

1 Genome assembly of *Danaus chrysippus* and comparison with the
2 Monarch *Danaus plexippus*

3

4 Kumar Saurabh Singh^{1*}, Rishi De-Kayne^{2*}, Kennedy Saitoti Omufwoko³, Dino J. Martins^{3,4},
5 Chris Bass⁵, Richard ffrench-Constant⁵ & Simon H. Martin^{2†}

6

7 ¹Bioinformatics Group, Wageningen University, Wageningen, The Netherlands

8 ²Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, United Kingdom

9 ³Department of Ecology and Evolutionary Biology, Princeton University, Princeton, USA

10 ⁴Mpala Research Centre, Nanyuki, Kenya

11 ⁵Centre for Ecology and Conservation, University of Exeter, Penryn Campus, Penryn, United
12 Kingdom

13

14 * These authors contributed equally

15 † Corresponding author (simon.martin@ed.ac.uk)

16

17 Running head: Genome of the African Monarch butterfly

18

19 Keywords: African Monarch; African Queen; Plain Tiger; repeat content; intron length; genome
20 size; butterfly

21 **Abstract**

22
23 Milkweed butterflies in the genus *Danaus* are studied in a diverse range of research fields
24 including the neurobiology of migration, biochemistry of plant detoxification, host-parasite
25 interactions, evolution of sex chromosomes, and speciation. We have assembled a nearly
26 chromosomal genome for *Danaus chrysippus* (known as the African Monarch, African Queen,
27 and Plain Tiger) using long read sequencing data. This species is of particular interest for the
28 study of genome structural change and its consequences for evolution. Comparison with the
29 genome of the North American Monarch *Danaus plexippus* reveals generally strong synteny,
30 but highlights three inversion differences. The three chromosomes involved were previously
31 found to carry peaks of intra-specific differentiation in *D. chrysippus* in Africa, suggesting that
32 these inversions may be polymorphic and associated with local adaptation. The *D. chrysippus*
33 genome is over 40% larger than that of *D. plexippus*, and nearly all of the additional ~100
34 Megabases of DNA comprises repeats. Future comparative genomic studies within this genus
35 will shed light on the evolution of genome architecture.

36

37 **Introduction**

38

39 The genus *Danaus* is perhaps best known for the iconic Monarch butterfly *Danaus plexippus*
40 and its extraordinary migrations in North America. Genomic studies of the Monarch have shed
41 light on host plant detoxification (Tan *et al.* 2019), sex chromosome evolution (Mongue *et al.*
42 2017; Gu *et al.* 2019) and the genetic basis of migratory behaviour (Zhan *et al.* 2014). Its
43 relative *Danaus chrysippus* is found throughout Africa, the Mediterranean and south Asia, and is
44 known as the African Monarch, African Queen, and Plain Tiger butterfly in different parts of its
45 range. *D. chrysippus* is emerging as a useful study system in evolutionary genomics. Several
46 subspecies of *D. chrysippus* with distinct warning patterns occupy distinct geographic ranges
47 separated by broad hybrid zones (Smith *et al.* 1997; Lushai *et al.* 2003). Patterns of genetic
48 differentiation suggest a role for chromosomal rearrangements in maintaining these differences
49 (Martin *et al.* 2020). In the east African hybrid zone, a neo-W sex chromosome has emerged in
50 the past few thousand years and is associated with infection by a male-killing endosymbiont
51 *Spiroplasma* (Smith *et al.* 2016; Martin *et al.* 2020). This species therefore has great potential
52 for future research on the evolutionary impacts of genome structural change.

53

54 Here we describe the generation of a chromosome-level assembly for *D. chrysippus* based on
55 Pacific Biosciences long read sequencing data. This serves to replace a previous reference
56 genome based on short read sequences and mate-pair libraries, which had low contiguity
57 ($N_{50}=0.63\text{ Mb}$, Martin *et al.* 2020). Our new assembly has an N_{50} of 11.45 Mb. Nineteen of the
58 30 chromosomes are represented by a single contig, and the remaining eleven by two contigs
59 each. At 354 Mb, this genome is average in size for a butterfly, but about 1.4 times larger than
60 that of *D. plexippus* (~250 Mb). Comparative analyses indicate that this difference is largely
61 explained by increased repeat content, but *D. chrysippus* also has larger introns, implying that
62 these species have experienced different selection pressures acting on non-essential DNA.

63

64 **Materials and Methods**

65

66 *DNA sequencing*

67 High-molecular-weight DNA was extracted from a single female pupa from a captive butterfly
68 stock using the Qiagen Blood & Cell Culture DNA Mini Kit following the manufacturer's
69 guidelines. Long read Pacific Biosciences sequencing was performed using seven PacBio
70 Sequel SMRT cells on the Sequel platform, yielding approximately 9.7 gigabases (Gb) per
71 SMRT cell. The 3.8 million PacBio reads totalled 67.6 Gb and had an N50 of 27.3 kb. In
72 addition, we generated Illumina sequencing data for the same individual on the Novaseq 600
73 platform (118 million paired-end reads of 150 bp with an insert size of 350 bp) totalling 35 Gb.
74

75 *Genome assembly*

76 We assembled the long reads using both Canu (Koren *et al.* 2017) and Falcon (Chin *et al.*
77 2016), and then merged these assemblies to maximize the genome completeness using
78 quickmerge -v 0.3 (Chakraborty *et al.* 2016). Redundant contigs or haplotigs were removed
79 using Purge_haplotigs -v 1.0.4 (Roach *et al.* 2018) with the -align_cov (Percent cutoff for
80 identifying a contig as haplotig) value of 65. Before merging, assemblies were polished
81 iteratively using three rounds of Pilon -v 1.22 in diploid mode (Walker *et al.* 2014; using a
82 trimmed version of the short-read data; reads were trimmed using Trim_Galore -v 0.4.0;
83 Krueger 2012), and Racon -v 1.3.1 (Vaser *et al.* 2017; using the long-read data). Illumina and
84 PacBio raw reads are archived under European Nucleotide Archive project accession:
85 PRJEB47812.

