAs present in each
808 species, some of which will be shared with other species and others will not.

809

810 **Methods**

811 **Materials and sequencing**

812 We collected 21 males of the grasshopper *Oedaleus decorus* in Cortijo Shambala
813 (Sierra Nevada, Granada, Spain; 36.96111 N, 3.33583 W) on 6 July 2015. They were
814 anaesthetized with ethyl-acetate vapours prior to dissection, and testes were fixed in 3:1
815 ethanol-acetic acid and stored at 4°C for subsequent fluorescent in situ hybridization
816 (FISH) analysis. Body remains were immersed in liquid nitrogen and stored at -80 °C
817 for molecular analysis and DNA sequencing. We then extracted genomic DNA from a

818 hind leg from one male, using the GenElute Mammalian Genomic DNA Miniprep kit
819 (Sigma). Next we sent the purified DNA to Macrogen Inc. (South Korea) who built a
820 genomic library with ~180 bp insert size, using the Illumina Truseq nano DNA kit, and
821 sequenced it in an Illumina HiSeq2000 platform (2x101 nt) yielding about 9 Gb of reads.
822 We deposited this library in the Sequence Read Archive (SRA) under accession number
823 SRR9649806 [86].

824 For the *Locusta migratoria* satellitome, we used the results generated in Ruiz-
825 Ruano et al. [22], including some new analyses of the same Illumina libraries obtained
826 from a Spanish individual lacking B chromosomes (SRA library SRR2911427 [87]),
827 satDNA FISH location, and their consensus sequences (GenBank accession numbers
828 KU056702–KU056808). During these new analyses, we detected a previous mistake in
829 the assembly of the LmiSat01A-193 subfamily, consisting of a false tandem duplication
830 of 8 nt in the consensus monomer. We amended this mistake and renamed the (new)
831 sequence as LmiSat01A-185 (GenBank accession number KU056702.2). We thus
832 performed a new analysis of abundance and divergence for the whole satellitome,
833 considering this modification that implied only slight changes.

834 In addition, we generated an Oxford Nanopore library for *L. migratoria* using
835 the MinION system with a flow cell version R9. We constructed the library using 5 µg
836 of DNA without fragmentation step applying the the Nanopore Genomic Kit version
837 SQK-LSK108 and the CleanNGS magnetic beads for washes. After applying the
838 localbase-calling program from Nanopore, we got 63,346 reads summing up 130 Mb
839 (~0.02x of coverage).

840

841 **Bioinformatic and sequence analyses**

842 We characterized the *O. decorus* satellitome applying the satMiner protocol [22].
843 Briefly, this protocol begins with a run of RepeatExplorer [88] and the elimination of
844 homologous reads with Deconseq [89] to perform a new round of RepeatExplorer with
845 the remaining reads. We started with 100,000 read pairs and performed five additional
846 rounds, subsequently duplicating the number of read pairs. Then we identified clusters
847 in each RepeatExplorer round showing spherical or ring-shaped graphs, which are
848 typical for satDNA. We checked the structure of their contigs with a dot-plot using
849 Geneious v4.8.5 [90] to test if they were tandemly repeated, and only those that met this
850 condition were considered as satDNA. Every satDNA family was named with three
851 letters alluding to species name (*L. migratoria* or *O. decorus*) followed by "Sat", a
852 catalogue number (in decreasing order of abundance) and monomer length, following
853 our previous suggestion in Ruiz-Ruano et al. [22]. For instance, the most abundant
854 satDNA families in the two species analyzed here were LmiSat01-185 and OdeSat01-
855 287. The different subfamilies within a same family were alphabetically named with
856 capital letters in order of decreasing abundance.

857 Considering their level of sequence identity, we classified every collection of
858 homologous sequences into subfamilies (identity>95%), families (>80%), and
859 superfamilies (>40%). Next, we randomly selected 5 million read pairs with SeqTK
860 (<https://github.com/lh3/seqtk>) and aligned them against the reference sequences with
861 RepeatMasker v4.0.5 [91]. With these results, we estimated total abundance and
862 average divergence and generated a repeat landscape. Finally, we numbered the satellite
863 families in descending order of abundance. We deposited sequences for satellite DNAs
864 characterized in *O. decorus* in GenBank with accession numbers MT009035-
865 MT009125.

866 We then searched for homology between *L. migratoria* and *O. decorus*
867 satellitomes with the `rm_homolgy` script [22] that makes all-to-all alignments with
868 RepeatMasker [91]. We aligned homologous satellites with Muscle v3.6 [92]
869 implemented in Geneious v4.8.5 [90] and reviewed them manually. Then we generated
870 minimum spanning trees (MST) with Arlequin v3.5 [93] (Excoffier and Lischer 2010)
871 and visualized them with HapStar v0.7 [94]. We used the same alignments to estimate
872 the divergence between satDNA families of *L. migratoria* and *O. decorus*. To estimate a
873 consensus turnover rate (CTR) of satDNA sequences, we performed alignments of
874 consensus sequences using ClustalX [95]. Sequence divergence between species was
875 calculated according to the Kimura two-parameter model (K2P; [96]), using MEGA6
876 [97]. When orthologous satDNA families were composed of several subfamilies, all
877 consensus sequences from each subfamily were aligned and the average of all pairwise
878 distances between the two species was computed. Finally, CTR was calculated using the
879 $CTR = K/2T$ equation, where T = divergence time between species and K = K2P
880 divergence (Kimura 1980). Turnover rates were estimated considering that the *Oedaleus*
881 and *Locusta* genera split 22.81 Ma [45].

882 To get some insights on array length, we analyzed our MinION library obtained
883 from *L. migratoria* gDNA (see above). For this purpose, we performed an alignment of
884 these reads against the consensus sequences of the *L. migratoria* satellitome using
885 RepeatMasker [91]. However, due to the lack of resolution at subfamily level due to the
886 high level of sequencing errors in these long reads, we only performed this analysis only
887 for the most abundant subfamily in each family, i.e, that noted with the letter “A”. We
888 then analyzed the length of all arrays found for each family to recorded the maximum
889 array length (MAL) for subsequent analysis. For this purpose, we only considered
890 arrays showing length higher than 1.5 repeat units, i.e. at least dimers, and the observed

891 figures for MAL in the 56 satDNA families analyzed in *L. migratoria* ranged between
892 62 and 20,180 repeat units. In addition, we considered 3 nt as the maximum inter-array
893 distance to collapse two consecutive TR arrays into a same array, in order to partly
894 counteract the splitting effect of short insertions or deletions due to replication slippage.
895 These calculations were implemented in a custom script
896 (https://github.com/mmarpe/satION/blob/master/dis_bed_max.py).

897

898 **Analysis of tandem structure**

899 We developed a method to estimate the degree of tandem structure in satDNA using a
900 pipeline that we made publicly available throughout repository
901 (<https://github.com/fjruiaruano/SatIntExt>). This method is based on scoring the number
902 of Illumina read pairs containing repeat units for a given satDNA family in the two
903 reads (onwards named "homogeneous read pairs") and the number of read pairs
904 containing such a repeat in only one member of the read pair (onwards named
905 "heterogeneous read pairs"). The proportion of homogeneous read pairs indicates the
906 degree at which a satDNA family is tandemly structured (tandem structure index = TSI).
907 This index underestimates the true value by the equivalent to the half of the number of
908 arrays (since each array has two external units). However, as the number of repeat units
909 is much higher than the number of arrays, we consider that this underestimation may be
910 low at the genomic level. To validate TSI, we analyzed Oxford Nanopore MinION long
911 reads in *L. migratoria*, by annotating all satDNA variants found in them and scoring the
912 number of repeat units constituting the longest array found for each satDNA family.
913 Despite low coverage of the MinION reads, these longest arrays showed significant
914 positive correlation with TSI (Spearman rank correlation: $r_s = 0.42$, $N = 55$, $t = 3.36$, $P =$
915 0.001), indicating that TSI is a valid estimator for the degree of tandem structure of

916 satDNA. In addition, we tried to annotate the external read of every heterogeneous read
917 pair with the database of repetitive elements of *L. migratoria* generated in Ruiz-Ruano
918 et al. [98] with RepeatMasker. Thus, we found homology of the elements adjacent to the
919 satDNA arrays with satDNAs, transposable elements, rDNAs, snDNAs, tRNAs,
920 histones, mitochondrial DNA and unknown elements in some read pairs, and counted
921 the number of occurrences. This analysis is also integrated in the above-mentioned
922 pipeline.

923

924 **Homogenization and degeneration indices**

925 SatDNA homogenization, i.e., the degree of intraspecific similarity between its
926 tandemly structured monomers, is conceptually inverse to average sequence divergence.
927 Therefore, a homogenization index should be negatively correlated with the K2P
928 divergence. Trying to get such an index, we built repeat landscapes for each satDNA
929 subfamily (90 in *O. decorus* and 103 in *L. migratoria*) and searched for divergence
930 peaks, i.e., those divergence values showing the highest abundance in the repeat
931 landscape (DIVPEAK) (Fig. 1). Then, we summed up the abundances of all satDNA
932 sequences at $\pm 2\%$ divergence from the DIVPEAK class to calculate abundance in the
933 5% peak or PEAK-SIZE (Fig. 1). The logic was to get a collection of sequences
934 diverging 5% or less to the consensus sequence, thus coinciding with our criterion to
935 define subfamilies, as they probably derived from the same amplification event (see
936 Ruiz-Ruano et al. [22] for details). Finally, we calculated relative peak size (RPS) as the
937 quotient between PEAK-SIZE and total abundance (see Fig. 1), which measures the
938 proportion of repeat units being part of the last amplification event. To calculate RPS at
939 the family level in those families showing two or more subfamilies, we followed the
940 same procedure including all subfamily satDNA sequences, so that each subfamily

941 weighted in proportion to its abundance. RPS serves as an index of homogenization
942 because it is expected to increase with satDNA amplification, as the new units derived
943 from tandem duplication will initially show identical sequences, thus increasing global
944 identity. DIVPEAK serves as an index of degeneration because it will increase by
945 mutation accumulation and is thus proportional to the time passed since the last
946 amplification. Specifically, DIVPEAK is the value of divergence (from 0% onwards) at
947 which a given satDNA shows its maximum abundance, and increases when mutational
948 decay move its abundance peak away from complete homogenization (divergence=0)
949 where it arrived after its last major amplification event. The values for average
950 divergence, total abundance, maximum abundance, maximum divergence, RPS and
951 DIVPEAK for every satDNA family were estimated from with a custom script using the
952 divsum files from RepeatMasker
953 (https://github.com/fjruiaruano/SatIntExt/blob/main/divsum_stats.py).

