

1 A butterfly chromonome reveals selection dynamics during extensive and
2 cryptic chromosomal reshuffling

3

4

5 Authors

6

7 Jason Hill^{1,†}, Ramprasad Neethiraj¹, Pasi Rastas², Nathan Clark³, Nathan Morehouse⁴, Maria
8 Celorio¹, Jofre Carnicer Cols⁵, Heinrich Dircksen¹, Camille Meslin³, Kristin Sikkink⁶, Maria Vives⁵,
9 Heiko Vogel⁸, Christer Wiklund¹, Joel Kingsolver⁹, Carol Boggs⁷, Soren Nylin¹, Christopher Wheat^{1,†}

10

11 Affiliations:

12 1 Population Genetics, Department of Zoology, Stockholm University, Stockholm, Sweden

13 2 Department of Zoology, University of Cambridge, Cambridge, UK

14 3 Department of Computational and Systems Biology, University of Pittsburgh, Pittsburgh, USA

15 4 Department of Biological Sciences, University of Cincinnati, Cincinnati, USA

16 5 Departament d'Ecologia, Universitat de Barcelona (UB), Barcelona, 08028, Catalonia, Spain

17 6 Department of Ecology, Evolution and Behavior, University of Minnesota, St Paul, MN, USA

18 7 Department of Biological Sciences University of South Carolina, Columbia, SC, USA

19 8 Department of Entomology, Max Planck Institute for Chemical Ecology, D-07745 Jena, Germany

20 9 Department of Biology, University of North Carolina, Chapel Hill, North Carolina, USA

21

22 [†] Authors for correspondence:

23 Jason Hill <jason.hill@zoologi.su.se>

24 Christopher Wheat <chris.wheat@zoologi.su.se>

25

26 **Abstract**

27 Taxonomic Orders vary in their degree of chromosomal conservation with some having high
28 rates of chromosome number turnover despite maintaining some core sets of gene order (e.g.
29 Mammalia) and others exhibiting rapid rates of gene-order reshuffling without changing
30 chromosomal count (e.g. Diptera). However few clades exhibit as much conservation as the
31 Lepidoptera where both chromosomal count and gene collinearity (synteny) are very high over the
32 past 140 MY. In contrast, here we report extensive chromosomal rearrangements in the genome of
33 the green-veined white butterfly (*Pieris napi*, Pieridae, Linnaeus, 1758). This unprecedented
34 reshuffling is cryptic, micro-synteny and chromosome number do not indicate the extensive
35 rearrangement revealed by a chromosome level assembly and high resolution linkage map.
36 Furthermore, the rearrangement blocks themselves appear to be non-random, as they are
37 significantly enriched for clustered groups of functionally annotated genes revealing that the
38 evolutionary dynamics acting on Lepidopteran genome structure are more complex than previously
39 envisioned.

40 **Introduction**

41 The role of chromosomal rearrangements in adaptation and speciation has long been
42 appreciated and recent work has elevated the profile of supergenes in controlling complex adaptive
43 phenotypes¹⁻⁴. Chromosome number variation has also been cataloged for many species but analysis
44 of the adaptive implications have mostly been confined to the consequences of polyploidy and
45 whole genome duplication^{5,6}. The identification of pervasive fission and fusion events throughout
46 the genome is relatively unexplored since discovery of this pattern requires chromosome level
47 assemblies, leaving open the possibility of cryptic chromosomal dynamics taking place in many
48 species for which this level of genome assembly has not been achieved. As chromosomal levels
49 assemblies become more common, uncovering a relationship between such dynamics and
50 adaptation or speciation can be assessed.

51 Here we focus upon the Lepidoptera, the second most diverse animal group with over 160,000
52 extant species in more than 160 families. Butterflies and moths exist in nearly all habitats and have
53 equally varied life histories yet show striking similarity in genome architecture, with the vast
54 majority having a haploid chromosome number of $n=31^{7-9}$ (Ahola et al 2014; Lukhtanov, V. A. Sex
55 chromatin and sex chromosome systems in nonditrysian Lepidoptera (Insecta). J. Zool. Syst. Evol.
56 Res. 38, 73–79 (2000); Robinson R. Lepidopteran genetics (Pergamon Press, 1971)). While haploid
57 chromosome number can vary from $n = 5$ to $n = 223^{10-12}$, gene order and content is remarkably
58 similar within chromosomes (i.e. displays macro-synteny) and within these chromosomes the
59 degree of synteny between species separated by up to 140 My is astounding as illustrated by recent
60 chromosomal level genomic assemblies^{7,13}, as well as previous studies¹⁴⁻¹⁷. This ability of
61 Lepidoptera to accommodate such chromosomal rearrangements, yet maintain high levels of macro
62 and micro-synteny (i.e. collinearity at the scale of 10s to 100's of genes) is surprising. While a
63 growing body of evidence indicates that gene order in eukaryotes is non-random along
64 chromosomes, with upwards of 12% of genes organized into functional neighborhoods of shared
65 function and expression patterns¹⁸, to what extent this may play a role¹⁹ in the chromosomal
66 evolution of Lepidoptera is an open question.

67 Variation in patterns of synteny across clades must arise due to an evolutionary interaction between
68 selection and constraint²⁰, likely at the level of telomere and centromere performance. *Drosophila*,
69 and likely all Diptera, differ from the previously mentioned non-insect clades in lacking the
70 telomerase enzyme, and instead protect their chromosomal ends using retrotransposons²¹. This
71 absence of telomerase is posited to make evolving novel telomeric ends more challenging, limiting
72 the appearance of novel chromosomes and thereby resulting in high macro-synteny via constraint²².

73 In contrast, Lepidoptera like most Metazoans use telomerase to protect their chromosomal ends
74 which allows for previously internal chromosomal DNA to become subtelomeric in novel
75 chromosomes^{7,13}. Additionally all Lepidoptera have holocentric chromosomes in which the
76 decentralized kinetochore allows for more rearrangements by fission, fusion, and translocation of

77 chromosome fragments than monocentric chromosomes²³. Thus, Lepidoptera should be able to
78 avoid the deleterious consequences of large scale chromosomal changes.

79 Here we present the chromosome level genome assembly of the green-veined white butterfly (*Piers*
80 *napi*). Our analysis reveals large scale fission and fusion events similar to known dynamics in other
81 Lepidopteran species but at an accelerated rate and without a change in haploid chromosome count.
82 The resulting genome wide breakdown of the chromosome level synteny is unique among
83 Lepidoptera. While we are unable to identify any repeat elements associated with this cryptic
84 reshuffling, we find the chromosomal ends reused and the collinearity of functionally related genes.
85 These finding support a reinterpretation of the chromosomal fission dynamics in the Lepidoptera.

86

87 Results

88 The *P. napi* genome was generated using DNA from inbred siblings from Sweden, a genome
89 assembly using variable fragment size libraries (180 bp to 100 kb; N50-length of 4.2 Mb and a total
90 length of 350 Mb), and a high density linkage map across 275 full-sib larva, which placed 122
91 scaffolds into 25 linkage groups, consistent with previous karyotyping of *P. napi*^{24,25}. After
92 assessment and correction of the assembly, the total chromosome level assembly was 299 Mb
93 comprising 85% of the total assembly size and 114% of the k-mer estimated haploid genome size,
94 with 2943 scaffolds left unplaced (**Supplementary Note 3**). Subsequent annotation predicted
95 13,622 gene models, 9,346 with functional predictions (**Supplementary Note 4**).