86

87 *Whole genome alignment and synteny assessment*

88 To assess synteny and putatively assign contigs to chromosomes, we aligned the *D. chrysippus*
89 assembly to two *Danaus plexippus* assemblies: 'Dplex_v4', a chromosome-level assembly
90 produced by scaffolding 4115 scaffolds using chromatin conformation (Hi-C) data
91 (GCA_009731565.1, (Gu *et al.* 2019)) and 'MEX_DaPlex', a long-read based assembly
92 consisting of 66 scaffolds, of which 38 (97% of total sequence) have been assigned to
93 chromosomes (GCA_018135715.1, (Ranz *et al.* 2021)). Alignments were generated using both
94 MUMmer's nucmer tool version 3.1 (Marçais *et al.* 2018), with default parameters except the
95 'maxGap' parameter set to 1000, and with minimap2 v2.17 (Li 2018), using the 'asm20'
96 parameter preset, designed for whole genome alignment of species with sequence divergence
97 below 20%. Nucmer alignments were explored using the interactive alignment visualisation tool
98 Dot (<https://github.com/dnanexus/dot>) and final alignment plots based on minimap2 alignments
99 were generated using Asynt (<https://github.com/simonhmartin/asynt>).
100

101 *Correcting putative misassemblies*

102 Visualisation of whole genome alignment to both *Danaus plexippus* assemblies (described
103 above) revealed two putatively misassembled contigs that had portions aligning confidently to
104 two different chromosomes. Although these could theoretically represent real translocation or

105 fusion products, we took the conservative decision to split these contigs, as additional long-read
 106 assemblies generated in a related study (Kim *et al.* 2021) showed no evidence for translocations
 107 or fusions. Optimal split points were identified by visual inspection of the alignments, as well as
 108 additional BLASTn alignments made between the two genomes. The original unsplit assembly,
 109 along with details of split points, is available at <https://doi.org/10.5281/zenodo.5731560>.

110
 111 Despite having performed automated removal of redundant contigs using Purge_haplotigs,
 112 visual exploration of the alignments identified a further four contigs that appeared to be
 113 redundant (i.e. representing a part of the genome already represented by a larger contig).
 114 These included one of the split products described above. To confirm this, we aligned the
 115 Illumina reads for the assembled individual back to the assembly using BWA MEM (Li and
 116 Durbin 2010) using default parameters, and computed read depth using Samtools depth (Li *et*
 117 *al.* 2009). Visualisation of median read depth averaged in 50 kb windows confirmed that these
 118 four contigs had 50% depth, so they were removed from the assembly. Finally, two contigs
 119 included portions that appeared to be redundant in the alignments as well as read depth plots.
 120 These were therefore split at the point in the alignment where the redundancy began, and the
 121 redundant fragment was removed from the final assembly. The full original assembly, along with
 122 details of all splits and portions retained to produce the final Dchry2.2 assembly, is available at
 123 <https://doi.org/10.5281/zenodo.5731560>. To assess the base-level accuracy of our assembly we
 124 calculated the consensus quality (QV), comparing the frequency of k-mers present in the raw
 125 Illumina reads with those present across the final assembly (all 83 contigs) using Merqury v.1.3
 126 (Rhee *et al.* 2020).

127

128 *Repeat Annotation*

129 To assess the repeat content of the assembly, the genome was masked using a custom repeat
 130 library. First, a repeat library was produced using the finished genome assembly, using
 131 RepeatModeler v2.0.1 (Smit and Hubley 2008), and this library was then combined with a broad
 132 Lepidoptera repeat database extracted using RepeatMasker v.4.1.0 (Smit *et al.* 2015). Repeat
 133 masking of the genome was then carried out using RepeatMasker (Smit *et al.* 2015). To
 134 determine the prevalence of expanding transposable element families within *Danaus chrysippus*
 135 the scripts calDivergenceFromAlign.pl and createRepeatLandscape.pl from RepeatModeler
 136 (Smit and Hubley 2008) were used to produce a repeat landscape for the assembly. To facilitate
 137 a comparison with other *Danaus* genome assemblies, this repeat masking process was
 138 repeated using the same custom repeat library for two well-assembled *Danaus plexippus*
 139 genome assemblies (NCBI accessions GCF_009731565.1 and GCA_018135715.1 (Ranz *et al.*
 140 2021)). The resulting softmasked assemblies were then used for genome annotation.

141

142 *Gene Annotation*

143 Due to a lack of RNAseq data, a preliminary genome annotation was carried out using two
 144 protein sets from the close relative to *D. chrysippus*, *D. plexippus*, the Monarch butterfly. This
 145 combined protein set was produced by collating protein information from two different, *D.*
 146 *plexippus* assemblies, the first a proteome downloaded from uniprot under the accession
 147 UP000596680 (associated with the Dplex_v4 assembly), and the second taken by extracting

148 amino acid sequences from the annotation of the 'MEX_DaPlex' *Danaus plexippus* assembly
 149 GCF_009731565.1 (Ranz *et al.* 2021) (both protein sets had high BUSCO scores, indicative of
 150 high-quality annotation). This combined protein set was then used as input for the BRAKER2
 151 pipeline (Hoff *et al.* 2018) to annotate each of the three soft masked genome assemblies
 152 produced above (specifying --gff3 --softmasking --prot_seq=protein_set.fasta --prg=gth --
 153 gth2traingenes --trainFromGth). GenomeTools (Gremme *et al.* 2013) was then used to sort and
 154 tidy the annotation output (gt gff -sort -tidy -retainids -fixregionboundaries) and calculate
 155 summary statistics of the annotation (gt stat -genelengthdistri -genescoredistri -exonlengthdistri
 156 -exonnumberdistri -intronlengthdistri -cdslengthdistri). Functional annotation for the resulting
 157 *Danaus chrysippus* protein set was carried out using Pannzer2 (Törönen *et al.* 2018). To
 158 determine variation in intron and exon length between *D. chrysippus* and *D. plexippus*, introns
 159 and exons were extracted from our corresponding annotation file for each of the three
 160 assemblies.
 161

162 *Genome comparison and assembly validation*

163 To assess the quality of annotation, BUSCOs were calculated for both the full *D. chrysippus*
 164 assembly and the protein sequences resulting from annotation specifying the insecta_odb10
 165 lineage BUSCO set in BUSCO v.5.0.0 (Simão *et al.* 2015). Additionally, the annotation was
 166 compared against those of both published *D. plexippus* annotations. To compare each of the
 167 assemblies, and in turn the consistency of genome structure across *Danaus* species we plotted
 168 the distribution of intron lengths for annotations from each of the three assemblies. This was
 169 carried out by extracting introns and exons from the longest annotated transcript for each gene
 170 within each of the annotations (using the BRAKER2 re-annotations for both *D. plexippus*
 171 assemblies to ensure lengths were comparable across assemblies).

