1    A butterfly chromonome reveals selection dynamics during extensive and

2    cryptic chromosomal reshuffling

3

4

5    Authors

6

7    Jason Hill[1,‡], Ramprasad Neethiraj[1], Pasi Rastas[2], Nathan Clark[3], Nathan Morehouse[4], Maria de la

8    Paz Celorio-Mancera[1], Jofre Carnicer Cols[5,6], Heinrich Dircksen[10], Camille Meslin[3], Kristin

9    Sikkink[7], Maria Vives[5,6], Heiko Vogel[9], Christer Wiklund[1], Carol L. Boggs[8], Sören Nylin[1],

10   Christopher Wheat[1,‡]

11   Affiliations:

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

13   2 Institute of Biotechnology, (DNA Sequencing and Genomics), University of Helsinki, Helsinki,

14   Finland

15   3 Department of Computational and Systems Biology, University of Pittsburgh, Pittsburgh, USA

16   4 Department of Biological Sciences, University of Cincinnati, Cincinnati, USA

17   5  Department of Evolutionary Biology, Ecology and Environmental Sciences, University of

18   Barcelona, 08028 Barcelona, Spain

19   6  CREAF, Global Ecology Unit, Autonomous University of Barcelona, 08193 Cerdanyola del

20   Vallès, Spain

21   7 Department of Ecology, Evolution and Behavior, University of Minnesota, St Paul, MN, USA

22   8 Department of Biological SciencesUniversity of South Carolina, Columbia, SC, USA

23   9 Department of Entomology, Max Planck Institute for Chemical Ecology, D-07745 Jena, Germany

24   10 Functional Morphology, Department of Zoology, Stockholm University, Stockholm, Sweden

25

26    ᴵ Authors for correspondence:

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

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

29

# Abstract

31        Taxonomic orders vary in their degree of chromosomal conservation with some having high

32    rates of chromosome number turnover despite maintaining some core sets of ordered genes (e.g.

33    Mammalia) and others exhibiting rapid rates of gene-order reshuffling without changing

34    chromosomal count (e.g. Diptera). However few clades exhibit as much conservation as the

35    Lepidoptera for which both chromosomal count and gene colinearity (synteny) are very high over

36    the past 140 MY. In contrast, here we report extensive chromosomal rearrangements in the  genome

37    of the green-veined white butterfly (*Pieris napi,* Pieridae, Linnaeus, 1758). This unprecedented

38    reshuffling is cryptic: microsynteny and chromosome number do not indicate the extensive

39    rearrangement revealed by a chromosome level assembly and high-resolution linkage map.

40    Furthermore, the rearrangement blocks themselves appear to be non-random, as they are

41    significantly enriched for clustered groups of functionally annotated genes revealing that the

42    evolutionary dynamics acting on Lepidopteran genome structure are more complex than previously

43    envisioned.

# Introduction

45        The role of chromosomal rearrangements in adaptation and speciation has long been

46    appreciated and recent work has elevated the profile of supergenes in controlling complex adaptive

47    phenotypes[1–4]. Chromosome number variation has also been cataloged for many species but

48    analyses of the adaptive implications have mostly been confined to the consequences of polyploidy

49    and whole genome duplication[5,6]. The identification of pervasive fission and fusion events

50    throughout the genome is relatively unexplored since discovery of this pattern requires chromosome

51  level assemblies. This leaves open the possibility of cryptic chromosomal dynamics taking place in

52  many species for which this level of genome assembly has not been achieved. As chromosomal

53  levels assemblies become more common, uncovering a relationship between such dynamics and

54  adaptation or speciation can be assessed.

55  Here we focus upon the Lepidoptera, the second most diverse animal group with over 160,000

56  extant species in more than 160 families. Butterflies and moths exist in nearly all habitats and have

57  equally varied life histories yet show striking similarity in genome architecture, with the vast

58  majority having a haploid chromosome number of n=31[7–9]. While haploid chromosome number can

59  vary from n = 5 to n = 223[10–12], gene order and content is remarkably similar within chromosomes

60  (i.e. displays macrosynteny), regardless of haploid chromosome number. The degree of such

61  synteny between species separated by up to 140 My is astounding as illustrated by recent

62  chromosomal level genomic assemblies[7,13], as well as previous studies of the sequence and structure

63  of lepidopteran genomes[14–17]. This ability of Lepidoptera to accommodate such chromosomal

64  rearrangements, yet maintain  high levels of macro and microsynteny (i.e. collinearity at the scale of

65  10s to 100's of genes) is surprising. While a growing body of evidence indicates that gene order in

66  eukaryotes is non-random along chromosomes, with upwards of 12% of genes organized into

67  functional neighborhoods of shared function and expression patterns[18,19], to what extent this may

68  play a role in the chromosomal evolution is an open question.