96 Single copy orthologs (SCOs) in common between *P. napi* and the first Lepidopteran
97 genome, the silk moth *Bombyx mori*, were identified and revealed an unexpected deviation in gene
98 order and chromosomal structure in *P. napi* relative to *B. mori* as well as another lepidopteran
99 genome with a linkage map and known chromosomal structure *Heliconius melpomene* (Fig 1a).
100 Large scale rearrangements that appeared to be the fission and subsequent fusion of fragments in
101 the megabase scale were found to be present on every *P. napi* chromosome relative to *B. mori*, *H.*

102 *melpomene*, and *Meliteae cinxia* (fig 1b). We characterized the size and number of large scale
103 rearrangements between *P. napi* and *B. mori* using shared SCOs to identify 99 clearly defined
104 blocks of co-linear gene order (hereafter referred to as “syntenic blocks”), with each syntenic block
105 having an average of 69 SCOs. Each *P. napi* chromosome contained an average of 3.96 (SD = 1.67)
106 syntenic blocks, which derived from on average 3.5 different *B. mori* chromosomes. In *P. napi*, the
107 average syntenic block length was 2.82 Mb (SD = 1.97 Mb) and contained 264 genes (SD = 219).

108 The indication that *P. napi* diverged radically from the thus far observed chromosomal
109 structure of Lepidopterans raised questions about how common a *P. napi* like chromosomal
110 structure is observed vs. the structure reported in the highly syntenic *B. mori*, *H. melpomene*, and
111 *M. cinxia* genomes. We accessed 22 publicly available Lepidopteran genome assemblies
112 representing species diverged up to 140 MYA as well as their gene annotations to identify the genes
113 corresponding to the SCO’s used in previous analyses and blastx (Diamond v0.9.10)²⁶ to place those
114 genes on their native species scaffolds. With informations about each SCO’s location on the *P. napi*
115 chromosomes and the *B. mori* chromosomes we recorded how often a scaffold contained a cluster
116 of genes whose orthologs resided on two *P. napi* chromosomes or two *B. mori* chromosomes. If two
117 *P. napi* chromosomes were represented but only as single *B. mori* chromosome the scaffold was
118 marked as containing an mori-like join. Conversely if two *B. mori* chromosomes were represented
119 but only a single *P. napi* chromosome the scaffold was marked as containing a napi-like join. In
120 total we found for 20 species have more mori-like joins, and two species of *Pieris* represented by 3
121 assemblies have more napi-like joins (Fig 2a). While this type of assessment is noisy the indication
122 is that the genome structure described here is novel to the *Pieris* genus.

123 We validated this novel chromosomal reorganization using four complementary but
124 independent approaches to assess our scaffold joins. First, we generated a second linkage map for *P.*
125 *napi*, which confirmed the 25 linkage groups and the ordering of scaffold joins along chromosomes
126 (Fig. 3; Supplementary Fig. 2). Second, as the depth of the MP reads spanning joins indicated by the
127 first linkage map provides an independent assessment of the join validity, we quantified MP reads

128 spanning each base pair position along a chromosome (Fig. 3; Supplementary Fig. 2, Note 7),
129 finding strong support the scaffold joins. Third, we aligned the scaffolds of a recently constructed
130 genome of *P. rapae*²⁷ to *P. napi*, looking for *P. rapae* scaffolds that spanned the chromosomal level
131 scaffold joins within *P. napi*, finding support for 71 of the 97 joins (Supplementary Fig. 5). Fourth,
132 by considering *B. mori* syntenic blocks that spanned a scaffold join within a *P. napi* chromosome as
133 support for that *P. napi* chromosome assembly, we found that 62 of the 97 scaffold joins were
134 supported by *B. mori* (Supplementary Fig. 2, Note 8,9).

135 To assess this, we investigated the ordering and content of these syntenic blocks in *P. napi*.
136 First, we tested whether telomeric ends of chromosomes were at all conserved between species
137 despite the extensive chromosomal reshuffling (Fig. 4a). We found significantly more syntenic
138 blocks sharing telomere facing orientations between species than expected ($P < 0.01$, two tailed;
139 Fig. 4b). We also identified a significant enrichment for SCOs in *B. mori* and *P. napi* to be located at
140 roughly similar distance from the end of their respective chromosomes (Fig. 4c). Both of these
141 findings are consistent with the ongoing use of telomeric ends, indicating strong selection dynamics
142 acted upon their retention over evolutionary time. Second, we tested for gene set functional
143 enrichment within the observed syntenic blocks by investigating the full gene set of *P. napi* genes
144 within them. We found that 57 of the 99 block regions in the *P. napi* genome contained at least three
145 genes with a shared GO term that occurred with a $p < 0.01$ relative to the rest of the genome
146 (Supplementary fig. 3). We then tested whether the observed enrichment in the syntenic blocks of *P.*
147 *napi* was greater than expected by randomly assigning the genome into similarly sized blocks. The
148 mean number of GO enriched fragments in each of the simulated 10,000 genomes was 38.8
149 (variance of 46.6 and maximum of 52), which was significantly lower than observed ($P = 0$).

150 To assess the possible cause of the reshuffling, we surveyed the distribution of different
151 repeat element classes across the genome, looking for enrichment of specific categories near the
152 borders of syntenic blocks. While Class 1 transposons were found to be at higher density at near the
153 ends of chromosomes relative to the distribution internally (Supplementary fig. 4), no repeat

154 elements were enriched relative to the position of syntenic block regions. We therefore investigated
155 whether any repeat element classes had expanded within *Pieris* compared to other sequenced
156 genomes by assessing the distribution of repeat element classes and genome size among sequenced
157 Lepidoptera genomes. In accordance with other taxa²⁸ we find an expected strong relationship
158 between genome size and repetitive element content in *Pieris* species. Thus, while repetitive
159 elements such as transposable elements are likely involved in the reshuffling, our inability to find
160 clear elements involved suggests these events may be old and their signal decayed.

161 **Methods**

162 **Sample collection and DNA extraction.** Pupal DNA was isolated from a 4th generation inbred
163 cohort that originated from a wild caught female collected in Skåna, Sweden, using a standard salt
164 extraction²⁹.

165 **Illumina genome sequencing.** Illumina sequencing was used for all data generation used in
166 genome construction. A 180 paired end (PE) and the two mate pair (MP) libraries were constructed
167 at Science for Life Laboratory, the National Genomics Infrastructure, Sweden (SciLifeLab), using 1
168 PCR-free PE DNA library (180bp) and 2 Nextera MP libraries (3kb and 7kb) all from a single
169 individual. All sequencing was done on Illumina HiSeq 2500 High Output mode, PE 2x100bp by
170 Scilife. An additional two 40kb MP fosmid jumping libraries were constructed from a sibling used
171 in the previous library construction. Genomic DNA, isolated as above, was shipped to Lucigen Co.
172 (Middleton, WI, USA) for the fosmid jumping library construction and sequencing was performed
173 on an Illumina MiSeq using 2x250bp reads³⁰. Finally, a variable insert size libraries of 100 bp –
174 100,000 bp in length were generated using the Chicago and HiRise method³¹. Genomic DNA was
175 again isolated from a sibling of those used in previous library construction. The genomic DNA was
176 isolated as above and shipped to Dovetail Co. (Santa Cruz, CA, USA) for library construction,
177 sequencing and scaffolding. These library fragments were sequenced by Centrillion Biosciences
178 Inc. (Palo Alto, CA, USA) using Illumina HiSeq 2500 High Output mode, PE 2x100bp.

179 **Data Preparation and Genome assembly.** Nearly 500 M read pairs of data were generated,
180 providing ~ 285 X genomic coverage (Supplemental Table 1). The 3kb and 7kb MP pair libraries
181 were filtered for high confidence true mate pairs using Nextclip v0.8³². All read sets were then
182 quality filtered, the ends trimmed of adapters and low quality bases, and screened of common
183 contaminants using bbdduk v37.51 (bbtools, Brian Bushnell). Insert size distributions were plotted to
184 assess library quality, which was high (Supplementary Fig. 1). The 180bp, 3kb, and 7kb, read data
185 sets were used with AllpathsLG r50960³³ for initial contig generation and scaffolding
186 (Supplementary Note 1). AllpathsLG was run with haploidify = true to compensate for the high
187 degree of heterozygosity. Initial contig assembly's conserved single copy ortholog content was
188 assessed at 78% for *P. napi* by CEGMA v2.5³⁴. A further round of superscaffolding using the 40kb
189 libraries alongside the 3kb and 7kb libraries was done using SSPACE v2³⁵. Finally, both assemblies
190 were ultrascaffolded using the Chicago read libraries and the HiRise software pipeline. These steps
191 produced a final assembly of 3005 scaffolds with an N50-length of 4.2 Mb and a total length of 350
192 Mb (Supplementary Note 1).