172

173 *Data availability*

174 The assembly and annotation are available at the European Nucleotide Archive project
 175 accession: PRJEB47812. Additional data files are provided at
 176 <https://doi.org/10.5281/zenodo.5731560>: purged haplotigs, assembly before manual edits,
 177 details of manual edits made to the assembly, and repeat library and functional annotation files.
 178 Scripts for genome assembly are available at <https://github.com/kumarsaurabh20/DChry2.1> and
 179 scripts for the genome annotation and analysis of introns and exons at
 180 https://github.com/RishiDeKayne/Danaus_Dchry2.2_annotation.

181

182 **Results and Discussion**

183

184 *Genome assembly*

185 67.6 Gb of long-read data was assembled into 83 contigs. Manual splitting of three putatively
 186 misassembled contigs and removal of several remaining redundant fragments (see Materials
 187 and Methods) left 83 contigs with an N50 of 11.45 Mb and L50 of thirteen contigs, giving a total
 188 genome size of 354 Mb. Alignment with two different *D. plexippus* assemblies allowed us to
 189 confidently assign 41 contigs representing 97% of the sequence length to chromosomes (Figure
 190 1). Of the 30 *D. chrysippus* chromosomes, 19 are represented by a single contig and the rest by
 191 two contigs each. The contiguity of our assembly is therefore comparable to that of the *D.*

192 *plexippus* 'MEX_DaPlex' assembly (Ranz *et al.* 2021), for which 38 out of 66 scaffolds (97% of
193 the assembly) were assigned to chromosomes, of which 23 are represented by a single
194 scaffold. Among the 42 *D. chrysippus* contigs that were not assigned to a chromosome (3% of
195 the genome) it is possible that some represent fragments of the female-specific W chromosome.
196 However, given that butterfly W chromosomes are highly repetitive and have low inter-specific
197 homology (Lewis *et al.* 2021), further work comparing male and female genomes is required to
198 test this hypothesis. The genome-wide consensus quality of the assembly (QV; representing a
199 log-scaled probability of error for each base in our assembly) was 36.2373, suggesting a high
200 level of accuracy (equating to an accuracy between 99.9% and 99.99%).
201

202 Synteny and genome size comparison

203 The genomes of *D. chrysippus* and *D. plexippus* are largely syntenic (Figure 1). Our assembly
204 supports the earlier finding that the Z sex chromosome of *Danaus* species represents a fusion
205 between the ancestral lepidopteran Z chromosome and autosome 21, which occurred in a
206 recent ancestor of the genus (Mongue *et al.* 2017). We numbered chromosomes according to
207 their homologs in the most recent *D. Plexippus* assembly (Ranz *et al.* 2021), which follows the
208 growing convention of using the chromosome numbering system introduced for *Melitaea cinxia*,
209 the first assembled lepidopteran genome that retains the ancestral karyotype of 31 (Ahola *et al.*
210 2014; Cicconardi *et al.* 2021; Ranz *et al.* 2021; Lewis *et al.* 2021). As such, the *Danaus*
211 karyotype lacks an autosome 21, as this is now part of the Z sex chromosome.
212

213 Several putative inversion differences can be identified between the two *Danaus* species, most
214 notably on chromosomes 4, 17 and 30 (Figure 1). We note that all three of these chromosomes
215 were found to carry sharp peaks of intra-specific differentiation between subspecies of *D.*
216 *chrysippus* in Africa, against a background of very low genetic differentiation (Martin *et al.*
217 2020), suggesting that these putative inversions may be polymorphic and associated with local-
218 adaptation in *D. chrysippus*. In addition, chromosomes 15, 26 and 29 all carry large
219 duplicated/repetitive portions relative to *D. plexippus*. One of these, on chromosome 15, was
220 identified previously as the site of a large expansion in gene copy number through multiple
221 duplications and is associated with subspecies differentiation and colour pattern variation in *D.*
222 *chrysippus* (Martin *et al.* 2020). Further work to dissect this genomic region and compare
223 chromosome structure among *D. chrysippus* subspecies is ongoing.
224

225 In total, the 354 Mb *D. chrysippus* genome is 42-44% larger than that of *D. plexippus* (245-248
226 Mb). This difference is consistent for all autosomes, but most dramatic for the three
227 chromosomes carrying large repetitive/duplicated tracts: namely 15, 26 and 29 (Figure 2). By
228 contrast, the Z sex chromosome is nearly identical in size in the two species. This difference in
229 autosome sizes could be explained either via a systematic size reduction in the lineage leading
230 to *D. plexippus*, or a systematic increase in the lineage leading to *D. chrysippus*. These
231 hypotheses can be distinguished by comparison with assemblies of other members of the genus
232 or outgroups in the future.
233

234 *Transposable element and repeat content*

235 In total the *D. chrysippus* genome comprises 35.5% repeats, with the largest proportion of these
 236 being retroelements which make up 11.9% of the genome sequence (Figure 3). Repeat
 237 masking of each of the *D. plexippus* assemblies revealed that a substantially lower proportion of
 238 the genomes of these close relatives comprise repeats, only between 11.2% and 14.3%. Each
 239 of the main classified repeat families are more abundant in *D. chrysippus* compared to *D.*
 240 *plexippus*, with the largest difference between the species observed for the rolling-circle family
 241 which represents 1.8% and 2% of the *D. plexippus* genome sequence, compared to 7.6% of the
 242 *D. chrysippus* sequence. The repeat landscape of *D. chrysippus* (Figure 4) highlights a number
 243 of expanding repeat families, most strikingly the rolling-circle repeats RC/Helitron (pink; Figure
 244 4). The increased prevalence of repetitive elements within the *D. chrysippus* genome (91-98Mb
 245 more than *D. plexippus*) largely explains the larger genome size of *D. chrysippus* compared to
 246 *D. plexippus* (an increase of 106-109Mb). Although the repeat content of genomes across the
 247 Lepidoptera order has been shown to vary substantially (Talla *et al.* 2017) our results suggest
 248 that even within a genus a large amount of variation can be present. Although the genome of *D.*
 249 *chrysippus* is rather repetitive, even within Lepidoptera, (With repeats making up 35.5% of the
 250 genome), *D. plexippus* tends towards the lower end of repeat content (repeats make up 11.2-
 251 14.3% of the genome).