954

955 **Concerted evolution index and incomplete library sorting**

956 We calculated the divergence at intra- ($K2P_{intra}$) and inter-specific ($K2P_{inter}$) levels for
957 the 20 pairs of orthologous satDNA families, and calculated an index of concerted
958 evolution (CEI) as \log_2 the $K2P_{inter}/K2P_{intra}$ quotient.

959 The comparative analysis of RLs and MSTs revealed that the observed
960 differences between OSFs in CTR were due to the state of library sorting between
961 species. On this basis, we observed that the OSF showing the highest CTR was that
962 showing a best separation between species for all families and subfamilies of satDNA.
963 We then gave 1 to the sorting state of this OSF and then divided all CTR values by this
964 maxCTR to obtain an index of the relative sorting for each OSF. One minus the

965 obtained value thus indicated the degree of incomplete library sorting (ILibS) for each
966 OSF.

967

968 **Analysis of conserved motifs and curvature**

969 We analyzed the consensus sequences of shared and non-shared satDNAs between the
970 two species looking for functional signatures. We used the ETANDEM, EINVERTED,
971 and PALINDROME programs from the EMBOSS suite of bioinformatics tools [99] for
972 the detection of internal repeats (direct or inverted) and palindromes. Short internal
973 direct repeats indicate the presence of functional motifs within the satDNA repeats.
974 Dyad symmetries, many of them associated with thermodynamically stable secondary
975 structures, are predicted to adopt non-B DNA conformations, such as stem-loops or
976 cruciforms, which might have a role as targets for protein binding. Thus, as an
977 additional test on the propensity to form non-B DNA conformations, we checked all
978 satDNA families using the Mfold web server
979 (<http://www.unafold.org/mfold/applications/rna-folding-form-v2.php>) for nucleic acid
980 folding prediction [100], estimating Gibbs free energy (dG) of the predicted secondary
981 structures [101]. We also checked the consensus sequences of both types of satDNAs
982 for sequence-dependent bendability/curvature propensity of repeats. We produced the
983 bendability/curvature propensity plots with the bend.it server at
984 http://pongor.itk.ppke.hu/dna/bend_it.html#/bendit_intro [102], using the DNase I based
985 bendability parameters of Brukner et al. [103] and the consensus bendability scale [104].
986 Finally, we used the sliding windows option of the DnaSP v.5.10 program [105] for the
987 analysis of nucleotide diversity (π) per position for every shared satDNA in order to
988 detect DNA conserved motifs. For this, we use multiple alignments of several dozens of
989 monomer repeats selected per each satDNA.

990

991 **Chromosomal location of the *O. decorus* satDNAs**

992 To compare the chromosomal location of orthologous satDNA families in these species,
993 we performed fluorescent in situ hybridization (FISH) for 14 satDNA families in *O.*
994 *decorus* which showed sequence homology with 20 families in *L. migratoria*. For this
995 purpose, we designed divergent primers for these 14 satDNA families in *O. decorus*
996 using Primer3 [106] with a $T_m \sim 60$ °C, to generate FISH probes as described in Cabrero
997 et al. [107] and Ruiz-Ruano et al. [22].

998

999 **Statistical analysis**

1000 To investigate distribution fitting of RPS and DIVPEAK, we used the chi-square test,
1001 and the normality of other variable distributions was tested by the Shapiro-Wilks test,
1002 and, when this condition was not met, we used the non-parametric Spearman rank
1003 correlation test. In the case of turnover rate, we performed forward stepwise multiple
1004 regression to analyze its dependence on other variables. In this case, we calculated
1005 variance inflation factors (VIFs) to test for multicollinearity, and the fit of standardized
1006 residuals of this regression to a normal distribution was tested by means of the Shapiro-
1007 Wilks test. All these analyses were performed using the Statistica software (Statsoft
1008 Inc.). Two-group comparisons were performed by the Gardner-Altman estimation plot
1009 method devised by Ho et al. [108] following the design in Gardner and Altman [109], as
1010 implemented in <https://www.estimationstats.com>. This analysis calculates the effect size
1011 by the mean difference between groups, for independent samples, or else by the paired
1012 mean difference in case of paired samples. The effect size is then evaluated by the 95%
1013 confidence interval (95% CI) and whether it includes or not the zero value. Contingency
1014 tests were performed by the RXC program, which employs the Metropolis algorithm to

1015 obtain an unbiased estimate of the exact p-value [110]. In all cases 20 batches of 2,500
1016 replicates were performed.

1017

1018 **Abbreviations**

1019 B-pattern: Banded pattern (pattern in FISH analyses)

1020 CEI: Concerted Evolution Index

1021 CI: Confidence Interval

1022 CTR: Consensus Turnover Rate

1023 dG: Gibbs free energy

1024 DIVPEAK: Divergence Peak

1025 FISH: Fluorescence *In Situ* Hybridization

1026 ILibS: Incomplete Library Sorting

1027 K2P: Kimura Two-Parameter (substitution model)

1028 Lmi: *Locusta migratoria*

1029 NS-pattern: No signal pattern (in FISH analyses)

1030 MAL: Maximum Array Length (observed in MinIon reads of *L. migratoria*)

1031 MST: Minimum Spanning Tree

1032 Ode: *Oedaleus decorus*

1033 OSF: Orthologous Superfamily

1034 RL: Repeat Landscape

1035 RPS: Relative peak size

1036 RUL: Repeat Unit Length

1037 satDNA: satellite DNA

1038 SF: Superfamily

1039 TSI: Tandem Structure Index

1040 VIF: Variance inflation factors

1041

1042 **Declarations**

1043 **Ethics approval and consent to participate**

1044 Not applicable.

1045 **Consent for publication**

1046 Not applicable.

1047 **Availability of data and materials**

1048 The Illumina libraries used for this article are available in the Sequence Read Archive
1049 (SRA) with accession numbers SRR9649806 [86] and SRR2911427 [87]. Main data
1050 generated or analyzed during this study are included in this published article and its
1051 supplementary information files.

1052 **Competing interests**

1053 The authors declare no competing interests.

1054 **Funding**

1055 FJRR was also supported by a postdoctoral fellowship from Sven och Lilly Lawskis
1056 fond (Sweden) and a Marie Skłodowska-Curie Individual Fellowship (grant agreement
1057 875732, European Union).

1058 **Acknowledgments**

1059 Not applicable.

1060 **Authors' contributions**

1061 Conceptualization: JPMC, JC, MDLL, MMP, FP, MAGR, FJRR; experimental
1062 design: JPMC, JC, MDLL, MMP, FP, MAGR, FJRR; sampling: JPMC and JC;
1063 cytogenetic analyses: JPMC, JC, MDLL; data analysis: JPMC, MMP, MAGR, FJRR.
1064 All authors read and approved the manuscript.

1065

1066 **References**

- 1067 1. Kit S. Equilibrium sedimentation in density gradients of DNA preparations from
1068 animal tissues. *J Mol Biol.* 1961;3:711–6.
- 1069 2. Waring M, Britten RJ. Nucleotide sequence repetition: A rapidly reassociating
1070 fraction of mouse DNA. *Science.* 1966;154:791–4.
- 1071 3. Smith GP. Evolution of repeated DNA sequences by unequal crossover. *Science.*
1072 1976;191:528–35.
- 1073 4. Kimura M, Ohta T. Population genetics of multigene family with special reference to
1074 decrease of genetic correlation with distance between gene members on a
1075 chromosome. *Proc Nat Acad Sci USA.* 1979;76:4001–5.
- 1076 5. Ohta T. Genetic variation in small multigene families. *Genet Res.* 1981;37:133–49.
- 1077 6. Ohta T. On the evolution of multigene families. *Theor Popul Biol.* 1983;23:216–40.
- 1078 7. Ohta T, Kimura M. Some calculations on the amount of selfish DNA. *Proc Natl Acad*
1079 *Sci USA.* 1981;78:1129–32.
- 1080 8. Stephan W. Recombination and the evolution of satellite DNA. *Genet Res.*
1081 1986;47:167–74.
- 1082 9. Stephan W. Quantitative variation and chromosomal location of satellite DNAs.
1083 *Genet Res.* 1987;50:41–52.
- 1084 10. Stephan W. Tandem-repetitive non coding DNA: forms and forces. *Mol Biol Evol.*
1085 1989;6:198–212.
- 1086 11. Charlesworth B, Langley CH, Stephan W. The evolution of restricted recombination
1087 and the accumulation of repeated DNA sequences. *Genetics.* 1986;112(4):947–62.
- 1088 12. Garrido-Ramos MA. Satellite DNA: An evolving topic. *Genes.* 2017;8:230.