193 **Linkage Map.** RAD-seq data of 5463 SNP markers from 275 full-sib individuals, without parents,
194 was used as input into Lep-MAP2³⁶. The RAD-seq data was generated from next-RAD technology
195 by SNPsaurus (Oregon, USA)(Supplemental note 10). To obtain genotype data, the RAD-seq data
196 was mapped to the reference genome using BWA mem³⁷ and SAMtools³⁸ was used to produce
197 sorted bam files of the read mappings. Based on read coverage (samtools depth), Z chromosomal
198 regions were identified from the genome and the sex of offspring was determined. Custom scripts³⁹
199 were used to produce genotype posteriors from the output of SAMtools mpileup.

200 The parental genotypes were inferred with Lep-MAP2 ParentCall module using parameters
201 "ZLimit=2 and ignoreParentOrder=1", first calling Z markers and second calling the parental
202 genotypes by ignoring which way the parents are informative (the parents were not genotyped so
203 we could not separate maternal and paternal markers at this stage). Scripts provided with Lep-
204 MAP2 were used to produce linkage file from the output of ParentCall and all single parent

205 informative markers were converted to paternally informative markers by swapping parents, when
206 necessary. Also filtering by segregation distortion was performed using Filtering module.
207 Following this, the SeparateChromosomes module was run on the linkage file and 25 chromosomes
208 were identified using LOD score limit 39. Then JoinSingles module was run twice to add more
209 markers on the chromosomes with LOD score limit of 20. Then SeparateChromosomes was run
210 again but only on markers informative on single parent with LOD limit 10 to separate paternally
211 and maternally informative markers. 51 linkage groups were found and all were ordered using
212 OrderMarkers module. Based on likelihood improvement of marker ordering, paternal and maternal
213 linkage groups were determined. This was possible as there is no recombination in female
214 (achiasmatic meiosis), thus order of the markers does not improve likelihood on the female map.
215 The markers on the corresponding maternal linkage groups were converted to maternally
216 informative and OrderMarkers was run on the resulting data twice for each of 25 chromosomes
217 (without allowing recombination on female). The final marker order was obtained as the order with
218 higher likelihood of the two runs.

219

220 **Chromosomal assembly.** The 5463 markers that composed the linkage map were mapped to the *P.*
221 *napi* ultrascaffolds using bbmap⁴⁰ with sensitivity = slow. Reads that mapped uniquely were used to
222 identify misassemblies in the ultra-scaffolds and arrange those fragments into chromosomal order.
223 54 misassemblies were identified and overall 115 fragments were joined together into 25
224 chromosomes using a series of custom R scripts (supplemental information) and the R package
225 Biostrings⁴¹. Scaffold joins and misassembly corrections were validated by comparing the number
226 of correctly mapped mate pairs spanning a join between two scaffolds. Mate pair reads from the
227 3kb, 7kb, and 40kb libraries were mapped to their respective assemblies with bbmap (po=t,
228 ambig=toss, kbp=t). SAM output was filtered for quality and a custom script was used to tabulate
229 read spanning counts for each base pair in the assembly.

230 **Synteny Comparisons Between *P. napi*, *B. mori*, and *H. melpomene*.** A list of 3100 single copy
231 orthologs (SCO) occurring in the Lepidoptera lineage curated by OrthoDB v9.1⁴² was used to
232 extract gene names and protein sequences of SCOs in *Bombyx mori* from
233 KaikoBase⁴³ (Supplemental Note 5) using a custom script. Reciprocal best hits (RBH) between gene
234 sets of *Pieris napi*, *Pieris rapae*, *Heliconius melpomene*, *Melitea Cinxia*, and *Bombyx mori* SCOs
235 were identified using BLASTP⁴⁴ and custom scripts. Gene sets of *H. melpomene* v2.5 and *M.*
236 *cinxia* v1 were downloaded from LepBase v4⁴⁵. Coordinates were converted to chromosomal
237 locations and visualized using Circos⁴⁶ and custom R scripts.

238 **Synteny Comparison Within Lepidoptera.** Genome assemblies and annotated protein sets were
239 downloaded for 24 species of Lepidoptera from LepBase v4⁴⁷ and other sources (Supplemental
240 Table 4). Each target species protein set was aligned to its species genome as well as to the *Pieris*
241 *napi* protein set using Diamond v0.9.10²⁶ with default options. The protein-genome comparison was
242 used to assign each target species gene to one of its assembled scaffolds, while the protein-protein
243 comparison was used to identify RBHs between the protein of each species and its ortholog in *P.*
244 *napi*, and *B. mori*. Using this information we used a custom R script to examine each assembly
245 scaffold for evidence of synteny to either *P. napi* or *B. mori*. First, each scaffold of the target species
246 genome was assigned genes based on the protein-genome blast results, using its own protein set and
247 genome. A gene was assigned to a scaffold if at least 3 HSPs of less than 200bp from a gene aligned
248 with $\geq 95\%$ identity. Second, if any of these scaffolds then contained genes whose orthologs
249 resided on a single *B. mori* chromosome but two *P. napi* chromosomes, and those same two *P. napi*
250 chromosome segments were also joined in the *B. mori* assembly, that was counted as a ‘mori-like
251 join’. Conversely if a target species scaffold contained genes whose orthologs resided on a single *P.*
252 *napi* chromosome but two *B. mori* chromosomes, and those same two *B. mori* chromosome
253 segments were also joined in the *P. napi* assembly, that was counted as a ‘napi-like join’.

254

255 **Pieridae chromosomal evolution.**

256 Reconstruction of the chromosomal fusions and fissions were estimated across the family Pieridae
257 by placing previously published karyotype studies of haploid chromosomal counts into their
258 evolutionary context. There are approximately 1000 species in the 85 recognized genera of Pieridae
259 and recently we reconstructed a robust fossil-calibrated chronogram for this family at the genus
260 level^{48,49}. Upon this time calibrated phylogeny we then placed the published chromosomal counts
261 for 201 species^{9,50}, with ancestral chromosomal reconstructions for chromosome count, treated as a
262 continuous character, used the contMap function of the phytools R package⁵¹

263 **Second Linkage Map for *P. napi*.** A second linkage map was constructed from a different family
264 of *P. napi* in which a female from Abisko, Sweden was crossed with a male from Catalonia, Spain.
265 Genomic DNA libraries were constructed for the mother, father, and four offspring (2 males, 2
266 females). RNA libraries were constructed for an additional 6 female and 6 male offspring. All
267 sequencing was performed on a Illumina HiSeq 2500 platform using High Output mode, with PE
268 2x100bp reads at SciLifeLab (Stockholm, Sweden). Both DNA and RNA reads were mapped to the
269 genome assembly with bbmap. Samtools was used to sort read mappings and merge them into an
270 mpileup file (Supplemental Note 6). Variants were called with BCFtools⁵² and filtered with
271 VCFtools⁵³. Linkage between SNPs was assessed with PLINK⁵⁴. A custom script was used to assess
272 marker density and determine sex-specific heterozygosity.