252

253 *Gene content*

254 In total 16,260 protein-coding genes were annotated in the assembly by BRAKER2 (with 19,639
 255 protein-coding mRNAs annotated - accounting for multiple transcripts/isoforms of the same
 256 gene), which included 136,694 exons and 117,106 introns. This number of genes is similar to
 257 that of the published annotations for each of the *D. plexippus* assemblies, which annotated
 258 15,006 (Dplex_v4) and 15,995 (MEX_DaPlex) genes (as well as our re-annotated versions of
 259 these assemblies which annotated 18,663 and 21,311 genes; in both cases our annotation
 260 involved annotating additional smaller scaffolds not annotated in the original assemblies - 284
 261 vs 30 scaffolds for Dplex_v4 and 64 vs 55 for MEX_DaPlex). An analysis of BUSCOs using the
 262 *insecta_odb10* benchmarking set shows that the full genome sequences and annotated protein
 263 set for *D. chrysippus* are 98.2% and 96.3% complete for BUSCOs, respectively. This
 264 percentage completeness is close to that of both published *D. plexippus* annotations which have
 265 94.6% (Dplex_v4) and 97.1% (MEX_DaPlex) complete BUSCOs (Table 1). Pannzer2 allowed
 266 us to add functional annotation to 9,567 of the full 16,260 gene set (functional annotation
 267 available at <https://doi.org/10.5281/zenodo.5731560>).

268

269 *Intron and exon length*

270 Exon length is relatively consistent across the three genomes, ranging from 217bp (Dplex_v4
 271 re-annotation) to 238bp (MEX_DaPlex re-annotation) (Figure 5A). However, mean intron length
 272 in the *D. Chrysippus* assembly (975bp) is higher than that in the two *D. plexippus* assemblies
 273 (665bp and 738bp, respectively) (Figure 5B). This substantial increase in intron length in *D.*
 274 *chrysippus* likely explains the remaining variation in genome size between *D. chrysippus* and *D.*
 275 *plexippus*. This difference may represent a neutral increase in introns in *D. chrysippus* or a

276 selection-mediated reduction in intron size in *D. plexippus*. These hypotheses may be resolved
277 by comparison with genomes of other members of the genus in the future.

278

279 *Conclusions*

280 We have assembled a nearly chromosome-level genome for *D. chrysippus*, which is highly
281 comparable in its quality to the best available assembly for *D. plexippus*. Although the two
282 species retain strong synteny, the *D. chrysippus* genome is >40% larger, with more repetitive
283 content and larger introns on average. This implies stronger selection to limit non-essential DNA
284 in *D. plexippus*. Future comparative studies involving other members of the genus could shed
285 light on the processes and forces driving the evolution of genome size. The *D. chrysippus*
286 genome will also serve as a reference for population genomic studies to test hypotheses about
287 the evolution of warning colouration, host-parasite interactions and the consequences of
288 chromosomal rearrangements.

289

290 **Acknowledgements**

291

292 We thank Alexander Mackintosh for providing advice on genome assembly and annotation. This
293 work was supported by a Royal Society University Research Fellowship (URF\R1\180682) and
294 Enhancement Award (RGF\EA\181071) awarded to SHM, a Swiss National Science Foundation
295 Early Postdoc Mobility Fellowship (P2BEP3_195567) awarded to RDK, and a European
296 Research Council Horizon 2020 research and innovation programme grant 646625 awarded to
297 CB.

298

299 **Literature Cited**

- 300 Ahola, V., R. Lehtonen, P. Somervuo, L. Salmela, P. Koskinen *et al.*, 2014 The Glanville fritillary
301 genome retains an ancient karyotype and reveals selective chromosomal fusions in
302 Lepidoptera. *Nat. Commun.* 5: 4737.
- 303 Chakraborty, M., J. G. Baldwin-Brown, A. D. Long, and J. J. Emerson, 2016 Contiguous and
304 accurate de novo assembly of metazoan genomes with modest long read coverage.
305 *Nucleic Acids Res.* 44: e147.
- 306 Chin, C.-S., P. Peluso, F. J. Sedlazeck, M. Nattestad, G. T. Concepcion *et al.*, 2016 Phased
307 diploid genome assembly with single-molecule real-time sequencing. *Nat. Methods* 13:
308 1050–1054.
- 309 Cicconardi, F., J. J. Lewis, S. H. Martin, R. D. Reed, C. G. Danko *et al.*, 2021 Chromosome

- 310 fusion affects genetic diversity and evolutionary turnover of functional loci, but consistently
311 depends on chromosome size. Mol. Biol. Evol. 38: 4449–4462
- 312 Gremme, G., S. Steinbiss, and S. Kurtz, 2013 GenomeTools: a comprehensive software library
313 for efficient processing of structured genome annotations. IEEE/ACM Trans. Comput. Biol.
314 Bioinform. 10: 645–656.
- 315 Gu, L., P. F. Reilly, J. J. Lewis, R. D. Reed, P. Andolfatto *et al.*, 2019 Dichotomy of Dosage
316 Compensation along the Neo Z Chromosome of the Monarch Butterfly. Curr. Biol. 29:
317 4071–4077.
- 318 Hoff, K. J., A. Lomsadze, M. Stanke, and M. Borodovsky, 2018 BRAKER2: incorporating protein
319 homology information into gene prediction with GeneMark-EP and AUGUSTUS. Plant and
320 Animal Genomes XXVI.
- 321 Kim, K.W., R. De-Kayne, I. J. Gordon, K. S. Omufwoko, D. J. Martins, R. ffrench-Constant, S. H.
322 Martin, 2021 Stepwise evolution of a butterfly supergene via duplication and inversion.
323 Biorxiv BIORXIV/2021/471392.
- 324 Kolmogorov, M., J. Yuan, Y. Lin, and P. A. Pevzner, 2019 Assembly of long, error-prone reads
325 using repeat graphs. Nat. Biotechnol. 37: 540–546.
- 326 Koren, S., B. P. Walenz, K. Berlin, J. R. Miller, N. H. Bergman *et al.*, 2017 Canu: scalable and
327 accurate long-read assembly via adaptive k-mer weighting and repeat separation. Genome
328 Res. 27: 722–736.
- 329 Krueger, F., 2012 Trim Galore: a wrapper tool around Cutadapt and FastQC to consistently
330 apply quality and adapter trimming to FastQ files, with some extra functionality for Mspl-
331 digested RRBS-type (Reduced Representation Bisulfite-Seq) libraries. URL http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/. (Date of access: 28/04/2016).
- 333 Lewis, J. J., F. Cicconardi, S. H. Martin, R. D. Reed, C. G. Danko *et al.*, 2021 The Dryas iulia