- 1089 13. Walsh JB. Persistence of tandem arrays: implications for satellite and simple-
1090 sequence DNAs. *Genetics*. 1987;115:553–67.
- 1091 14. Cohen S, Agmon N, Yacobi K, Mislovati M, Segal D. Evidence for rolling circle
1092 replication of tandem genes in *Drosophila*. *Nucleic Acids Res*. 2005;33:4519–26.
- 1093 15. Cohen S, Agmon N, Sobol O, Segal D. Extrachromosomal circles of satellite repeats
1094 and 5S ribosomal DNA in human cells. *Mobile DNA*. 2010;1:11.
- 1095 16. Šatović E, Plohl M. Tandem Repeat-Containing MITEs in the Clam *Donax*
1096 *trunculus*. *Genome Biol Evol*. 2013;5:2549–59.
- 1097 17. Pavlek M, Gelfand Y, Plohl M, Meštrović N. Genome-wide analysis of tandem
1098 repeats in *Tribolium castaneum* genome reveals abundant and highly dynamic
1099 tandem repeat families with satellite DNA features in euchromatic chromosomal
1100 arms. *DNA Res*. 2015;22:387–401.
- 1101 18. Meštrović N, Mravinac B, Pavlek M, Vojvoda-Zeljko T, Šatović E, Plohl M.
1102 Structural and functional liaisons between transposable elements and satellite
1103 DNAs. *Chromosome Res*. 2015;23:583–96.
- 1104 19. Šatović E, Vojvoda Zeljko T, Luchetti A, Mantovani B, Plohl M. Adjacent
1105 sequences disclose potential for intra-genomic dispersal of satellite DNA repeats
1106 and suggest a complex network with transposable elements. *BMC Genom*.
1107 2016;17:997.
- 1108 20. Ruiz-Ruano FJ, Castillo-Martínez J, Cabrero J, Gómez R, Camacho JPM, López-
1109 León MD. High-throughput analysis of satellite DNA in the grasshopper
1110 *Pyrgomorpha conica* reveals abundance of homologous and heterologous higher-
1111 order repeats. *Chromosoma*. 2018;127:323–40.

- 1112 21. Fry K, Salser W. Nucleotide sequences of HS- α satellite DNA from kangaroo rat
1113 *Dipodomys ordii* and characterization of similar sequences in other rodents. Cell.
1114 1977;12:1069–84.
- 1115 22. Ruiz-Ruano FJ, López-León MD, Cabrero J, Camacho JPM. High-throughput
1116 analysis of the satellitome illuminates satellite DNA evolution. Sci Rep.
1117 2016;6:28333.
- 1118 23. Djupedal I, Kos-Braun IC, Mosher RA, Söderholm N, Simmer F, Hardcastle TJ,
1119 Fender A, Heidrich N, Kagansky A, Bayne E, et al. Analysis of small RNA in
1120 fission yeast; centromeric siRNAs are potentially generated through a structured
1121 RNA. EMBO J. 2009;28:3832–44.
- 1122 24. Schueler MG, Swanson W, Thomas PJ. NISC Comparative Sequencing Program &
1123 Green, E.D. Adaptive evolution of foundation kinetochore proteins in primates.
1124 Mol Biol Evol. 2010;27:1585–97.
- 1125 25. Fachinetti D, Han JS, McMahon MA, Ly P, Abdullah A, Wong AJ, Cleveland DW.
1126 DNA sequence-specific binding of CENP-B enhances the fidelity of human
1127 centromere function. Dev Cell. 2015;33:314–27.
- 1128 26. Arnason U, Grettarsdottir S, Widegren B. Mysticete (baleen whale) relationships
1129 based upon the sequence of the common cetacean DNA satellite. Mol Biol Evol.
1130 1992;9:1018–28.
- 1131 27. Garrido-Ramos MA, Jamilena M, Lozano R, Ruiz Rejón C, Ruiz Rejón M. The
1132 EcoRI centromeric satellite DNA of the Sparidae family (Pisces, Perciformes)
1133 contains a sequence motive common to other vertebrate centromeric satellite
1134 DNAs. Cytogenet. Cell Genet. 1995;71:345–51.
- 1135 28. Garrido-Ramos MA, de la Herrán R, Jamilena M, Lozano R, Ruiz Rejón C, Ruiz
1136 Rejón M. Evolution of centromeric satellite-DNA and its use in phylogenetic

- 1137 studies of the Sparidae family (Pisces, Perciformes). *Mol Phyl Evol.* 1999;12:200–
1138 4.
- 1139 29. de la Herrán R, Fontana F, Lanfredi M, Congiu L, Leis M, Rossi R, Ruiz Rejón C,
1140 Ruiz Rejón M, Garrido-Ramos MA.. Slow rates of evolution and sequence
1141 homogenization in an ancient satellite DNA family of sturgeons. *Molecular*
1142 *Biology and Evolution.* 2001;18:432–36.
- 1143 30. de La Herrán R, Ruiz Rejón C, Ruiz Rejón M, Garrido-Ramos MA. The molecular
1144 phylogeny of the Sparidae (Pisces, Perciformes) based on two satellite DNA
1145 families. *Heredity.* 2001;87:691–7.
- 1146 31. Mravinac B, Plohl M, Meštrović N, Ugarković Đ. Sequence of PRAT satellite DNA
1147 “frozen” in some Coleopteran species. *J Mol Evol.* 2002;54:774–83.
- 1148 32. Mravinac B, Plohl M, Ugarković Đ. Preservation and high sequence conservation of
1149 satellite DNAs suggest functional constraints. *J Mol Evol.* 2005;61:542–50.
- 1150 33. Robles F, de la Herrán R, Ludwig A, Ruiz Rejón C, Ruiz Rejón M, Garrido-Ramos
1151 MA. Evolution of ancient satellite DNAs in sturgeon genomes. *Gene.*
1152 2004;338:133–42.
- 1153 34. Cafasso D, Chinali G. An ancient satellite DNA has maintained repetitive units of
1154 the original structure in most species of the living fossil plant genus *Zamia*.
1155 *Genome.* 2014;57:125–35.
- 1156 35. Chaves, R., Ferreira, D., Mendes-Da-Silva, A., Meles, S., & Adegá, F. FA-SAT is
1157 an old satellite DNA frozen in several bilateria genomes. *Genome Biol Evol.*
1158 2017;9:3073–87.
- 1159 36. Harding RM, Boyce AJ, Clegg JB. The evolution of tandemly repetitive DNA:
1160 recombination rules. *Genetics.* 1992;132:847–59.

- 1161 37. Dover G. Molecular drive: a cohesive mode of species evolution. *Nature*.
1162 1982;299:111–7.
- 1163 38. Dover G. A molecular drive through evolution. *Bioscience*. 1982;32:526–33.
- 1164 39. Nei M, Rooney AP. Concerted and birth-and-death evolution of multigene families.
1165 *Annu Rev Genet*. 2005;39:121–52.
- 1166 40. Eirín-López JM, Rebordinos L, Rooney AP, Rozas J. The birth- and-death evolution
1167 of multigene families revisited. In: Garrido-Ramos MA, editor. *Repetitive DNA*.
1168 Basel: S. Karger AG. 2012. p. 170–96.
- 1169 41. Kuhn GCS, Küttler H, Moreira-Filho O, Heslop-Harrison JS. The 1.688 repetitive
1170 DNA of drosophila: Concerted evolution at different genomic scales and
1171 association with genes. *Mol Biol Evol*. 2012;29:7–11.
- 1172 42. Brajković J, Feliciello I, Bruvo-Madžarić B, Ugarković Đ. Satellite DNA-like
1173 elements associated with genes within euchromatin of the beetle *Tribolium*
1174 *castaneum*. *G3-Genes Genom Genet*. 2012;2:93141.
- 1175 43. Feliciello I, Akrap I, Brajković J, Zlatar I, Ugarković Đ. Satellite DNA as a driver of
1176 population divergence in the red flour beetle *Tribolium castaneum*. *Genome Biol*
1177 *Evol*. 2015;7:228–39.
- 1178 44. Garrido-Ramos MA. Satellite DNA in Plants: More than Just Rubbish. *Cytogenet*
1179 *Genome Res*. 2015;146:153–70.
- 1180 45. Song H, Amédégno C, Cigliano MM, Desutter-Grandcolas L, Heads SW, Huang
1181 Y, Otte D, Whiting MF. 300 million years of diversification: elucidating the
1182 patterns of orthopteran evolution based on comprehensive taxon and gene
1183 sampling. *Cladistics*. 2015;31:621–51.
- 1184 46. Blount ZD, Lenski RE, Losos JB. Contingency and determinism in evolution:
1185 Replaying life’s tape. *Science*. 2018;362:6415.

- 1186 47. Utsunomia R, Ruiz-Ruano FJ, Silva DMZA, Serrano EA, Rosa IF, Scudeler PES,
1187 Hashimoto DT, Oliveira C, Camacho JPM, Foresti F. A glimpse into the satellite
1188 DNA library in Characidae fish (Teleostei, Characiformes). *Front Genet.*
1189 2017;8:103.
- 1190 48. Palacios-Gimenez OM, Milani D, Song H, Marti DA, López-León MD, Ruiz-Ruano
1191 FJ, Camacho JPM, Cabral-de-Mello DC. Eight million years of satellite DNA
1192 evolution in grasshoppers of the genus *Schistocerca* illuminate the ins and outs of
1193 the library hypothesis. *Genome Biol Evol.* 2020;12:88–102.
- 1194 49. Ávila Robledillo L, Neumann P, Koblížková A, Novák P, Vrbová I, Macas J.
1195 Extraordinary Sequence Diversity and Promiscuity of Centromeric Satellites in the
1196 Legume Tribe Fabeae. *Mol Biol Evol.* 2020;37:2341–56.
- 1197 50. Zeni dos Santos R, Milan Calegari R, Silva DMZA, Ruiz-Ruano FJ, Melo S,
1198 Oliveira C, Foresti F, Uliano-Silva M, Porto-Foresti F, Utsunomia R. A long-term
1199 conserved satellite DNA that remains unexpanded in several genomes of
1200 Characiformes fish is actively transcribed. *Genome Biol Evol.* 2021;13:evab002.
- 1201 51. Goodsell DS, Dickerson RE. Bending and curvature calculations in B-DNA.
1202 *Nucleic Acids Res.* 1994;22:5497–503.
- 1203 52. Gabrielian A, Simoncsits A, Pongor S. Distribution of bending propensity in DNA
1204 sequences. *FEBS letters.* 1996;393:124–30.
- 1205 53. Plohl M, Meštrovic N, Mravinac B. Satellite DNA evolution. In: Garrido-Ramos
1206 MA, editor. *Repetitive DNA*. Basel: S. Karger AG; 2012. p. 126–52.
- 1207 54. Stephan W, Cho S. Possible role of natural selection in the formation of tandem-
1208 repetitive noncoding DNA. *Genetics.* 1994;136:333–41.