273

274 **Annotation of *Pieris napi* genome.** Genome annotation was carried out by the Bioinformatics
275 Short-term Support and Infrastructure (BILS, Sweden). BILS was provided with the chromosomal
276 assembly of *P. napi* and 45 RNAseq read sets representing 3 different tissues (head, fat body, and
277 gut) of 7 male and 8 female larva from lab lines separate from the one used for the initial
278 sequencing. Sequence evidence for the annotation was collected in two complementary ways. First,
279 we queried the Uniprot database⁵⁵ for protein sequences belonging to the taxonomic group of
280 Papilionoidea (2,516 proteins). In order to be included, proteins gathered in this way had to be
281 supported on the level of either proteomics or transcriptomics and could not be fragments. In

282 addition, we downloaded the Uniprot-Swissprot reference data set (downloaded on 2014-05-15)
283 (545,388 proteins) for a wider taxonomic coverage with high-confidence proteins. In addition, 493
284 proteins were used that derived from a *P. rapae* expressed sequence tag library that was Sanger
285 sequenced.

286

287 **Permutation test of syntenic block position within chromosomes.** Syntenic blocks (SBs)
288 identified as interior vs terminal and the ends of terminal blocks were marked as inward or outward
289 facing. SBs were reshuffled into 25 random chromosomes of 4 SBs in a random orientation and the
290 number of times that a terminal block occurred in a random chromosome with the outward end
291 facing outward was counted. This was repeated 10,000 times to generate a random distribution
292 expectation. The number of terminal outward facing SBs in *B. mori* that were also terminal and
293 outward facing in *P. napi* was compared to this random distribution to derive significance of
294 deviation from the expected value. To test the randomness of gene location within chromosomes,
295 orthologs were numbered by their position along each chromosome in both *B. mori* and *P. napi*.
296 10,000 random genomes were generated as above. Distance from the end of the new chromosome
297 and distance from the end of *B. mori* chromosome was calculated for each ortholog and the result
298 binned. P-values were determined by comparing the number of orthologs in a bin to the expected
299 distribution of genes in a bin from the random genomes. All test were done using a custom R script.

300 **Gene set enrichment analysis of syntenic blocks.** Gene ontology set enrichment was initially
301 tested for within syntenic blocks of the *P. napi* genome using topGO⁵⁶ with all 13,622 gene models
302 generated from the annotation. For each syntenic block within the genome, each GO term of any
303 level within the hierarchy that had at least 3 genes belonging to it was analyzed for enrichment. If a
304 GO term was overrepresented in a syntenic block compared to the rest of the genome at a p-value of
305 < 0.01 by a Fisher exact test, that block was counted as enriched. 57 of the 99 syntenic blocks in the
306 *P. napi* genome were enriched in this way. Because arbitrarily breaking up a genome and testing for
307 GO enrichment can yield results that are dependent on the distribution of the sizes used, we

308 compared the results of the previous analysis to the enrichment found using the same size genomic
309 regions, randomly selected from the *P. napi* genomes. The size distribution of the 99 syntenic
310 blocks were used to generate fragment sizes into which the genome was randomly assigned. This
311 resulted in a random genome of 99 fragments which in total contained the entire genome but the
312 content of a given fragment was random compared to the syntenic block that defined its size. This
313 random genome was tested for GO enrichment of the fragments in the same way as the syntenic
314 blocks in the original genome, and the number of enriched blocks counted. This was then repeated
315 10,000 times to generate a distribution of expected enrichment in genome fragments of the same
316 size as the *P. napi* syntenic blocks.

317

318 Discussion

319 While massive chromosomal fission events are well documented in butterflies (e.g.
320 *Leptidea* in Pieridae (n=28-103); *Agrodiaetus* in Lycaenidae (n=10-134)), their contribution to
321 Lepidopteran diversity appears to be minimal as all these clades are very young⁵⁷⁻⁵⁹. However, our
322 results challenge this interpretation. Rather, *P. napi* appears to represent a lineage that has
323 undergone an impressive reconciliation of an earlier series of rampant fission events. Moreover, the
324 subsequent fusion events exhibit a clear bias toward using ancient telomeric ends, as well as
325 returning gene clusters to their relative ancestral position within chromosomes even when the other
326 parts of the newly formed chromosome originated from other sources. Luckily these initial fission
327 events have been frozen in time as the reshuffled syntenic blocks, revealing the potential fitness
328 advantage of maintaining certain functional categories as syntenic blocks.

329 Thus, despite the potential for holocentric species to have relaxed constraint upon their
330 chromosomal evolution, we find evidence for selection actively maintaining ancient telomeric ends,
331 as well as gene order within large chromosomal segments. Together these observations suggest that

332 the low chromosome divergence in Lepidoptera over > 100 million generations is at least partially
333 due to purifying selection maintaining an adaptive chromosomal structure.