- 334 Genome Supports Multiple Gains of a W Chromosome from a B Chromosome in
335 Butterflies. *Genome Biol. Evol.* 13: evab128
- 336 Li, H., 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34: 3094–
337 3100.
- 338 Li, H., and R. Durbin, 2010 Fast and accurate long-read alignment with Burrows-Wheeler
339 transform. *Bioinformatics* 26: 589–595.
- 340 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence
341 Alignment/Map format and SAMtools. *Bioinformatics* 25: 2078–2079.
- 342 Lushai, G., D. A. S. Smith, I. J. Gordon, D. Goulson, J. A. Allen *et al.*, 2003 Incomplete sexual
343 isolation in sympatry between subspecies of the butterfly *Danaus chrysippus* (L.) and the
344 creation of a hybrid zone. *Heredity* 90: 236–246.
- 345 Marçais, G., A. L. Delcher, A. M. Phillippy, R. Coston, S. L. Salzberg *et al.*, 2018 MUMmer4: A
346 fast and versatile genome alignment system. *PLoS Comput. Biol.* 14: e1005944.
- 347 Martin, S. H., K. S. Singh, I. J. Gordon, K. S. Omufwoko, S. Collins *et al.*, 2020 Whole-
348 chromosome hitchhiking driven by a male-killing endosymbiont. *PLoS Biol.* 18: e3000610.
- 349 Mongue, A. J., P. Nguyen, A. Voleníková, and J. R. Walters, 2017 Neo-sex Chromosomes in
350 the Monarch Butterfly, *Danaus plexippus*. *G3* 7: 3281–3294.
- 351 Ranz, J. M., P. M. González, B. D. Clifton, N. O. Nazario-Yepiz, P. L. Hernández-Cervantes *et*
352 *al.*, 2021 A de novo transcriptional atlas in *Danaus plexippus* reveals variability in dosage
353 compensation across tissues. *Commun Biol* 4: 791.
- 354 Rhie, A., B. P. Walenz, S. Koren, & A. M. Phillippy, 2020 Merqury: reference-free quality,
355 completeness, and phasing assessment for genome assemblies. *Genome biology*, 21: 1–
356 27.
- 357 Roach, M. J., S. A. Schmidt, and A. R. Borneman, 2018 Purge Haplotype: allelic contig

- 358 reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* 19: 460.
- 359 Ruan, J., and H. Li, 2019 Fast and accurate long-read assembly with wtdbg2. *bioRxiv* 530972.
- 360 Simão, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, and E. M. Zdobnov, 2015
- 361 BUSCO: assessing genome assembly and annotation completeness with single-copy
- 362 orthologs. *Bioinformatics* 31: 3210–3212.
- 363 Smith, D. A. S., I. J. Gordon, W. Traut, J. Herren, S. Collins *et al.*, 2016 A neo-W chromosome
- 364 in a tropical butterfly links colour pattern, male-killing, and speciation. *Proc. Biol. Sci.* 283:
- 365 20160821
- 366 Smith, D. A. S., D. F. Owen, I. J. Gordon, and N. K. Lowis, 1997 The butterfly *Danaus*
- 367 *chrysippus* (L.) in East Africa: polymorphism and morph-ratio clines within a complex,
- 368 extensive and dynamic hybrid zone. *Zool. J. Linn. Soc.* 120: 51–78.
- 369 Smit, A. F. A., and R. Hubley, 2008 RepeatModeler Open-1.0. Available fom <http://www.repeatmasker.org>.
- 370
- 371 Smit, A. F. A., R. Hubley, and P. Green, 2015 RepeatMasker Open-4.0. 2013--2015.
- 372 Talla, V., A. Suh, F. Kalsoom, V. Dinca, R. Vila *et al.*, 2017 Rapid Increase in Genome Size as a
- 373 Consequence of Transposable Element Hyperactivity in Wood-White (Leptidea) Butterflies.
- 374 *Genome Biol. Evol.* 9: 2491–2505.
- 375 Tan, W.-H., T. Acevedo, E. V. Harris, T. Y. Alcaide, J. R. Walters *et al.*, 2019 Transcriptomics of
- 376 monarch butterflies (*Danaus plexippus*) reveals that toxic host plants alter expression of
- 377 detoxification genes and down-regulate a small number of immune genes. *Mol. Ecol.* 28:
- 378 4845–4863.
- 379 Törönen, P., A. Medlar, and L. Holm, 2018 PANNZER2: a rapid functional annotation web
- 380 server. *Nucleic Acids Res.* 46: W84–W88.
- 381 Vaser, R., I. Sović, N. Nagarajan, and M. Šikić, 2017 Fast and accurate de novo genome

- 382 assembly from long uncorrected reads. *Genome Res.* 27: 737–746.
- 383 Walker, B. J., T. Abeel, T. Shea, M. Priest, A. Abouelliel *et al.*, 2014 Pilon: an integrated tool for
384 comprehensive microbial variant detection and genome assembly improvement. *PLoS One*
385 9: e112963.
- 386 Zhan, S., W. Zhang, K. Niitepõld, J. Hsu, J. F. Haeger *et al.*, 2014 The genetics of monarch
387 butterfly migration and warning colouration. *Nature* 514: 317–321.

388

389 **Table 1.** BUSCOs for the *D. chrysippus* genome assembly and each *D. plexippus* assembly in addition to
 390 the protein sequences resulting from both the original and re-annotation of each of these assemblies
 391 (using the insecta_odb10 BUSCO set, n=1367).

392

		% complete	% single	% duplicated	% fragmented	% missing
	Dchry2.2	98.2	97.5	0.7	0.6	1.2
Genomes	Dplex_v4	98.7	98.2	0.5	0.7	0.6
	MEX_DaPlex	98.9	98.2	0.7	0.4	0.7
	Dchry2.2	96.3	79.7	16.6	1.1	2.6
	Dplex_v4	94.6	93.7	0.9	2.5	2.9
Protein sets	Dplex_v4 (re-annotated)	98.1	72.6	25.5	1.2	0.7
	MEX_DaPlex	97.1	87.1	10.0	1.3	1.6
	MEX_DaPlex (re-annotated)	98.8	72.6	26.2	0.7	0.5

393

394 **Figure 1.** (A) Four colour morphs of *D. chrysippus* (above) and *D. plexippus* (below). (B) Whole genome
 395 alignment between *D. chrysippus* and *D. plexippus* (MEX_DaPlex assembly). Points represent minimap2
 396 alignments greater than 5kb in length. Alignments in the same orientation are shown in turquoise and
 397 those in the reverse orientation are shown in red. Only contigs that were confidently assigned to
 398 chromosomes (97% of the total in both assemblies) are included. Alternating grey and white bars indicate
 399 separate chromosomes. See Figure S1 for the same plot including contig/scaffold labels. Butterfly images
 400 from top to bottom by Forest Jarvis (CC-BY-NC), Paul Dickson (CC-BY-NC), Claude Martin, Steven
 401 Schulting (CC-BY-NC) and Edward Perry IV (CC-BY-NC).