- 1209 55. Navajas-Perez R, de la Herrán R, Jamilena M, Lozano R, Ruiz Rejon C, Ruiz Rejon
1210 M, Garrido-Ramos MA. Reduced rates of sequence evolution of Y-linked satellite
1211 DNA in *Rumex* (Polygonaceae). *J Mol Evol.* 2005;60:391–9.
- 1212 56. Ruiz-Ruano FJ, Cabrero J, López-León MD, Camacho JPM. Satellite DNA content
1213 illuminates the ancestry of a supernumerary (B) chromosome. *Chromosoma.*
1214 2017;126:487–500.
- 1215 57. Melters DP, Bradnam KR, Young HA, Telis N, May MR, Ruby JG, Sebra R, Peluso
1216 P, Eid J, Rank D, Garcia JF. *Genome Biol.* 2013;14:R10.
- 1217 58. Zhang H, Koblížková A, Wang K, Gong Z, Oliveira L, Torres GA, Wu YF, Zhang
1218 W, Novák P, Buell CR, Macas J, Jiang J. Boom-bust turnovers of megabase-sized
1219 centromeric DNA in *Solanum* species: rapid evolution of DNA sequences
1220 associated with centromeres. *Plant Cell.* 2014;26:1436–47.
- 1221 59. Shang WH, Hori T, Toyoda A, Kato J, Pependorf K, Sakakibara Y, Fujiyama A,
1222 Fukagawa T. Chickens possess centromeres with both extended tandem repeats
1223 and short non-tandem-repetitive sequences. *Genome Res.* 2010;20:1219–1228.
- 1224 60. Iwata A, Tek AL, Richard MM, Abernathy B, Fonsêca A, Schmutz J, Chen NWG,
1225 Thareau V, Magdelenat G, Li Y, Murata M, Pedrosa-Harand A, Geffroy V, Nagaki
1226 K, Jackson SA. Identification and characterization of functional centromeres of the
1227 common bean. *Plant J.* 2013;76:47–60.
- 1228 61. Neumann P, Navrátilová A, Schroeder-Reiter E, Koblížková A, Steinbauerová V,
1229 Chocholová E, Novák P, Wanner G, Macas J. Stretching the rules: monocentric
1230 chromosomes with multiple centromere domains. *PLoS Genet.* 2012;8:e1002777.
- 1231 62. Masumoto H, Masukata H, Muro Y, Nozaki N, Okazaki T. A human centromere
1232 antigen (CENP-B) interacts with a short specific sequence in alphoid DNA, a
1233 human centromeric satellite. *J Cell Biol.* 1989;109:1963–73.

- 1234 63. Masumoto H, Nakano M, Ohzeki J. The role of CENP-B and alpha-satellite DNA:
1235 De novo assembly and epigenetic maintenance of human centromeres.
1236 Chromosome Res. 2004;12:543–56.
- 1237 64. Muro Y, Masumoto H, Yoda K, Nozaki N, Ohashi M, Okazaki T. Centromere
1238 protein B assembles human centromeric alpha-satellite DNA at the 17-bp sequence,
1239 CENP-B box. J Cell Biol. 1992;116:585–96.
- 1240 65. Haaf T, Mater AG, Wienberg J, Ward DC. Presence and abundance of CENP-B box
1241 sequences in great ape subsets of primate-specific alpha-satellite DNA. J Mol Evol.
1242 1995;41:487–491.
- 1243 66. Koch J. Neocentromeres and alpha satellite: a proposed structural code for
1244 functional human centromere DNA. Hum Mol Genet. 2000;9:149–54.
- 1245 67. Hall SE, Kettler G, Preuss D. Centromere satellites from Arabidopsis populations:
1246 maintenance of conserved and variable domains. Genome Res. 2003;13:195–205.
- 1247 68. Luchetti A, Cesari M, Carrara G, Cavicchi S, Passamonti M, Scali V, Mantovani B.
1248 Unisexuality and molecular drive: Bag320 sequence diversity in *Bacillus* taxa
1249 (Insecta Phasmatodea). J Mol Evol. 2003;56:587–96.
- 1250 69. Pezer Ž, Brajković J, Feliciello I, Ugarković Đ. Satellite DNA-mediated effects on
1251 genome regulation. In: Garrido-Ramos MA, editor. Repetitive DNA. Basel: S.
1252 Karger AG. 2012. p. 153–69.
- 1253 70. Kasinathan S, Henikoff S. Non-B-form DNA is enriched at centromeres. Mol Biol
1254 Evol. 2018;35:949–62.
- 1255 71. Talbert PB, Henikoff S. Transcribing centromeres: noncoding RNAs and
1256 kinetochore assembly. Trends Genet. 2018;34:587–99.

- 1257 72. Navajas-Pérez R, Schwarzacher T, de la Herrán R, Ruiz Rejón C, Ruiz Rejón M,
1258 Garrido-Ramos MA. The origin and evolution of the variability in a Y-specific
1259 satellite-DNA of *Rumex acetosa* and its relatives. *Gene*. 2006;368:61–71.
- 1260 73. Navajas-Pérez R, Quesada del Bosque ME, Garrido-Ramos MA. Effect of location,
1261 organization, and repeat-copy number in satellite-DNA evolution. *Mol Genet*
1262 *Genom*. 2009;282(4):395–406.
- 1263 74. Suárez-Santiago VN, Blanca G, Ruiz-Rejón M, Garrido-Ramos MA. Satellite-DNA
1264 evolutionary patterns under a complex evolutionary scenario: The case of
1265 *Acrolophus* subgroup (*Centaurea* L., Compositae) from the western Mediterranean.
1266 *Gene*. 2007;404:80–92.
- 1267 75. Quesada del Bosque ME, López-Flores I, Suárez-Santiago VN, Garrido-Ramos MA.
1268 Differential spreading of *Hinfl* satellite DNA variants during radiation in
1269 Centaureinae. *Ann Bot*. 2013;112:1793–1802.
- 1270 76. Quesada del Bosque ME, López-Flores I, Suárez-Santiago VN, Garrido-Ramos MA.
1271 Satellite-DNA diversification and the evolution of major lineages in Cardueae
1272 (Carduoideae Asteraceae). *J Plant Res*. 2014;127:575–83.
- 1273 77. Luchetti A, Marini M, Mantovani B. Non-concerted evolution of the RET76
1274 satellite DNA family in *Reticulitermes* taxa (Insecta, Isoptera). *Genetica*.
1275 2006;128:123–132.
- 1276 78. Lorite P, Muñoz-López M, Carrillo JA, Sanllorente O, Vela J, Mora P, Tinaut A,
1277 Torres MI, Palomeque T. Concerted evolution, a slow process for ant satellite
1278 DNA: study of the satellite DNA in the *Aphaenogaster* genus (Hymenoptera,
1279 Formicidae). *Org Divers Evol*. 2017;17:595–606.

- 1280 79. Willard HF, Waye JS. Chromosome-specific subsets of human alpha satellite DNA:
1281 analysis of sequence divergence within and between chromosomal subsets and
1282 evidence for an ancestral pentameric repeat. *J Mol Evol.* 1987;25:207–14.
- 1283 80. Warburton PE, Willard HF. Genomic analysis of sequence variation in tandemly
1284 repeated DNA: evidence for localized homogeneous sequence domains within
1285 arrays of α -satellite DNA. *J Mol Biol.* 1990;216:3–16.
- 1286 81. Feliciello I, Picariello O, Chinali G. Intra-specific variability and unusual
1287 organization of the repetitive units in a satellite DNA from *Rana dalmatina*:
1288 molecular evidence of a new mechanism of DNA repair acting on satellite DNA.
1289 *Gene.* 2006;383:81–92.
- 1290 82. Lynch M. Statistical inference on the mechanisms of genome evolution. *PLoS Genet.*
1291 2011;7:1–4.
- 1292 83. Lynch M, Bobay L-M, Catania F, Gout J-F, Rho M. The Repatterning of Eukaryotic
1293 Genomes by Random Genetic Drift. *Annu Rev Genomics Hum Genet.*
1294 2011;12:347–66.
- 1295 84. Escudeiro A, Adegas F, Robinson TJ, Heslop-Harrison JS, Chaves R. Conservation,
1296 divergence and functions of centromeric satellite DNA families in the Bovidae.
1297 *Genome Biol Evol.* 2019;11:1152–65.
- 1298 85. Gould SJ. *Wonderful life: the Burgess Shale and the nature of history.* Norton, New
1299 York; 1989.
- 1300 86. Camacho JPM; Cabrero J; López-León MD; Martín-Peciña M; Perfectti F; Garrido-
1301 Ramos MA; Ruiz-Ruano FJ. *Oedaleus decorus* genomic library. 2020.
1302 <https://www.ncbi.nlm.nih.gov/sra/?term=SRR9649806>.