334

335 **Bibliography**

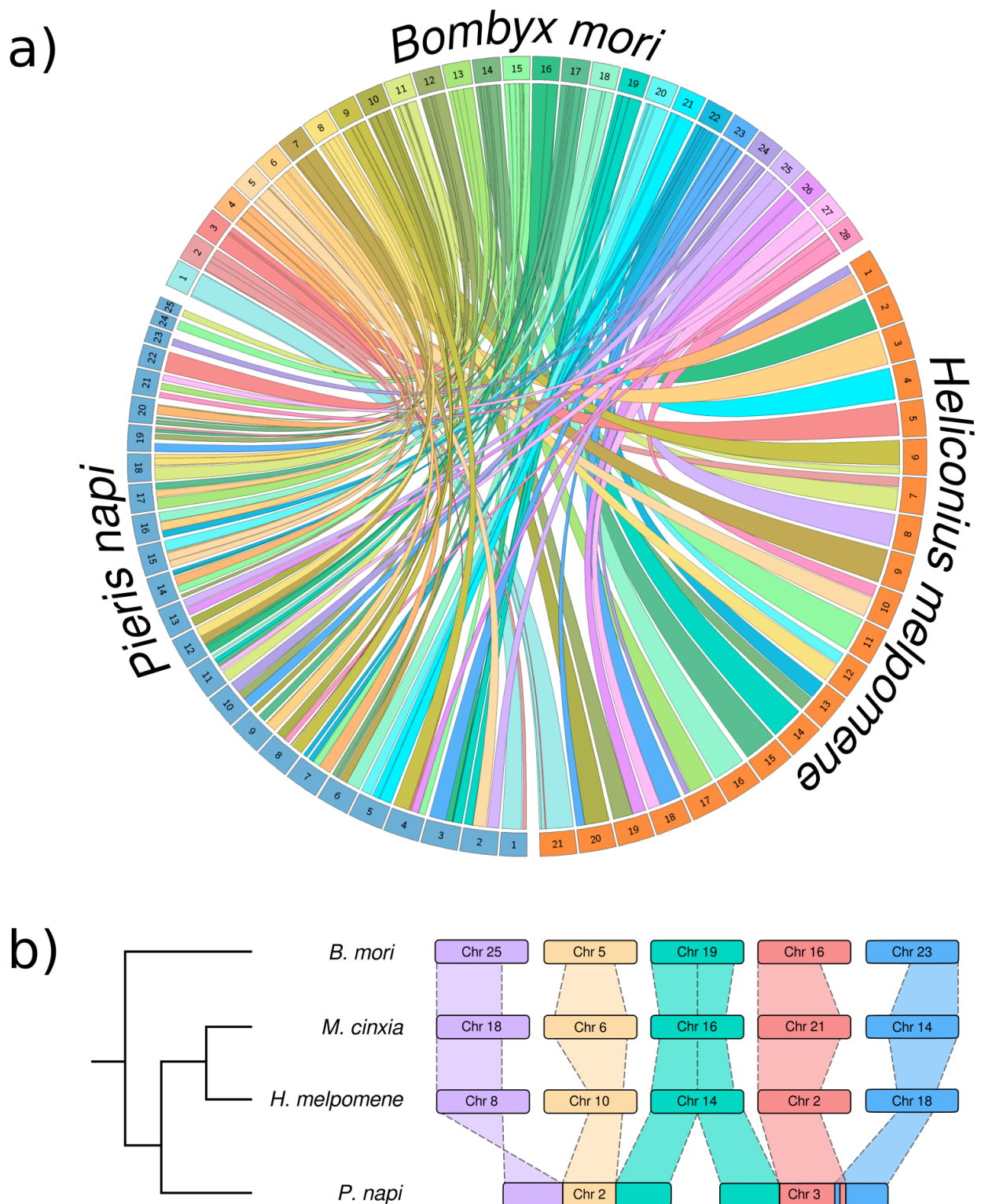
336

- 337 1. Schwander, T., Libbrecht, R. & Keller, L. Supergenes and complex phenotypes. *Curr. Biol.*
338 **24**, R288–R294 (2014).
- 339 2. Kunte, K. *et al.* Doublesex Is a Mimicry Supergene. *Nature* **507**, 229–232 (2014).
- 340 3. Fishman, L., Stathos, A., Beardsley, P. M., Williams, C. F. & Hill, J. P. Chromosomal
341 rearrangements and the genetics of reproductive barriers in mimulus (monkey flowers).
342 *Evolution (N. Y.)*. **67**, 2547–2560 (2013).
- 343 4. Lamichhaney, S. *et al.* Structural genomic changes underlie alternative reproductive
344 strategies in the ruff (*Philomachus pugnax*). *Nat. Genet.* **48**, 84–88 (2015).
- 345 5. Otto, S. P. & Whitton, J. Polyploid Incidence and Evolution. *Annu. Rev. Genet.* **34**, 401–437
346 (2000).
- 347 6. Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of polyploidy.
348 *Nat. Rev. Genet.* **18**, 411–424 (2017).
- 349 7. Ahola, V. *et al.* The Glanville fritillary genome retains an ancient karyotype and reveals
350 selective chromosomal fusions in Lepidoptera. *Nat Commun* **5**, 1–9 (2014).
- 351 8. Lukhtanov, V. A. Sex chromatin and sex chromosome systems in nonditrysian Lepidoptera
352 (Insecta). *J. Zool. Syst. Evol. Res.* **38**, 73–79 (2000).
- 353 9. ROBINSON, R. *Lepidoptera Genetics*. *Lepidoptera Genetics* (1971). doi:10.1016/B978-0-
354 08-006659-2.50017-1
- 355 10. Brown, K. S., Von Schoultz, B. & Suomalainen, E. Chromosome evolution in Neotropical
356 Danainae and Ithomiinae (Lepidoptera). *Hereditas* **141**, 216–236 (2004).
- 357 11. Kandul, N. P., Lukhtanov, V. A. & Pierce, N. E. Karyotypic diversity and speciation in
358 *Agrodiaetus* butterflies. *Evolution (N. Y.)*. **61**, 546–559 (2007).
- 359 12. Saura, A., Schoultz, B. Von, Saura, A. O. & Brown, K. S. Chromosome evolution in
360 Neotropical butterflies. *Hereditas* **150**, 26–37 (2013).
- 361 13. Davey, J. W. *et al.* Major Improvements to the *Heliconius melpomene* Genome Assembly
362 Used to Confirm 10 Chromosome Fusion Events in 6 Million Years of Butterfly Evolution.
363 *Genes|Genomes|Genetics* **6**, 695–708 (2016).
- 364 14. Yasukochi, Y. A Second-Generation Integrated Map of the Silkworm Reveals Synteny and
365 Conserved Gene Order Between Lepidopteran Insects. *Genetics* **173**, 1319–1328 (2006).

- 366 15. Yasukochi, Y. *et al.* A FISH-based chromosome map for the European corn borer yields
367 insights into ancient chromosomal fusions in the silkworm. *Heredity (Edinb)*. **116**, 75–83
368 (2016).
- 369 16. Beldade, P., Saenko, S. V, Pul, N. & Long, A. D. A gene-based linkage map for *Bicyclus*
370 anynana butterflies allows for a comprehensive analysis of synteny with the lepidopteran
371 reference genome. *PLoS Genet*. **5**, e1000366 (2009).
- 372 17. Šíchová, J. *et al.* Fissions, fusions, and translocations shaped the karyotype and multiple sex
373 chromosome constitution of the northeast-Asian wood white butterfly, *Leptidea amurensis*.
374 *Biol. J. Linn. Soc.* **118**, 457–471 (2016).
- 375 18. Al-Shahrour, F. *et al.* Selection upon genome architecture: Conservation of functional
376 neighborhoods with changing genes. *PLoS Comput. Biol.* **6**, (2010).
- 377 19. Gordon, J. L., Byrne, K. P. & Wolfe, K. H. Mechanisms of chromosome number evolution in
378 yeast. *PLoS Genet*. **7**, 0–3 (2011).
- 379 20. Coghlan, A., Eichler, E. E., Oliver, S. G., Paterson, A. H. & Stein, L. Chromosome evolution
380 in eukaryotes: A multi-kingdom perspective. *Trends Genet*. **21**, 673–682 (2005).
- 381 21. Levis, R. W., Ganesan, R., Houtchens, K., Tolar, L. A. & Sheen, F. miin. Transposons in
382 place of telomeric repeats at a *Drosophila* telomere. *Cell* **75**, 1083–1093 (1993).
- 383 22. Muller, H. J. The remaking of chromosomes. *Collect. net* **13**, 181–198 (1938).
- 384 23. Maddox, P. S., Oegema, K., Desai, A. & Cheeseman, I. M. ‘Holo’er than thou: Chromosome
385 segregation and kinetochore function in *C. elegans*. *Chromosom. Res.* **12**, 641–653 (2004).
- 386 24. Lorkovic, Z. The genetics and reproductive isolating mechanisms of the *Pieris napi* -
387 *bryoniae* group. *J. Lepid. Soc.* **16**, 5–19 (1955).
- 388 25. Maeki, K. & Kawazoe, A. On the Hybridization between Two Karyotype Lineages of *Pieris*
389 *napi* Linnaeus from Japan. *Cytologia (Tokyo)*. **59**, (1994).
- 390 26. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
391 DIAMOND. *Nat. Methods* **12**, 59–60 (2014).
- 392 27. Nallu, S. *et al.* The Molecular Genetic Basis of Herbivory between Butterflies and their Host-
393 Plants. *bioRxiv* (2017). doi:10.1101/154799
- 394 28. Clark, A. G. *et al.* Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* **450**,
395 203–218 (2007).
- 396 29. Aljanabi, S. M. & Martinez, I. Universal and rapid salt-extraction of high quality genomic
397 DNA for PCR-based techniques. *Nucleic Acids Res.* **25**, 4692–4693 (1997).
- 398 30. Bentley, D. R. *et al.* Accurate whole human genome sequencing using reversible terminator
399 chemistry. *Nature* **456**, 53–59 (2008).
- 400 31. Putnam, N. H. *et al.* Chromosome-scale shotgun assembly using an in vitro method for long-
401 range linkage. *Genome Res.* **26**, 342–350 (2016).