69  Variation in patterns of synteny across clades must arise due to an evolutionary interaction between

70  selection and constraint[20], likely at the level of telomere and centromere performance. *Drosophila*,

71  and likely all Diptera, differ from other eukaryotes studied to date in lacking the telomerase

72  enzyme, and instead protect their chromosomal ends using retrotransposons[21]. This absence of

73  telomerase is posited to make evolving novel telomeric ends more challenging, limiting the

74  appearance of novel chromosomes and thereby resulting in high macrosynteny via constraint[22].

75  In contrast, Lepidoptera like most Metazoans use telomerase to protect their chromosomal ends

76  which allows for previously internal chromosomal DNA to become subtelomeric in novel

77 chromosomes[7,13]. Additionally all Lepidoptera have holocentric chromosomes in which the

78 decentralized kinetochore allows for more rearrangements by fission, fusion, and translocation of

79 chromosome fragments than is the case for monocentric chromosomes[23]. Thus, Lepidoptera should

80 be able to avoid the deleterious consequences of large-scale chromosomal changes.

81 Here we present the chromosome level genome assembly of the green-veined white butterfly *P.*

82 *napi*. Our analysis reveals large-scale fission and fusion events similar to known dynamics in other

83 lepidopteran species but at an accelerated rate and without a change in haploid chromosome count.

84 The resulting genome-wide breakdown of the chromosome level synteny is unique among

85 Lepidoptera. While we are unable to identify any repeat elements associated with this cryptic

86 reshuffling, we find the chromosomal ends reused and the collinearity of functionally related genes.

87 These findings support a reinterpretation of the chromosomal fission dynamics in the Lepidoptera.

88

## Results

89

90 The *P. napi* genome was generated using DNA from inbred siblings from Sweden, a genome

91 assembly using variable fragment size libraries (180 bp to 100 kb; N50-length of 4.2 Mb and a total

92 length of 350 Mb), and a high density linkage map across 275 full-sib larva, which placed 122

93 scaffolds into 25 linkage groups, consistent with previous karyotyping of *P. napi*[24,25]. After

94 assessment and correction of the assembly, the total chromosome level assembly was 299 Mb

95 comprising 85% of the total assembly size and 114% of the k-mer estimated haploid genome size,

96 with 2943 scaffolds left unplaced (**Supplementary Note 3**). Subsequent annotation predicted

97 13,622 gene models, 9,346 with functional predictions (**Supplementary Note 4**).

98 Single copy orthologs (SCOs) in common between *P. napi* and the first sequenced

99 Lepidopteran genome, the silk moth *Bombyx mori* (Bombycidae), were identified. These revealed

100 an unexpected deviation in gene order and chromosomal structure in *P. napi* relative to *B. mori* as

101 well as another lepidopteran genome with a linkage map and known chromosomal structure, that of

102　*Heliconius melpomene* (Nymphalidae) (Fig 1a). Large-scale rearrangements that appeared to be the

103　fission and subsequent fusion of fragments on the scale of megabases were present on every *P. napi*

104　chromosome relative to *B. mori, H. melpomene,* and *Melitaea cinxia* (Nymphalidae) (fig 1b). We

105　characterized the size and number of large scale rearrangements between *P. napi* and *B. mori* using

106　shared SCOs to identify 99 clearly defined blocks of co-linear gene order (hereafter referred to as

107　"syntenic blocks"), with each syntenic block having an average of 69 SCOs. Each *P. napi*

108　chromosome contained an average of 3.96 (SD = 1.67) syntenic blocks, which derived on average

109　from 3.5 different *B. mori* chromosomes. In *P. napi*, the average syntenic block length was 2.82 Mb

110　(SD = 1.97 Mb) and contained 264 genes (SD = 219).