- 1303 87. Ruiz-Ruano FJ; López-León MD; Cabrero J; Camacho JPM. *Locusta migratoria* 0B
1304 gDNA Southern Lineage. 2016.
1305 <https://www.ncbi.nlm.nih.gov/sra/?term=SRR2911427>.
- 1306 88. Novák P, Neumann P, Pech J, Steinhaisl J, Macas J. RepeatExplorer: a Galaxy-
1307 based web server for genome-wide characterization of eukaryotic repetitive
1308 elements from next-generation sequence reads. *Bioinformatics*. 2013;29:792–3.
- 1309 89. Schmieder R, Edwards R. Fast identification and removal of sequence
1310 contamination from genomic and metagenomic datasets. *PLoS One*.
1311 2011;6:e17288.
- 1312 90. Drummond AJ, Ashton B, Buxton S, Cheung M, Cooper A, Heled J, Kearse M,
1313 Moir R, Stones-Havas S, Sturrock S, Thierer T, Wilson A. Geneious v. 4.8.
1314 Auckland, New Zealand: Biomatters Ltd. 2010.
- 1315 91. Smit AFA, Hubley R, Green P (2013) RepeatMasker Open-4.0.
1316 <http://www.repeatmasker.org>
- 1317 92. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high
1318 throughput. *Nucleic Acids Res*. 2004;32:1792–7.
- 1319 93. Excoffier L, Lischer HE. Arlequin suite ver 3.5: a new series of programs to
1320 perform population genetics analyses under Linux and windows. *Mol Ecol Resour*.
1321 2010;10:564–567.
- 1322 94. Teacher AGF, Griffiths DJ. HapStar: automated haplotype network layout and
1323 visualization. *Mol Ecol Resour*. 2011;11:151–3.
- 1324 95. Thompson JD, Gibson TJ, Plewniak F, Jeanmougin F, Higgins DG. The
1325 CLUSTAL_X windows interface: flexible strategies for multiple sequence
1326 alignment aided by quality analysis tools. *Nucleic Acids Res*. 1997;25:4876–82.

- 1327 96. Kimura M. A simple method for estimating evolutionary rates of base substitutions
1328 through comparative studies of nucleotide sequences. *J Mol Evol.* 1980;16:111–20.
- 1329 97. Tamura K, Stecher G, Peterson D, Filipinski A, Kumar S. MEGA6: molecular
1330 evolutionary genetics analysis version 6.0. *Mol Biol Evol.* 2013;30:2725–9.
- 1331 98. Ruiz-Ruano FJ, Cabrero J, López-León MD, Sánchez A, Camacho JPM.
1332 Quantitative sequence characterization for repetitive DNA content in the
1333 supernumerary chromosome of the migratory locust. *Chromosoma.* 2018;127:45–
1334 57.
- 1335 99. Rice P, Longden I, Bleasby A. EMBOSS: the European molecular biology open
1336 software suite. *Trends Genet.* 2000;16:276–7.
- 1337 100. Zuker M. Mfold web server for nucleic acid folding and hybridization prediction.
1338 *Nucleic Acids Res.* 2003;31:3406–15.
- 1339 101. SantaLucia Jr J. A unified view of polymer, dumbbell, and oligonucleotide DNA
1340 nearest-neighbor thermodynamics. *Proc Nat Acad Sci USA.* 1998;95:1460–5.
- 1341 102. Vlahovicek K, Kajan L, Pongor S. DNA analysis servers: plot. it, bend. it, model.
1342 it and IS. *Nucleic Acids Res.* 2003;31:3686–7.
- 1343 103. Brukner I, Sanchez R, Suck D, Pongor S. Sequence-dependent bending propensity
1344 of DNA as revealed by DNase I: parameters for trinucleotides. *The EMBO journal.*
1345 1995;14:1812–18.
- 1346 104. Gabrielian A, Pongor S. Correlation of intrinsic DNA curvature with DNA
1347 property periodicity. *FEBS letters.* 1996;393:65–8.
- 1348 105. Librado P, Rozas J. DnaSP v5: a software for comprehensive analysis of DNA
1349 polymorphism data. *Bioinformatics.* 2009;25:1451–2.
- 1350 106. Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen
1351 SG. Primer3—new capabilities and interfaces. *Nucleic Acids Res.* 2012;40:e115

- 1352 107. Cabrero J, Bakkali M, Bugrov A, Warchalowska-Sliwa E, López-León MD,
1353 Perfectti F, Camacho JPM. Multiregional origin of B chromosomes in the
1354 grasshopper *Eyprepocnemis plorans*. *Chromosoma*. 2003;112:207–11.
- 1355 108. Ho J, Tumkaya T, Aryal S, Choi H, Claridge-Chang A. Moving beyond P values:
1356 data analysis with estimation graphics. *Nat Methods*. 2019;16:565–6.
- 1357 109. Gardner MJ, Altman DG. Confidence intervals rather than P values: estimation
1358 rather than hypothesis testing. *Br Med J*. 1986;292:746–50.
- 1359 110. Rousset F, Raymond M. Testing heterozygote excess and deficiency. *Genetics*
1360 1995;140:1413–1419.

1361

1362 **Supplementary Information**

1363 ***Additional file 1 (.xls format):** Tables S1-S4.

1364 **Table S1.** Molecular and cytological properties of the satellitomes in *Oedaleus decorus*
1365 (*Ode*) and *Locusta migratoria* (*Lmig*). Note that telomeric DNA was also numbered in
1366 both species (no. 13 and 7, respectively) but are omitted here because they were not
1367 considered for this paper analyses. RUL= Repeat unit length. TSI= Tandem structure
1368 index. SF= Superfamily. RPS= Relative peak size. DIVPEAK= Divergence peak.
1369 MAL= Maximum array length observed in MinIon reads of *L. migratoria*. FISH= FISH
1370 pattern (B= banded, NS= No signal). Local= Localization (p= proximal, i= interstitial,
1371 d= distal). Motifs= Conserved motifs in the DNA sequence (0= Yes, 1= No).
1372 Curvature= Propensity to acquire stable structures (0= Yes, 1= No). dG= Gibbs free
1373 energy of the predicted secondary structure.

1374 **Table S2.** Homology between satDNA families found in *O. decorus* and *L. migratoria*.
1375 OSF= Orthologous superfamily. Those families chosen for comparisons between
1376 orthologous pairs are noted in bold-type letter.

1377 **Table S3.** Total number of external reads for each satellite family in *O. decorus* (Ode)
1378 and *L. migratoria* (Lmig) and its annotation. TSI= Tandem Structure Index.
1379 **Table S4.** Characteristics of the orthologous satDNA families analyzed in *O. decorus*
1380 (14) and *L. migratoria* (20). Each row includes one Ode and one Lmi satDNA families
1381 showing homology. Note that some Ode families showed homology with two or three
1382 Lmi ones. OSF= Orthologous superfamily, sf= number of subfamilies, SF= superfamily
1383 name, FISH= FISH pattern (B= banded, NS= no signal), RUL=Repeat unit length (bp),
1384 A+T= % A+T content, abun= abundance (% of the genome), div= divergence (%),
1385 peak_size= abundance of the 5% divergence classes around DIVPEAK, RPS= Relative
1386 peak size, DP= DIVPEAK, kur= kurtosis of repeat landscape distribution, TSI= Tandem
1387 structure index, dG= Free energy of repeat unit sequence, MAL= Maximum array
1388 length observed in MinIon reads of *L. migratoria*, CEI= Concerted evolution index (L=
1389 *L. migratoria*, O= *O. decorus*), Intid= Interspecific sequence identity (%), Intdiv=
1390 Interspecific divergence, CTR= Consensus turnover rate, ILibS= Incomplete library
1391 sorting. Negative CEI values and Int_id>95% are remarked in bold type letter.

1392

1393 ***Additional file 2 (.tif format):** Figure S1.

1394 **Figure S1.** Repeat landscape (RL) and minimum spanning tree (MST) of three
1395 orthologous superfamilies of satellite DNA in *O. decorus* and *L. migratoria* (OSF04,
1396 OSF05 and OSF07). a) RLs showed that OSF04 showed large peaks of amplification in
1397 both species but CTR values ranged between 1.16 and 1.6, presumably due to the
1398 incomplete library sorting (ILibS) evidenced by the MST (note how OdeSat32A and
1399 LmiSat51A connect with both species' sequences). b) OSF05 showed high CTR values,
1400 large amplification peaks in both species and ILibS for only OdeSat22C, which was the
1401 only sequence connected with sequences from both species. c) OSF07 showed the

1402 lowest CTR values and showed very small amplification peaks for OdeSat58 (green
1403 curves in the RL on the left) and higher ILibS, with three sequences being connected
1404 with both species' sequences (LmiSat45-274, LmiSat28A-263 and OdeSat58A-265).

1405

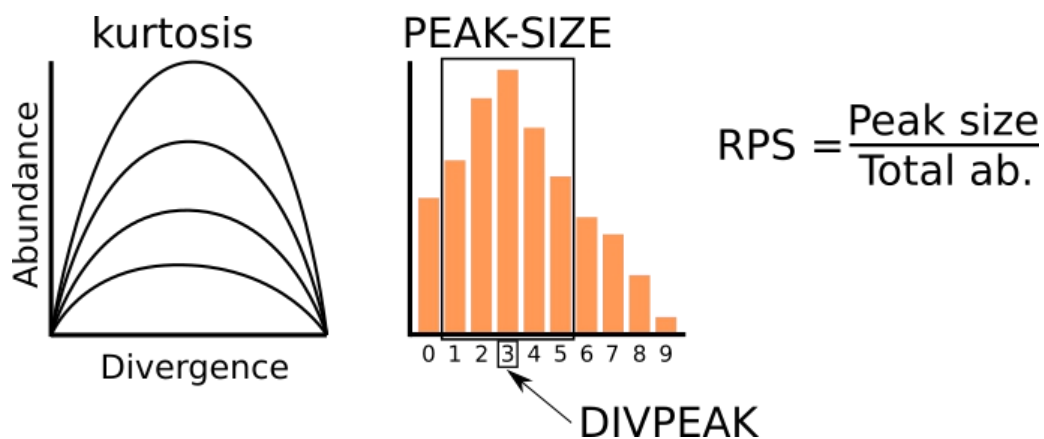
1406 * **Additional file 3 (.xls format):** Dataset S1.