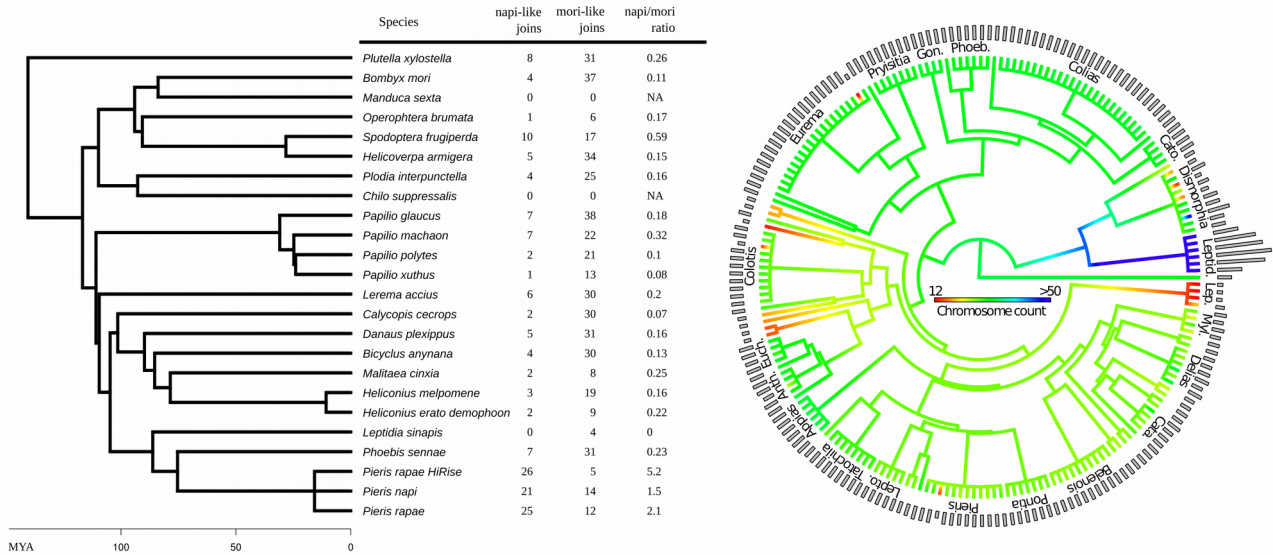
- 402 32. Leggett, R. M., Clavijo, B. J., Clissold, L., Clark, M. D. & Caccamo, M. NextClip: an
403 analysis and read preparation tool for Nextera Long Mate Pair libraries. *Bioinformatics* **30**,
404 566–8 (2014).
- 405 33. Gnerre, S. *et al.* High-quality draft assemblies of mammalian genomes from massively
406 parallel sequence data. *Proc. Natl. Acad. Sci.* **108**, 1513–1518 (2011).
- 407 34. Parra, G., Bradnam, K. & Korf, I. CEGMA: a pipeline to accurately annotate core genes in
408 eukaryotic genomes. *Bioinformatics* **23**, 1061–1067 (2007).
- 409 35. Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-
410 assembled contigs using SSPACE. *Bioinformatics* **27**, 578–579 (2011).
- 411 36. Rastas, P., Calboli, F. C. F., Guo, B., Shikano, T. & Merilä, J. Construction of Ultradense
412 Linkage Maps with Lep-MAP2: Stickleback F2 Recombinant Crosses as an Example.
413 *Genome Biol. Evol.* **8**, 78–93 (2015).
- 414 37. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. **0**,
415 1–3 (2013).
- 416 38. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–
417 2079 (2009).
- 418 39. Kvist, J. *et al.* Flight-induced changes in gene expression in the Glanville fritillary butterfly.
419 *Mol. Ecol.* **24**, 4886–4900 (2015).
- 420 40. Bushnell, B. BBTools. (2017). at <<https://jgi.doe.gov/data-and-tools/bbtools/>>
- 421 41. Pages H, Gentleman R, Aboyoun P, *et al.* Biostrings: String objects representing biological
422 sequences, and matching algorithms. *R Packag. version 2*, 2008 (2008).
- 423 42. Zdobnov, E. M. *et al.* OrthoDB v9.1: Cataloging evolutionary and functional annotations for
424 animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Res.* **45**, D744–
425 D749 (2017).
- 426 43. Shimomura, M. *et al.* KAIKObase: an integrated silkworm genome database and data mining
427 tool. *BMC Genomics* **10**, 486 (2009).
- 428 44. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database
429 search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
- 430 45. Kersey, P. J. *et al.* Ensembl Genomes 2016: More genomes, more complexity. *Nucleic Acids*
431 *Res.* **44**, D574–D580 (2016).
- 432 46. Krzywinski, M. *et al.* Circos: An information aesthetic for comparative genomics. *Genome*
433 *Res.* **19**, 1639–1645 (2009).
- 434 47. Challis, R. J., Kumar, S., Dasmahapatra, K. K., Jiggins, C. D. & Blaxter, M. Lepbase: The
435 Lepidopteran genome database. *bioRxiv* 56994 (2016). doi:10.1101/056994
- 436 48. Wahlberg, N., Rota, J., Braby, M. F., Pierce, N. E. & Wheat, C. W. Revised systematics and
437 higher classification of pierid butterflies (Lepidoptera: Pieridae) based on molecular data.
438 *Zool. Scr.* **43**, 641–650 (2014).

- 439 49. Edger, P. P. *et al.* The butterfly plant arms-race escalated by gene and genome duplications.
440 *Proc. Natl. Acad. Sci.* **112**, 8362–8366 (2015).
- 441 50. Lukhtanov, V. A. Karyotype evolution and systematics of higher taxa of Pieridae
442 (Lepidoptera) of the World. *Ent. Obozr.* **70** 619–636, 3 figs (1991).
- 443 51. Revell, L. J. phytools: An R package for phylogenetic comparative biology (and other
444 things). *Methods Ecol. Evol.* **3**, 217–223 (2012).
- 445 52. Narasimhan, V. *et al.* BCFtools/RoH: A hidden Markov model approach for detecting
446 autozygosity from next-generation sequencing data. *Bioinformatics* **32**, 1749–1751 (2016).
- 447 53. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158
448 (2011).
- 449 54. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based
450 Linkage Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 451 55. The UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.* **43**,
452 D204–12 (2015).
- 453 56. Alexa A and Rahnenfuhrer J. topGO: Enrichment Analysis for Gene Ontology. *R package*
454 *version 2.26.0.* (2016). at <<http://bioconductor.org/packages/release/bioc/html/topGO.html>>
- 455 57. Dincă, V., Lukhtanov, V. A., Talavera, G. & Vila, R. Unexpected layers of cryptic diversity in
456 wood white Leptidea butterflies. *Nat. Commun.* **2**, (2011).
- 457 58. Vila, R. *et al.* Phylogeny and palaeoecology of *Polyommatus* blue butterflies show Beringia
458 was a climate-regulated gateway to the New World. *Proc. R. Soc. B Biol. Sci.* **278**, 2737–
459 2744 (2011).
- 460 59. Lukhtanov, V. The blue butterfly *Polyommatus* (*Plebicula*) *atlanticus* (Lepidoptera,
461 Lycaenidae) holds the record of the highest number of chromosomes in the non-polyploid
462 eukaryotic organisms. *Comp. Cytogenet.* **9**, 683–690 (2015).