111　　　The indication that *P. napi* diverged radically from the thus far observed chromosomal

112　structure of lepidopterans raised questions about how frequently a *P. napi* like chromosomal

113　structure is observed *vs.* the structure reported in the highly syntenic *B. mori, H. melpomene,* and

114　*M. cinxia* genomes.  We accessed 22 publicly available lepidopteran genome assemblies and their

115　gene annotations representing species that diverged up to 140 MYA in order to identify the genes

116　corresponding to the SCO's used in the previous analyses. We used blastx (Diamond v0.9.10)[26] to

117　place those genes on their native species scaffolds. With informations about each SCO's location on

118　the *P. napi* and *B. mori* chromosomes we recorded how often a scaffold contained a cluster of genes

119　whose orthologs resided on two *P. napi* chromosomes or two *B. mori* chromosomes. If two *P. napi*

120　chromosomes were represented by only a single *B. mori* chromosome, then the scaffold was marked

121　as containing an mori-like join. Conversely if two *B. mori* chromosomes were represented but only

122　a single *P. napi* chromosome, then the scaffold was marked as containing a napi-like join. In total

123　we found that 20 species have more mori-like joins, and two species of *Pieris* represented by 3

124　assemblies have more napi-like joins (Fig 2a). While this type of assessment is preliminary the

125　indication is that the genome structure described here is novel to the genus *Pieris*.

126　　　We validated this novel chromosomal reorganization using four complementary but

127　independent approaches to assess our scaffold joins. First, we generated a second linkage map for *P.*

128 *napi,* which confirmed the 25 linkage groups and the ordering of scaffold joins along chromosomes

129 (Fig. 3; Supplementary Fig. 2). Second, the depth of the mate-pair (MP) reads spanning joins

130 indicated by the first linkage map provides an independent assessment of the join validity. We

131 therefore quantified MP reads spanning each base pair position along a chromosome (Fig. 3;

132 Supplementary Fig. 2, Note 7), finding strong support for the scaffold joins. Third, we aligned the

133 scaffolds of a recently constructed genome of *P. rapae*[27] to *P. napi,* looking for *P. rapae* scaffolds

134 that spanned the chromosomal level scaffold joins within *P. napi,* finding support for 71 of the 97

135 joins (Supplementary Fig. 5). Fourth, we considered *B. mori* syntenic blocks that spanned a scaffold

136 join within a *P. napi* chromosome as support for that *P. napi* chromosome assembly, and found that

137 62 of the 97 scaffold joins were supported by *B. mori* (Supplementary Fig. 2, Note 8,9).

138       To assess the novel chromosomal organization, we investigated the ordering and content of

139 these syntenic blocks in *P. napi.* First, we tested whether telomeric ends of chromosomes were at all

140 conserved between species despite the extensive chromosomal reshuffling (Fig. 4a). We found

141 significantly more syntenic blocks sharing telomere facing orientations between species than

142 expected ($P < 0.01$, two tailed t-test; Fig. 4b). We also identified a significant enrichment for SCOs

143 in *B. mori* and *P. napi* located at roughly similar distance from the end of their respective

144 chromosomes (Fig. 4c). Both of these findings are consistent with the ongoing use of telomeric

145 ends, indicating that strong selection dynamics have favored their retention over evolutionary time.

146 Second, we tested for gene set functional enrichment within the observed syntenic blocks by

147 investigating the full set of annotated *P. napi* genes. We found that 57 of the 99 block regions in the

148 *P. napi* genome contained at least three genes with a shared gene ontology (GO) term that was

149 significantly less frequent in the rest of the genome ($P < 0.01$, fisher) (Supplementary fig. 3). We

150 then tested whether the observed enrichment in the syntenic blocks of *P. napi* was greater than

151 expected by randomly assigning the genome into similarly sized blocks. The mean number of GO

152 enriched fragments in each of the simulated 10,000 genomes was 38.8 (variance of 46.6 and

153 maximum of 52), which was significantly lower than the observed ($P < 0.0001$).