402
 403 **Figure 2.** Chromosome length comparison between *D. chrysippus* and *D. plexippus* (MEX_DaPlex
 404 assembly). Chromosome lengths represent the sum total of the contigs/scaffolds assigned to each
 405 chromosome. Autosomes are shown in black and the Z sex chromosome in grey.

406
 407 **Figure 3.** Barplot showing the proportion of the genome of *Danaus chrysippus* (yellow) and *D. plexippus*
 408 (represented by both the MEX_DaPlex ,in red, and Dplex_v4, in turquoise, assemblies) made up of
 409 repetitive elements (as identified and masked by repeatmodeler and repeatmasker). Additionally, the
 410 proportion of each genome made up of specific repetitive element families is shown highlighting the
 411 increased proportion of repetitive elements in *Danaus chrysippus*.

412
 413 **Figure 4.** The repeat landscape of the Dchry2.2 assembly. In addition to unclassified repeats, rolling-
 414 circle (RC/Helitron), LINE and LTR families all appear to have expanded recently.

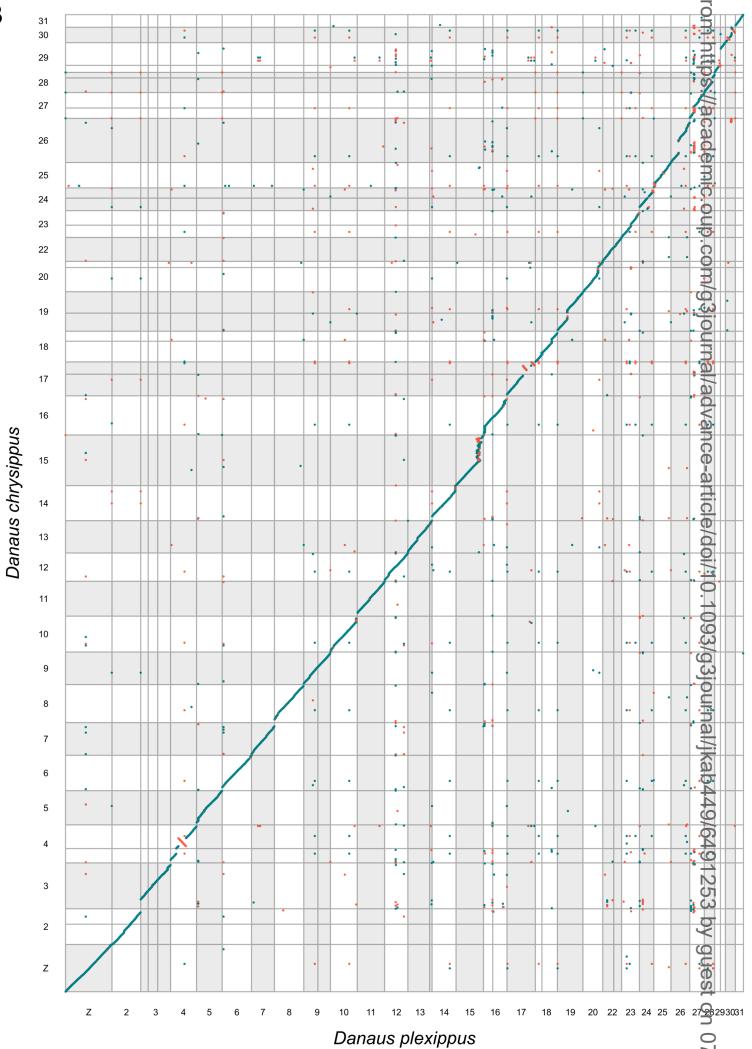
415
 416 **Figure 5.** Boxplots showing the distribution of (A) exon lengths and (B) intron lengths taken from the
 417 longest transcripts annotated with BRAKER2 for the three assemblies, *D. chrysippus* (Dchry2.2) and each
 418 of the re-annotated *D. plexippus* assemblies (MEX_DaPlex and Dplex_v4). Outlier points were omitted for
 419 clarity. Mean exon length is 226bp in Dchry2.2, 238bp for the re-annotated MEX_DaPlex assembly, and
 420 217bp for the re-annotated Dplex_v4 assembly. Mean intron length is 975bp in Dchry2.2, 665bp for the
 421 re-annotated MEX_DaPlex assembly, and 738bp for the re-annotated Dplex_v4 assembly.