1407 **Dataset 1a.** Data from the *Oedaleus decorus* repeat landscape indicating genomic
1408 abundance for each satellite DNA family and divergence interval.

1409 **Dataset 1b.** Data from the *Locusta migratoria* repeat landscape indicating genomic
1410 abundance for each satellite DNA family and divergence interval.

1411

1412 **Figures**



1414 **Figure 1.** Definition of satDNA parameters in respect to abundance and divergence.

1415 The distribution of the abundances of groups of sequences differing by 1% divergence

1416 constitutes a repeat landscape (RL). It may be seen as a curve (left) or an histogram

1417 (right). In addition of variation in kurtosis, represented by several curves on the left,

1418 three properties of satDNA can be defined on RLs: DIVPEAK is the divergence class

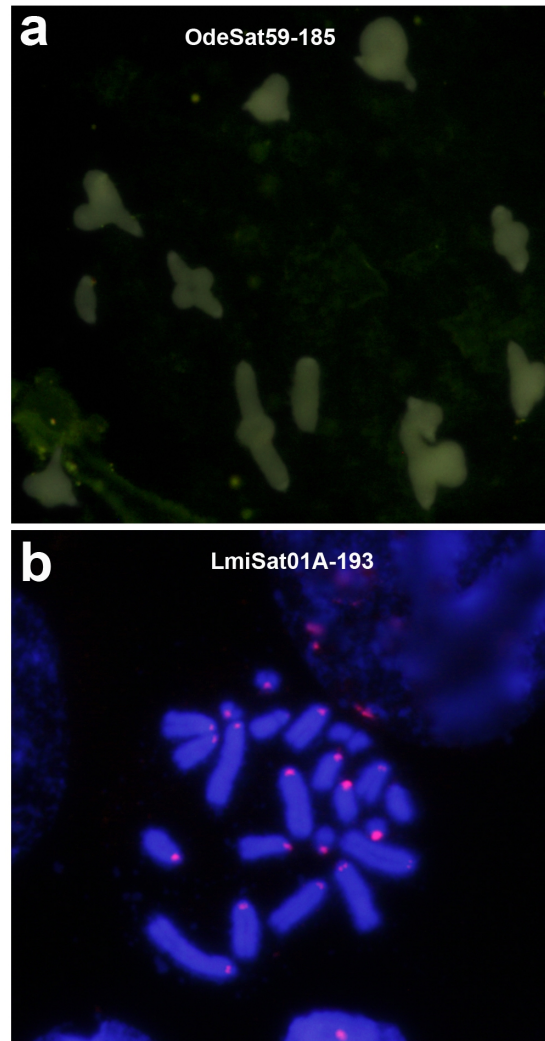
1419 showing the highest abundance (3% in the histogram); PEAK-SIZE is the sum of the

1420 abundances of the five classes included around DIVPEAK, thus constituting the sum of

1421 all sequences differing by less than 5%, thus coinciding with our definition of satDNA

1422 subfamily; RPS is the relative peak size and represents the fraction of abundance which
1423 is included in the 5% amplification peak.

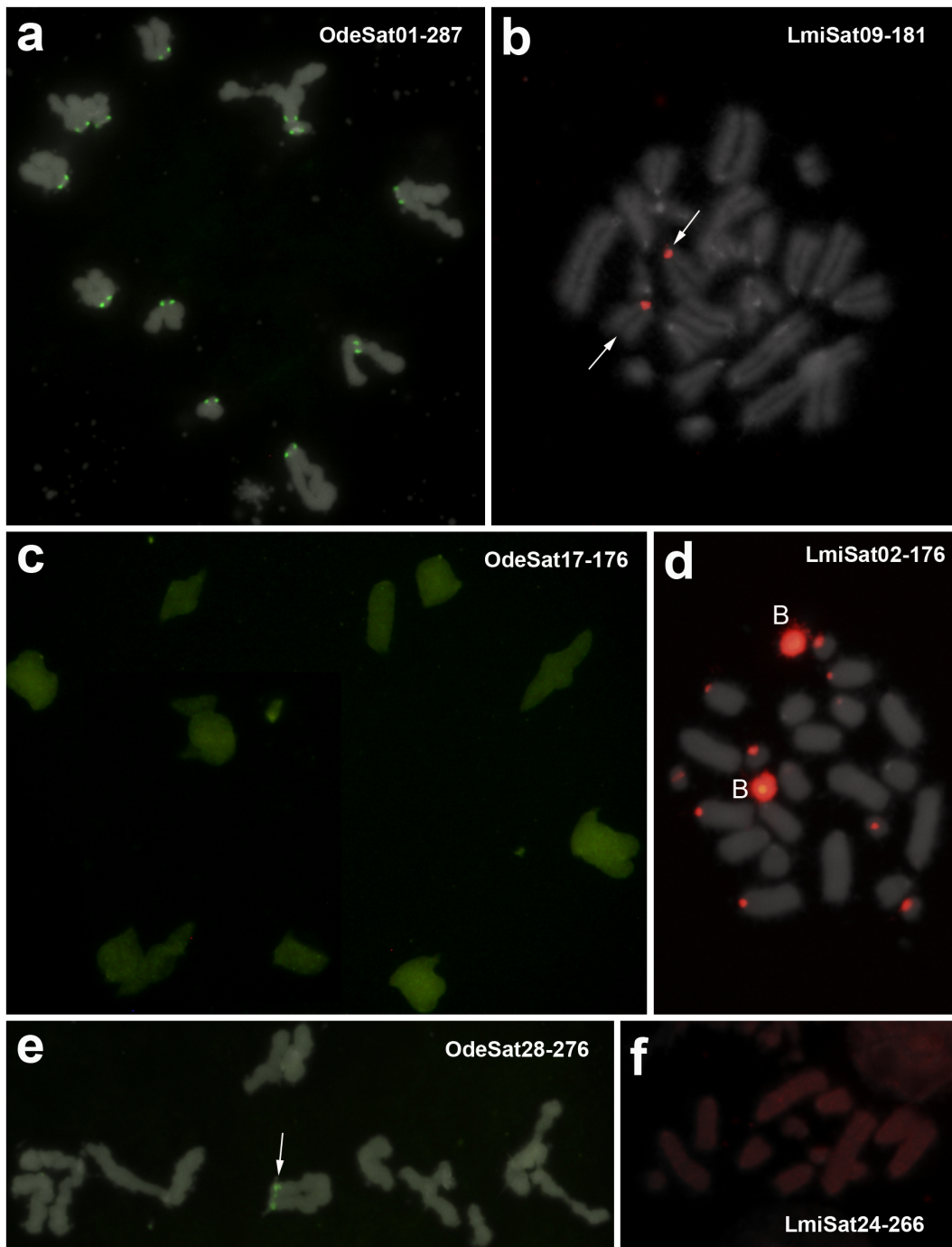
1424



1425

1426 **Figure 2.** FISH analysis of a pair of orthologous families, belonging to OSF12, in
1427 *O.decorus* (a) and *L. migratoria* (b). a) OdeSat59-185 showed no FISH bands on this
1428 meiotic metaphase I cell, thus showing the NS pattern. b) LmiSat01A-193 showed
1429 conspicuous pericentromeric FISH bands on most chromosomes of this embryonic mitotic
1430 metaphase cell, thus showing the B-pattern.

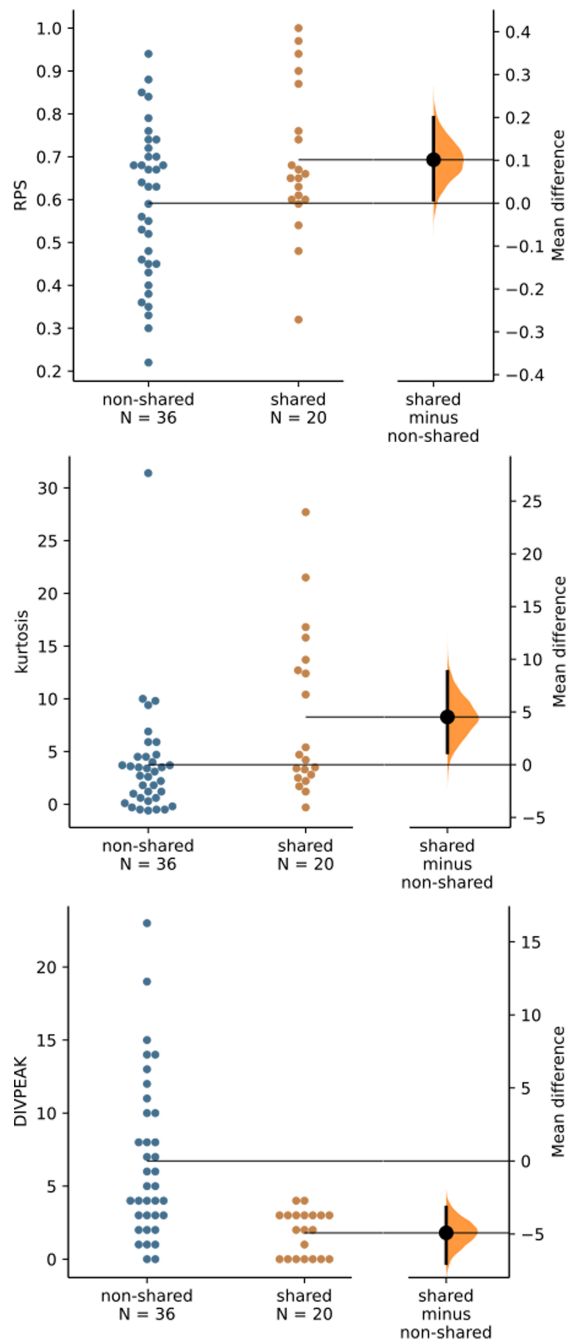
1431



1432

1433 **Figure 3.** FISH analysis of three pairs of orthologous families in *O.decorus* and *L.*
1434 *migratoria*: showing the B-pattern in both species for OSF1 (a and b), the NS and B
1435 patterns, respectively, for OSF3 (c and d), and the B and NS patterns, respectively,
1436 for OSF7 (e and f) (see also Table 2 for satDNA classification into OSFs). a) Presence
1437 of pericentromeric FISH bands for OdeSat01-287 on all chromosomes of this meiotic
1438 metaphase II cell of *O. decorus*. b) Note the presence of its orthologous family