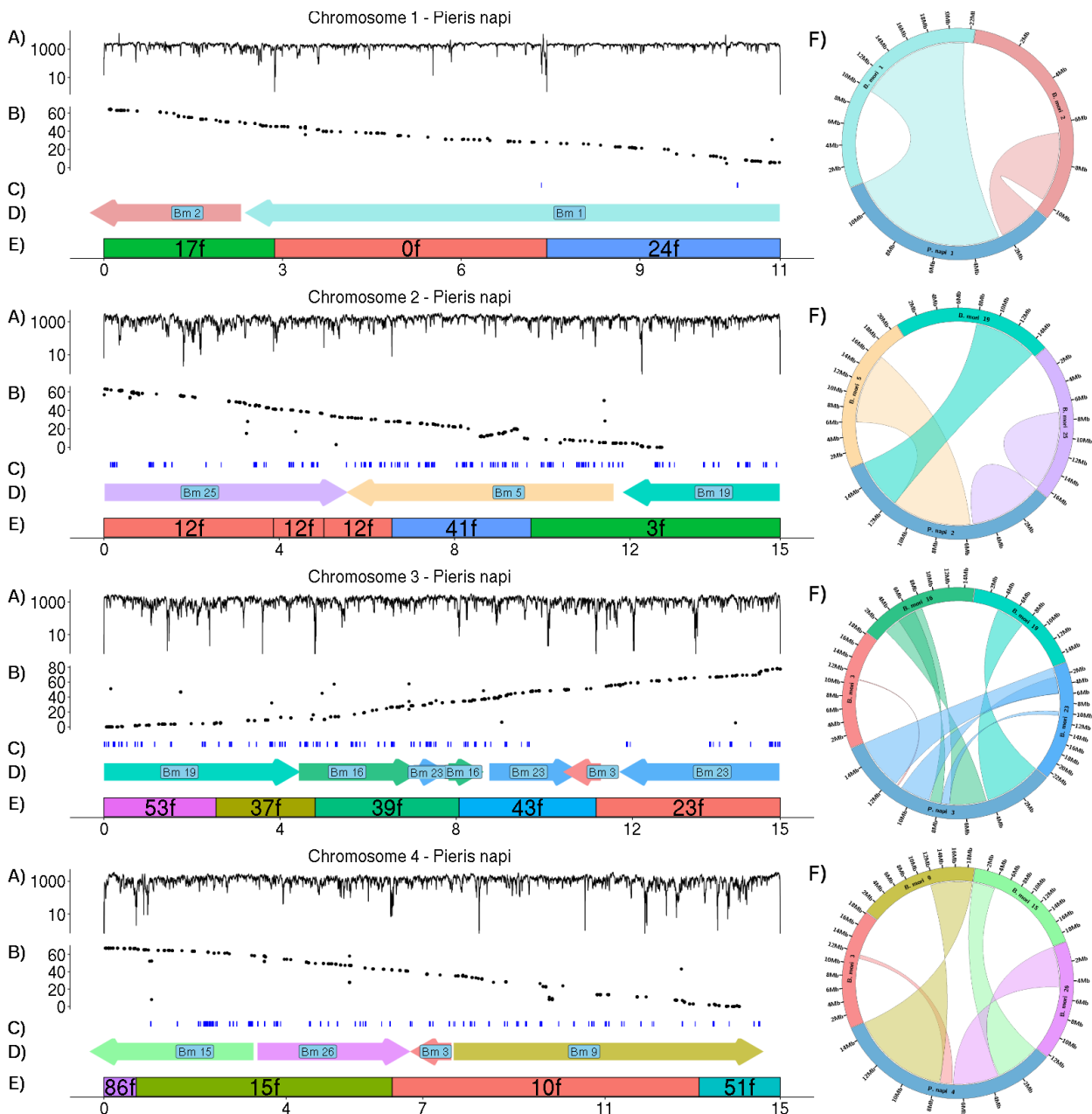


464 **Figure 1 a)** Chromosomal mapping between the moth *Bombyx mori* (Bombycoidea) and the
 465 butterflies *Pieris napi* (Pieridae) and *Heliconius melpomene* (Nymphalidae). These species last
 466 shared a common ancestor > 100 million generations ago⁴⁹. Depicted are the reciprocal best hit
 467 orthologs identified between *B. mori* and *P. napi* (n=2354) and between *B. mori* and *H. melpomene*
 468 (n=2771). Chromosome 1 is the Z chromosome in *B. mori* and *P. napi* and 21 is the Z chromosome
 469 in *H. melpomene*. Chromosomes 2-25 in *P. napi* are ordered in size from largest to smallest. Links
 470 between orthologs originate from the *B. mori* chromosome and are colored by their chromosome of
 471 origin, while *P. napi* chromosomes are colored blue and *H. melpomene* chromosomes are colored

472 orange. Links are clustered into blocks of synteny and each ribbon represents a contiguous block of
 473 genes spanning a region in both species. **b)** Two largest autosomes of *P. napi* and their synteny to
 474 other Lepidoptera and their phylogenetic relationship. The sister taxa and the more distant *B. mori*
 475 share a high degree of macro synteny while the *P. napi* genome required multiple chromosomal
 476 fusion and fission events to be patterned in the way that is observed. Band width for each species is
 477 proportional to the length of the inferred chromosomal region of ornithology, although the
 478 individual chromosomes are not to scale.
 479



480 **Figure 2 a)** A time calibrated phylogeny of currently available Lepidopteran genomes (n=24) and
 481 estimates their macrosynteny with *B. mori* and *P. napi*, with time in million years ago (MYA).
 482 Macrosynteny was estimated by quantifying the number of times a scaffold of a given species
 483 contained *B. mori* orthologs from two separate chromosomes and were one a single *P. napi*
 484 chromosome (napi-like join), or vice versa (mori-like joins)(see Supplemental Note for more
 485 details). **b)** A time calibrated ancestral state reconstruction of the chromosomal fusion and fission
 486 events across Pieridae (n=201 species). As only a time calibrated genus level phylogeny exists for
 487 Pieridae, all genera with > 1 species are set to an arbitrary polytomy at 5 MYA, while deeper
 488 branches reflect fossil calibrated nodes. The haploid chromosomal count of tips (histogram) and
 489 interior branches (color coding) are indicated, with the outgroup set to n=31 reflecting the butterfly
 490 chromosomal mode. Genus names are indicated for the larger clades (all tips labels in Supplemental
 491 Material).