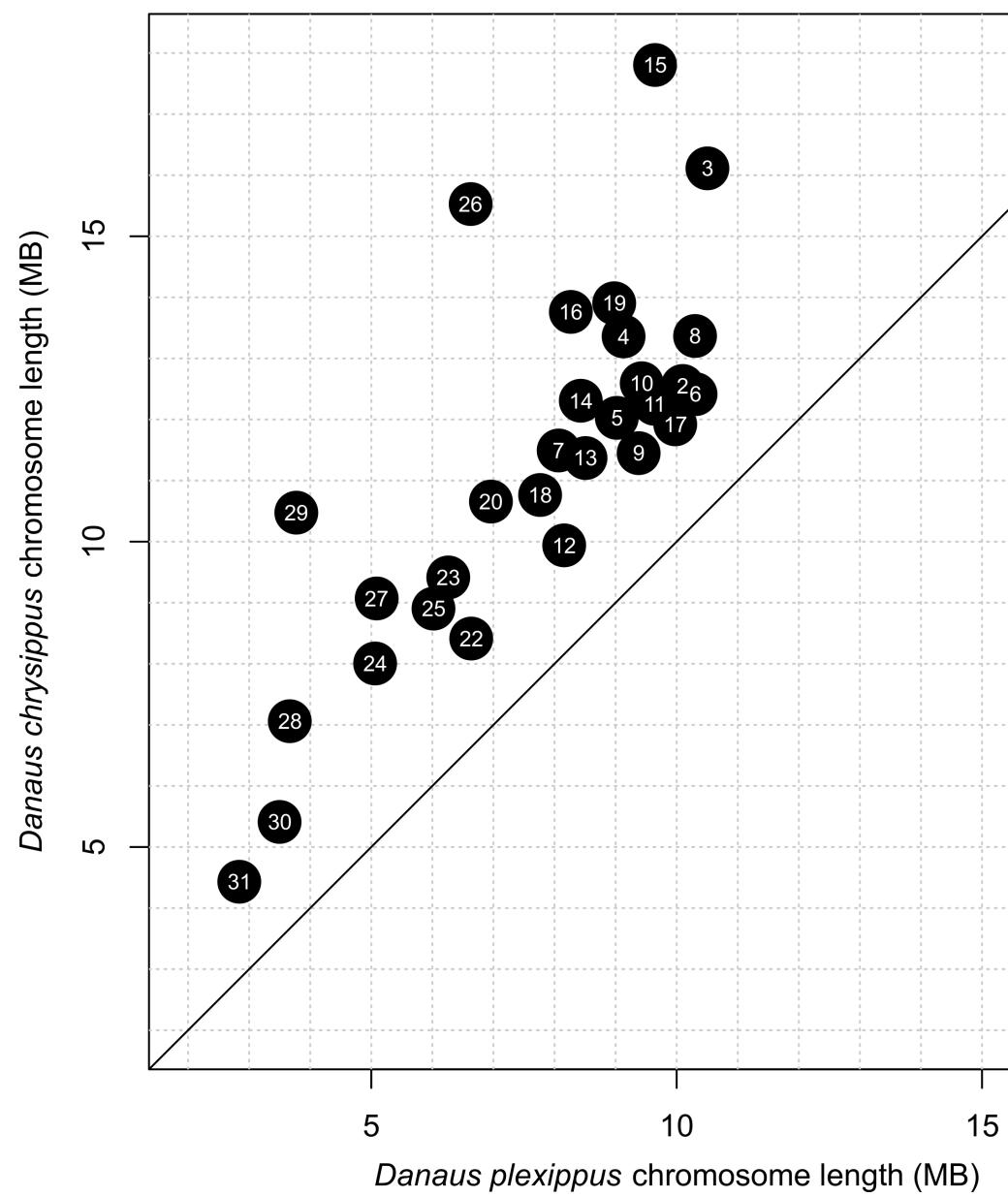
A

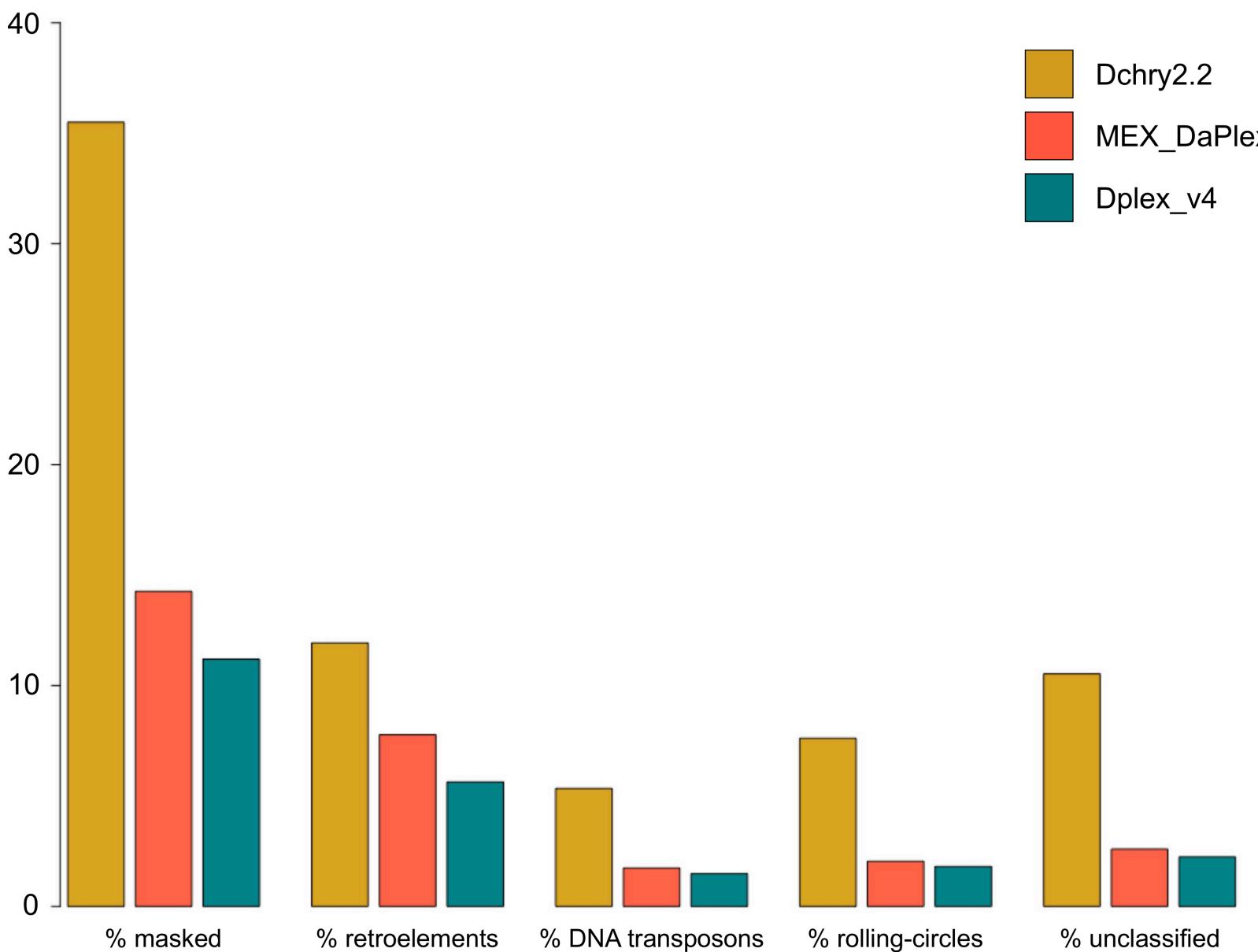
Danaus chrysippus