1439 (LmiSat09-181) on a single chromosome pair of this embryo mitotic metaphase cell of
1440 *L. migratoria*. c) Absence of FISH bands for OdeSat17-176 in a meiotic metaphase I
1441 cell of *O. decorus*. d) Presence of its orthologous LmiSat02-176 on pericentromeric
1442 regions of several chromosome pairs and on whole B chromosome length (B) of this
1443 embryo mitotic metaphase cell of *L. migratoria*. e) Presence of a pericentromeric FISH
1444 band on a single chromosome of the haploid set shown in this meiotic metaphase II cell
1445 of *O. decorus*. f) Absence of FISH bands for LmiSat24-266 on the haploid chromosome
1446 set shown in this embryo mitotic metaphase cell of *L. migratoria*.
1447



1448

1449 **Figure 4.** Repeat landscape (RL) and minimum spanning tree (MST) of two
1450 orthologous superfamilies of satellite DNA in *O. decorus* and *L. migratoria* (OSF02 and
1451 OSF12). a) OSF02 showed the highest consensus turnover rate (CTR= 2.86) found
1452 among the 20 values estimated between orthologous pairs of families in both species.
1453 Note that OSF02 showed large amplification peaks in both species (green curve in *O.*
1454 *decorus* and red curve in *L. migratoria*) and that the MST showed complete separation
1455 of OdeSat02 and LmiSat03 sequences. b) OSF12 showed the lowest CTR estimate (0.26

1456 between OdeSat59 and LmiSat01) and the MST (on the right) reveals that the consensus
1457 DNA sequences of these two satDNA families showed only two differences. Also note
1458 in the RL (on the left) that the OdeSat59 curve is very close to zero, as this is the
1459 satDNA family in *O. decorus* showing the lowest abundance, indicating that OSF12 is
1460 represented in this species as relict remains which, by chance, almost coincide in
1461 consensus sequence with the most abundant subfamily in *L. migratoria* (LmiSat01A),
1462 thus evidencing extreme incomplete lineage sorting (see other cases in Additional file 2:
1463 Fig. S1).

1464

1465

1466

1467

1468

1469

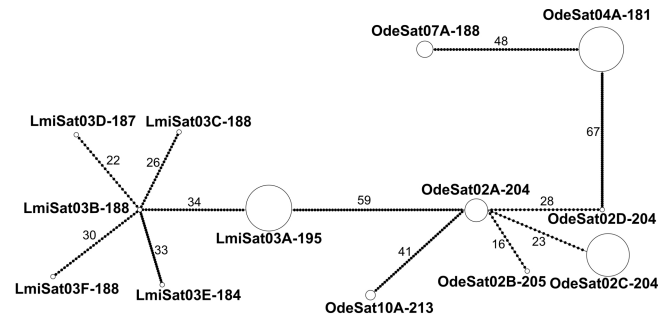
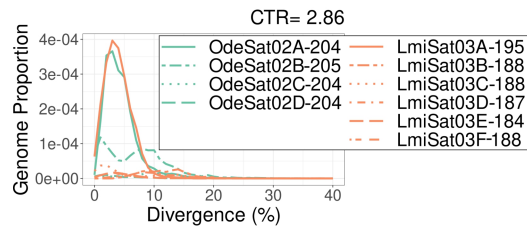
1470

1471

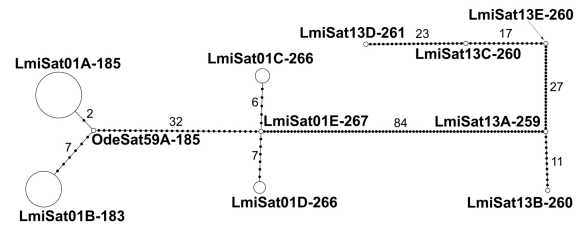
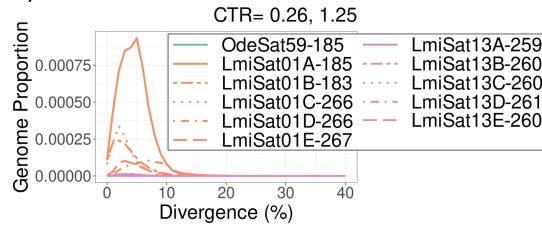
1472

1473

a) OSF02



b) OSF12



1474

1475 **Figure 5.** Gardner-Altman plots comparing RPS, kurtosis and DIVPEAK between

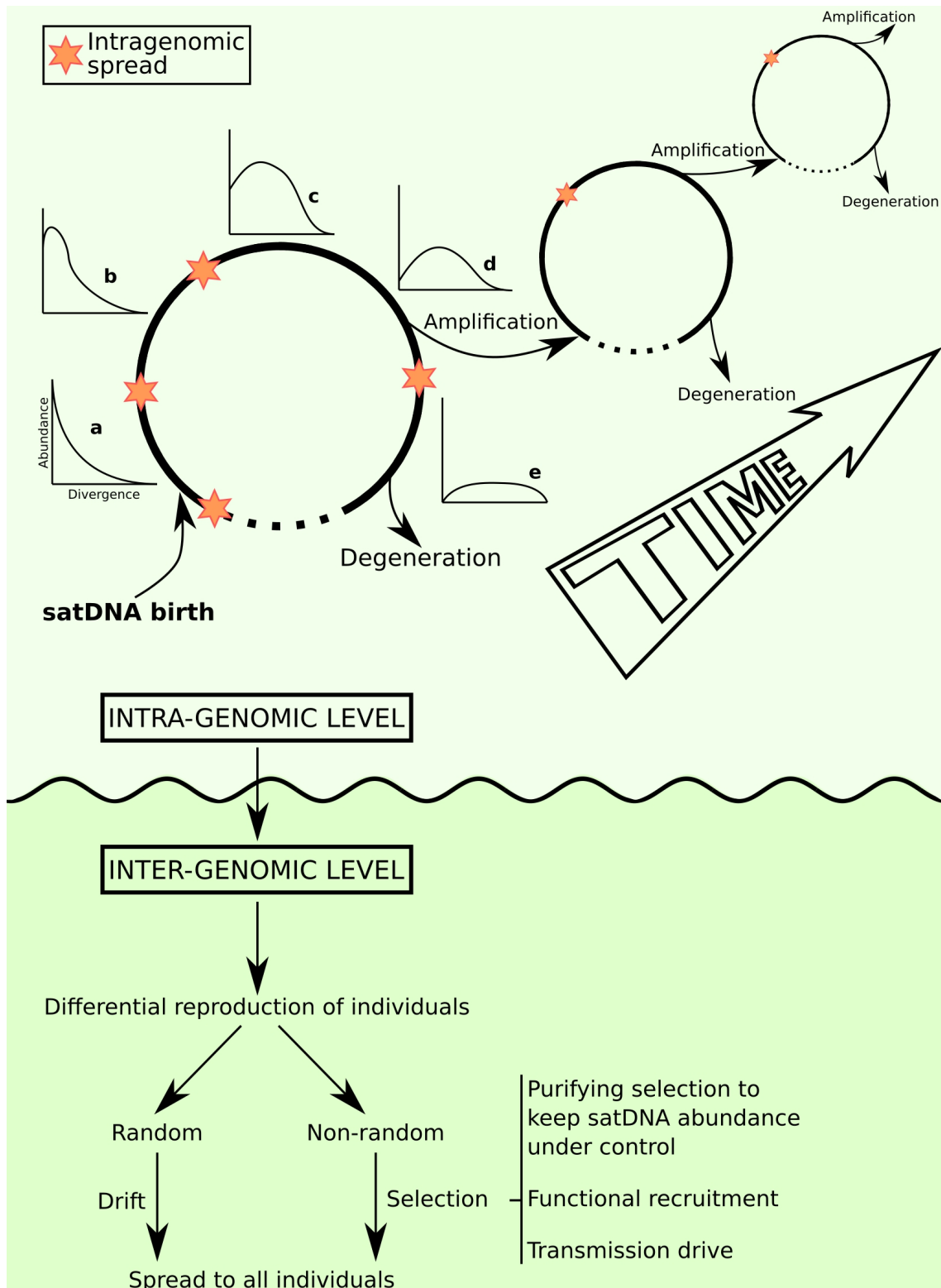
1476 the *L. migratoria* satDNA families being shared or non-shared with *O. decorus*. Note

1477 that shared satDNAs showed higher homogenization (higher RPS and kurtosis) and

1478 lower degeneration (5% effect size for mean difference in DIVPEAK) than non-shared

1479 ones, suggesting most recent amplification of the shared ones.