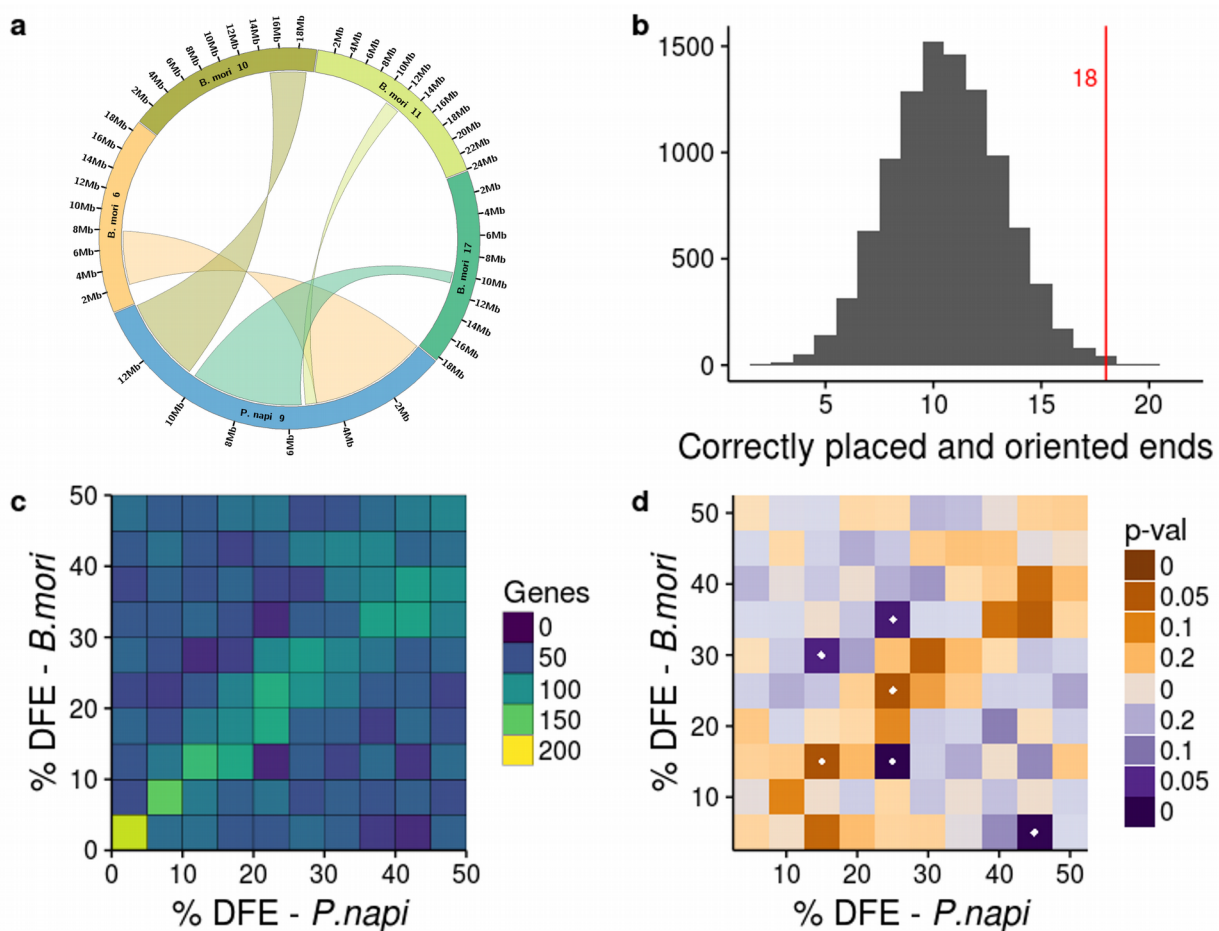


492

493

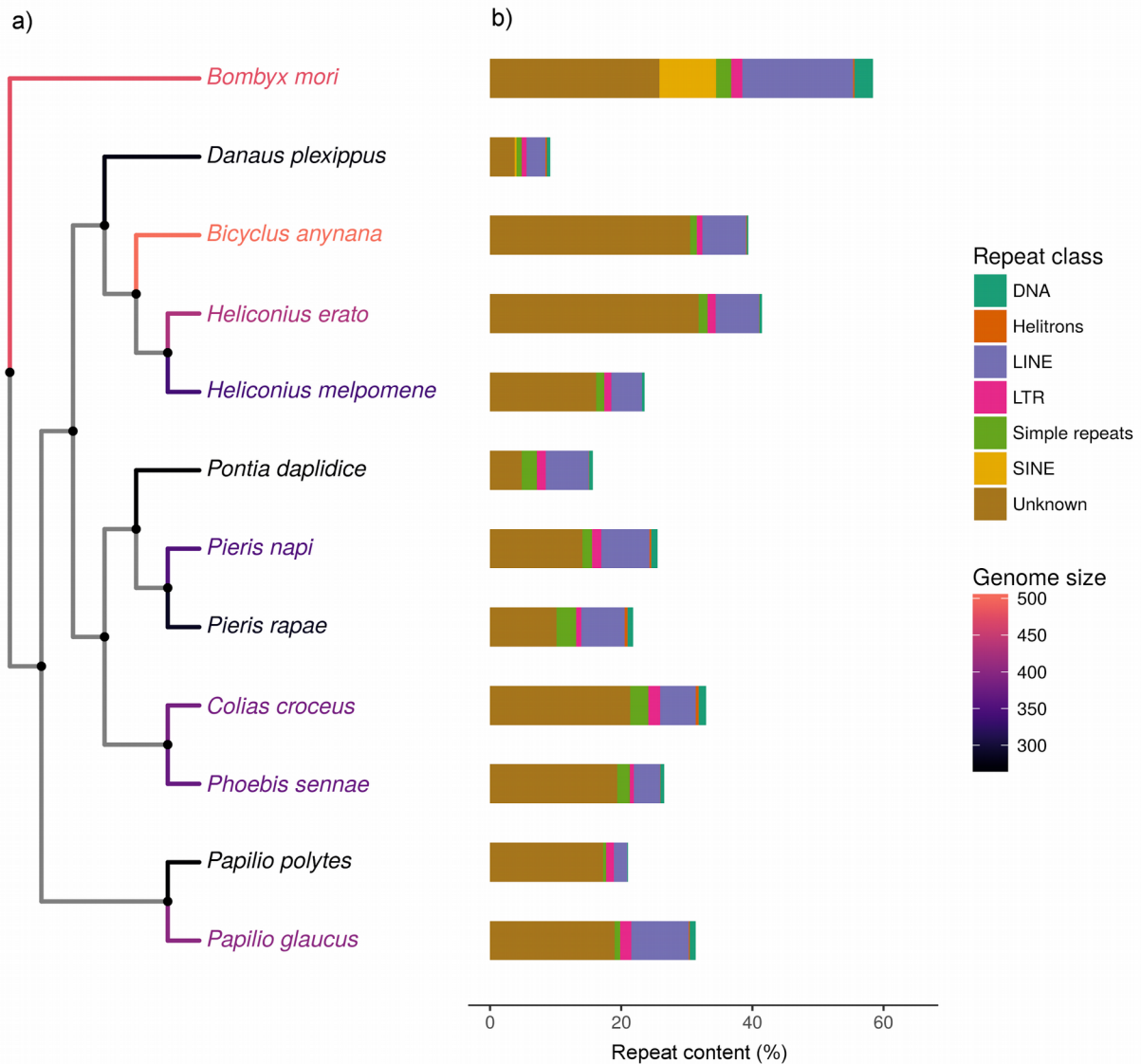
Figure 3 Validation of syntenic relationship between *B. mori* and first four *P. napi* chromosomes. (a) Mate pair spanning depth across each chromosome summed for the 3kb, 7kb, and 40kb libraries. Spanning depths averaged 1356 across the whole genome. Of the scaffold join positions 74 of 97 were spanned by > 50 properly paired reads (mean = 117.8, S.D. = 298.7) which we considered good evidence for correct assembly at scaffold boundaries while the remaining 23 scaffold joins had 0 mate pair spans. (b) RAD-seq linkage markers and recombination distance along chromosomes from the first linkage map that was used for genome assembly. (c) Results from the second linkage map of maternally inherited markers, using RNA-Seq and whole genome sequencing. All markers within a chromosome are completely linked due to suppressed recombination in females (i.e. recombination distance is not shown on Y axis). (d) Syntenic block origin and orientation colored and labeled by the *B. mori* chromosome containing the orthologs, as in Fig. 1 (e) Component scaffolds of each chromosome labeled to indicate scaffold number and orientation. (f) To the right of each *P. napi* chromosome is a circos plot showing the location and orientation of syntenic blocks within each *B. mori* chromosome that comprise a given *P. napi* chromosome. Ribbons representing the blocks of synteny are colored by their orthologs location in the *B. mori* genome. Relative orientation of a block is shown by whether the ribbon contains a twist. Remaining chromosomes shown in Supplementary Fig. 2.

510
511
512
513
514



515
516

Figure 4. Comparison of gene content of and chromosomal location of syntenic blocks between *Pieris napi* and *Bombyx mori* in observed and randomly generated expectation genomes. (a) Observed pattern of conserved syntenic block location within *P. napi* Chromosome 9, wherein telomere facing and interior syntenic blocks are conserved between species despite shuffling. (b) Histogram of the number of syntenic blocks that are terminal on the *B. mori* genome and also occur in the terminal position on chromosomes in a simulated genome, from 10,000 simulated genomes (average 10.7, std dev= 6.8). (c) Percentage distance from the end (DFE) of a chromosome of a single copy gene in *P. napi* vs. DFE of that gene's single copy ortholog (SCO) in *B. mori*. Counts binned on the color axis. (d) Comparison between the observed DFE distribution and the expected distribution generated from 10,000 genomes of 25 chromosomes constructed from the random fusion of syntenic blocks. Bins in which more genes occur in the observed genomes than the expected distribution are in orange, less genes in blue, $P < 0.05$ in either direction are denoted by a white dot. SCO spatial distribution was significantly higher than expected along the diagonal (two bins with $p < 0.05$), while significantly lower than expected off the diagonal (four bins with $p < 0.05$).



532
533 **Figure 5.** The genomic size and repeat content of Lepidopteran genomes placed in a phylogenetic
534 context. (a) Phylogenetic relationships represented as a cladogram, with terminal branches and
535 species names colored by genome size estimates from k-mer distributions of read data. (b) The
536 fraction of repeat content of each genome, color coded by repeat class.
537