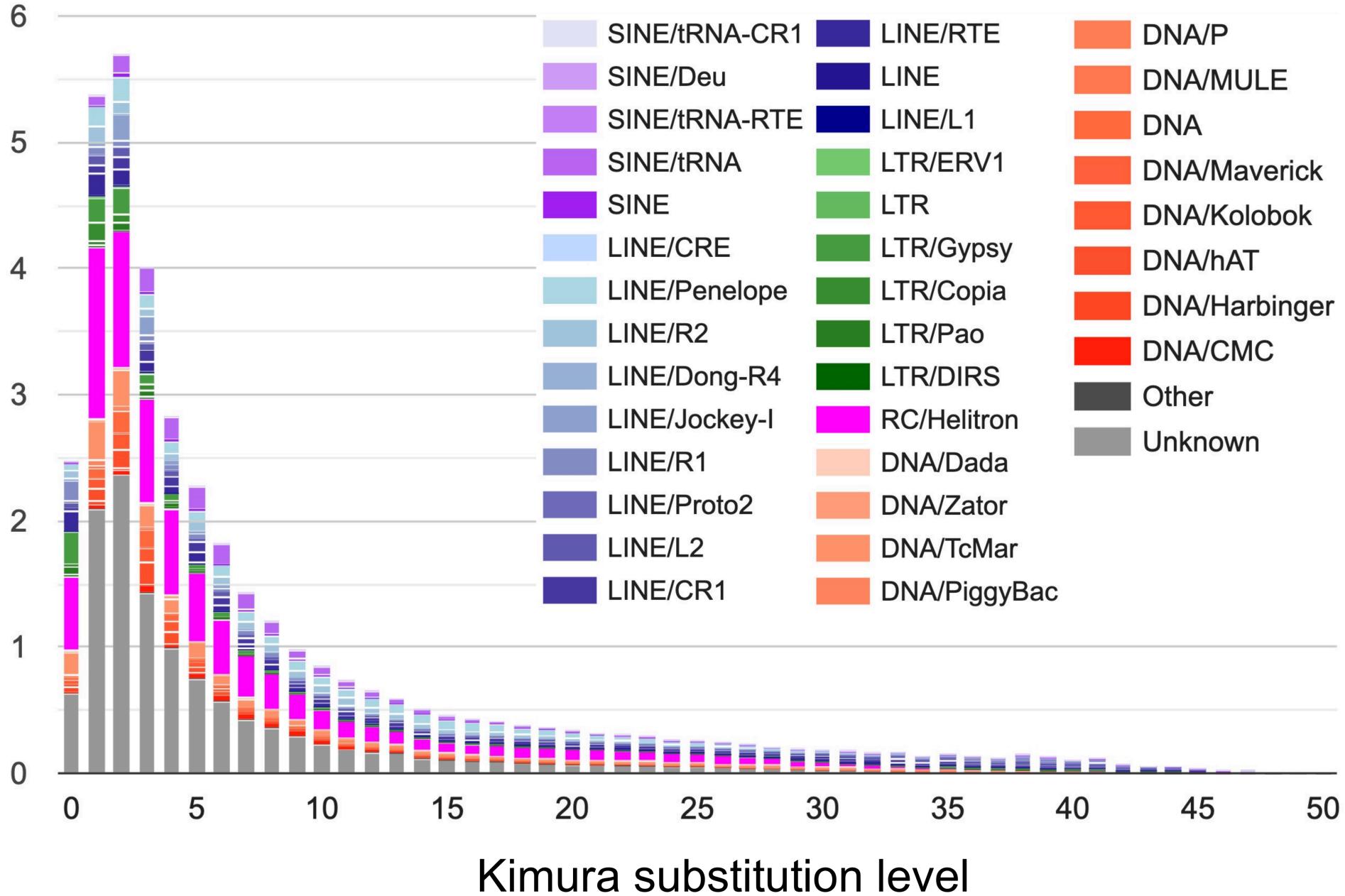
Danaus plexippus

**B**

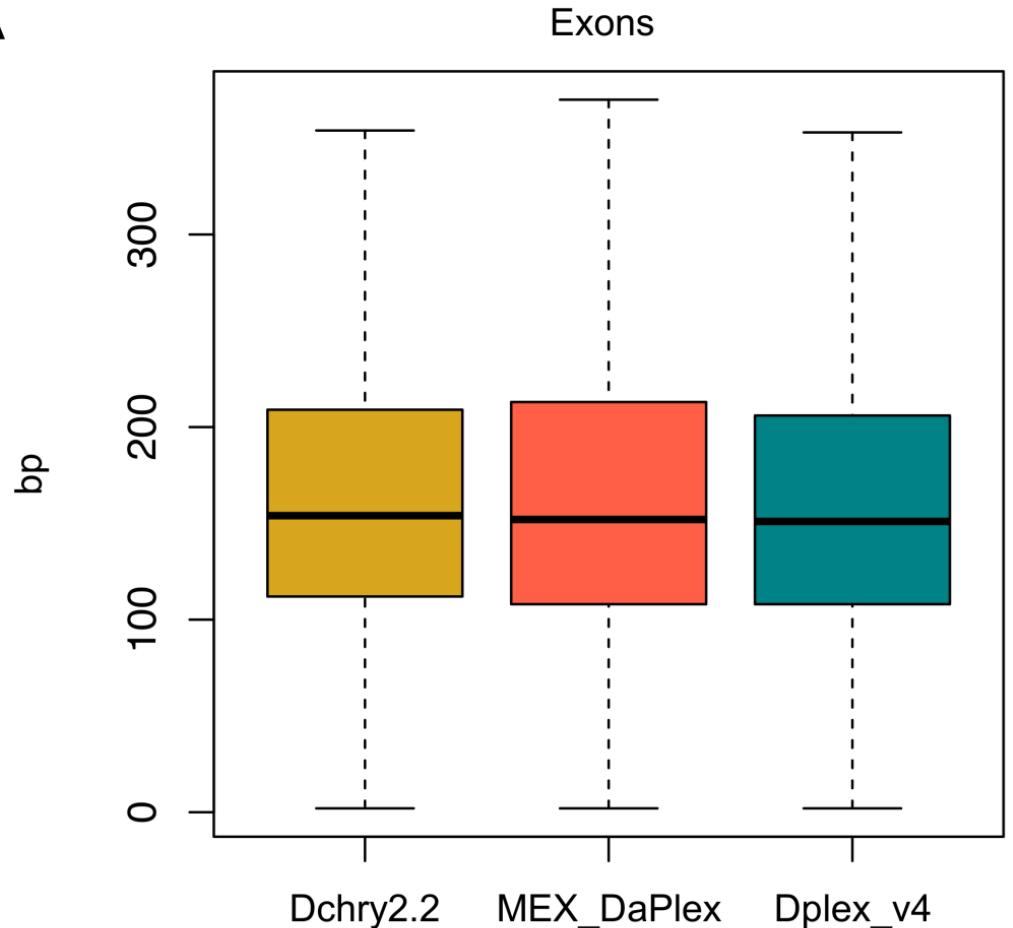




Percent of genome



A



B