1480



1481

1482

Figure 6. A model of satDNA evolution. We consider that evolutionary events are

1483

rather different at intra- and inter-genomic levels. At intra-genomic level, tandem

1484

duplication gives birth to a new tandem repeat and its reiteration yields many copies of

1485

identical non-coding sequences (satDNA amplification). The newly amplified satDNA

1486 displays RLs sharply leptokurtic (a). As time goes by, point mutation increases
1487 divergence among the amplified sequences and the curve progressively is flattened (b-e)
1488 and DIVPEAK (i.e. the divergence value showing the higher abundance) increases (i.e.
1489 the peak moves to the right in the a-e graphs). At any moment of this first amplification-
1490 degeneration cycle, another sequence undergoes amplification and begins a new cycle.
1491 This sets the satDNA family farther from degeneration and extinction because its
1492 average divergence decreases and now predominates a newly amplified subfamily with
1493 leptokurtic RL (we represent here three successive cycles of amplification; note that the
1494 differences in size among cycles are to facilitate drawing and have nothing to do with
1495 amplification level). In parallel, an intra-genomic spread of the satDNA can occur at
1496 higher or lower extent (brown stars). A conceivable exit of these cycles is satDNA
1497 degeneration, when homology with the original sequence is lost. At inter-genomic level,
1498 individual reproduction will mark the destiny of the different satDNA sequences in
1499 populations. When reproduction is differential, albeit random (drift) or non-random
1500 (selection), some sequences may become prevalent above others. At this respect, the
1501 mutational-hazard hypothesis is applicable to explain the limits to purifying selection in
1502 some species showing extremely high abundance of satDNA. Finally, we cannot rule
1503 out that, in some case, transmission drive could help satDNA to prosper and, even that
1504 positive selection may recruit satDNA for important functions, such as telomeric or
1505 centromeric functions.

1506

1507

1508

1509

1510 **Tables**

Table 1. Comparison of satellitome characteristics between *O. decorus* and *L. migratoria* (Southern Lineage), by means of estimation graphics using DABEST (Ho et al. 2019). 95% CI= Confidence interval. RUL= Repeat unit length. * means that 95% CI does not include the zero value.

means that 95% CI does not include the zero value.

Comparison	Item	Mean (SE)		Effect size			Includes zero?
		<i>O. decorus</i> (N= 58)	<i>L. migratoria</i> (N=56)	Unpaired mean difference	CI_low	CI_high	
All satDNAs	RUL	201.5 (13.6)	152.7 (14)	48.8	12.1	86.6	*
	A+T (%)	55.7 (1.2)	54.4 (1.1)	1.27	-1.81	4.38	
	Abundance (%)	0.044 (0.013)	0.038 (0.019)	0.0055	-0.0557	0.0415	
	Divergence	7.19 (0.56)	7.09 (0.61)	0.093	-1.55	1.75	
Shared satDNAs		<i>O. decorus</i> (N= 21)	<i>L. migratoria</i> (N= 20)				
	RUL	212.8 (12.6)	216.5 (14.1)	-3.69	-39.4	33.3	
	A+T (%)	58.3 (1.1)	58.0 (1.1)	0.333	-2.8	3.27	
	Abundance (%)	0.071 (0.033)	0.091 (0.052)	-0.0196	-0.171	0.0715	
	Divergence	8.08 (1.22)	4.90 (0.50)	3.18	1.19	6.34	*
Non-shared satDNAs		<i>O. decorus</i> (N= 37)	<i>L. migratoria</i> (N= 36)				
	RUL	195.1 (20.2)	117.2 (17.8)	77.9	26.7	129	*
	A+T (%)	54.2 (1.7)	52.5 (1.6)	1.76	-2.75	6.21	
	Abundance (%)	0.028 (0.01)	0.009 (0.002)	0.019	0.00635	0.0496	*
	Divergence	6.68 (0.53)	8.31 (0.84)	-1.63	-3.64	0.244	
<i>O. decorus</i>		Shared (N= 21)	Non-shared (N= 37)				
	RUL	212.8 (12.6)	195.1 (20.2)	17.7	-34.4	58.3	
	A+T (%)	58.3 (1.1)	54.2 (1.7)	4.11	0.299	8.19	*
	Abundance (%)	0.071 (0.033)	0.028 (0.01)	0.0434	-	0.139	
	Divergence	8.08 (1.22)	6.68 (0.53)	1.4	-0.699	4.63	
<i>L. migratoria</i>		Shared (N= 20)	Non-shared (N= 36)				
	RUL	216.5 (14.1)	117.2 (17.8)	99.3	50	139	*
	A+T (%)	58.0 (1.1)	52.5 (1.6)	5.45	1.95	9.43	*
	Abundance (%)	0.091 (0.052)	0.009 (0.002)	0.082	0.018	0.261	*
	Divergence	4.90 (0.50)	8.31 (0.84)	-3.41	-5.42	-1.59	*

1511

1512

1513

1514

Table 2. Characteristics of the orthologous satDNA families analyzed in *O. decorus* (14) and *L. migratoria* (20). Each row includes one Ode and one Lmi satDNA families showing homology between them. Note that some Ode families showed homology with two or three Lmi ones. OSF= Orthologous superfamily, sf= number of subfamilies, FISH= FISH pattern (B= banded, NS= no signal), abun= abundance (% of the genome), RPS= Relative peak size, DP= DIVPEAK, MAL= Maximum array length observed in Minlon reads of *L. migratoria*, CEI= Concerted evolution index (L= *L. migratoria*, O= *O. decorus*), Intid= Interspecific sequence identity (%), Intdiv= Interspecific divergence, CTR= Consensus turnover rate, ILibS= Incomplete library sorting. Negative CEI values and Int_id>95% are remarked in bold type letter. See Table S4 to complete data with repeat unit length, A+T content, divergence (%), peak size, kurtosis of the repeat landscape, tandem structure index and Gibbs free energy of the secondary structure.

<i>O. decorus</i>							<i>Locusta migratoria</i>							Interspecific comparisons						
OSF	Name	sf	FISH	abun	RPS	DP	Name	sf	FISH	abun	RPS	DP	MAL	CEI_O	CEI_L	Int_id	Int_div	CTR	ILibS	
1	OdeSat01-287	1	B	6.2E-03	87%	1	LmiSat09-181	5	B	3.0E-04	65%	0	4417	88.4	85.6	68.9	90.8	1.990	0.30	
2	OdeSat02-204	4	B	3.3E-03	51%	2	LmiSat03-195	6	B	3.0E-03	63%	3	13447	124.5	125.1	60.6	130.4	2.858	0	
3	OdeSat17-176	1	NS	2.0E-04	29%	27	LmiSat02-176	1	B	3.6E-03	68%	4	20180	-24.6	-5.1	99.4	0.6	0.013	1.00	
4	OdeSat21-228	3	NS	1.5E-04	58%	3	LmiSat51-241	1	B	2.9E-05	61%	3	1708	67.0	66.5	71.8	72.8	1.596	0.44	
4	OdeSat32-238	2	B	8.5E-05	36%	2	LmiSat26-240	2	B	1.0E-04	60%	3	1455	40.5	47.8	77.7	53.1	1.164	0.59	
4							LmiSat37-238	1	B	4.6E-05	59%	3	2454	54.4	59.5	75.6	67	1.469	0.49	
5	OdeSat22-267	3	B	1.4E-04	59%	1	LmiSat12-273	3	B	1.3E-04	74%	1	2948	90.6	94.8	75	98.1	2.150	0.25	
5							LmiSat16-278	1	B	1.4E-04	87%	2	1965	89.5	94.6	72.6	97	2.126	0.26	
6	OdeSat26-180	1	B	1.3E-04	88%	2	LmiSat41-180	1	B	5.1E-05	94%	3	515	29.2	28.2	74.4	31.7	0.695	0.76	
7	OdeSat28-276	1	B	1.2E-04	56%	5	LmiSat24-266	1	NS	5.9E-05	90%	0	1378	49.4	53.4	67.9	55.8	1.223	0.57	
7							LmiSat45-274	1	B	2.5E-05	54%	2	945	19.0	16.8	79.7	25.4	0.557	0.81	
7							LmiSat54-272	1	B	1.6E-05	65%	0	2073	58.7	60.2	66.3	65.1	1.427	0.50	
7	OdeSat58-265	2	NS	9.5E-06	88%	0	LmiSat28-263	2	B	6.0E-05	97%	0	2821	30.1	32.4	77.5	33.9	0.743	0.74	
7							LmiSat43-231	1	B	3.9E-05	100%	0		39.3	42.7	69.3	43.1	0.945	0.67	
8	OdeSat39-185	2	NS	6.8E-05	67%	4	LmiSat06-185	4	B	4.9E-04	66%	3	19168	14.9	16.1	84.3	21	0.460	0.84	
9	OdeSat41-75	1	NS	6.1E-05	29%	18	LmiSat27-57	1	NS	5.4E-05	32%	0	712	-2.4	7.1	92.7	16.2	0.355	0.88	
10	OdeSat56-249	1	NS	2.0E-05	93%	0	LmiSat32-261	1	B	3.9E-05	60%	0	1489	31.5	26.4	77.2	32.9	0.721	0.75	
11	OdeSat57-75	1	NS	1.4E-05	40%	4	LmiSat17-75	1	B	1.2E-04	48%	2	3194	-1.3	2.7	92	8.5	0.186	0.93	
12	OdeSat59-185	1	NS	5.8E-06	36%	3	LmiSat01-185	5	B	9.8E-03	46%	3	17619	-0.9	7.2	98.9	11.8	0.259	0.91	
12							LmiSat13-259	5	B	1.5E-04	76%	4	1379	44.1	52.3	63.3	56.8	1.245	0.56	
																Mean	77.3	50.6	1.109	61%
																SD	11.1	34.7	0.76	27%
																CV	14%	69%	69%	44%