154    To assess the possible cause of the reshuffling, we surveyed the distribution of different

155    repeat element classes across the genome, looking for enrichment of specific categories near the

156    borders of syntenic blocks. While Class 1 transposons were found to be at higher density at near the

157    ends of chromosomes relative to the distribution internally (Supplementary fig. 4), no repeat

158    elements were enriched relative to the position of syntenic block regions. We therefore investigated

159    whether any repeat element classes had expanded within *Pieris* compared to other sequenced

160    genomes by assessing the distribution of repeat element classes and genome size among sequenced

161    Lepidoptera genomes. In accordance with other taxa[28] we find an expected strong relationship

162    between genome size and repetitive element content in *Pieris* species. Thus, while repetitive

163    elements such as transposable elements are likely to have been involved in the reshuffling, our

164    inability to find clear elements involved suggests these events may be old and their signal decayed.

# Methods

165

166    **Sample collection and DNA extraction.** Pupal DNA was isolated from a 4th generation inbred

167    cohort that originated from a wild caught female collected in Skåne, Sweden, using a standard salt

168    extraction[29].

169    **Illumina genome sequencing.** Illumina sequencing was used for all data generation used in

170    genome construction. A 180 paired end (PE) and the two mate pair (MP) libraries were constructed

171    at Science for Life Laboratory, the National Genomics Infrastructure, Sweden (SciLifeLab), using 1

172    PCR-free PE DNA library (180bp) and 2 Nextera MP libraries (3kb and 7kb) all from a single

173    individual. All sequencing was done on Illumina HiSeq 2500 High Output mode, PE 2x100bp by

174    SciLifeLab. An additional two 40kb MP fosmid jumping libraries were constructed from a sibling

175    used in the previous library construction. Genomic DNA, isolated as above, was shipped to Lucigen

176    Co. (Middleton, WI, USA) for the fosmid jumping library construction and sequencing was

177    performed on an Illumina MiSeq using 2x250bp reads [30]. Finally, a variable insert size library of

178    100 bp – 100,000 bp in length were generated using the Chicago and HiRise method[31]. Genomic

179    DNA was again isolated from a sibling of those used in previous library construction. The genomic

180    DNA was isolated as above and shipped to Dovetail Co. (Santa Cruz, CA, USA) for library

181    construction, sequencing and scaffolding. These library fragments were  sequenced by Centrillion

182    Biosciences Inc. (Palo Alto, CA, USA) using Illumina HiSeq 2500 High Output mode, PE 2x100bp.

183    **Data Preparation and Genome assembly.** Nearly 500 M read pairs of data were generated,

184    providing ~ 285 X genomic coverage (Supplemental Table 1). The 3kb and 7kb MP pair libraries

185    were filtered for high confidence true mate pairs using Nextclip v0.8[32]. All read sets were then

186    quality filtered, the ends trimmed of adapters and low quality bases, and screened of common

187    contaminants using bbduk v37.51 (bbtools, Brian Bushnell). Insert size distributions were plotted to

188    assess library quality, which was high (Supplementary Fig. 1). The 180bp, 3kb, and 7kb, read data

189    sets were used with AllpathsLG r50960[33] for initial contig generation and scaffolding

190    (Supplementary Note 1). AllpathsLG was run with haploidify = true to compensate for the high

191    degree of heterozygosity. The initial contig assembly's conserved single copy ortholog content was

192    assessed at 78% for *P. napi*  by CEGMA v2.5[34]. A further round of superscaffolding using the 40kb

193    libraries alongside the 3kb and 7kb libraries was done using SSPACE v2[35]. Finally, both assemblies

194    were Ultascaffolded using the Chicago read libraries and the HiRise software pipeline. These steps

195    produced a final assembly of 3005 scaffolds with an N50-length of 4.2 Mb and a total length of 350

196    Mb (Supplementary Note 1).

197    **Linkage Map**. RAD-seq data of 5463 SNP markers from 275 full-sib individuals, without parents,

198    was used as input into Lep-MAP2[36]. The RAD-seq data was generated from next-RAD technology

199    by SNPsaurus (Oregon, USA)(Supplemental note 10). To obtain genotype data, the RAD-seq data

200    was mapped to the reference genome using BWA mem[37] and SAMtools[38] was used to produce

201    sorted bam files of the read mappings. Based on read coverage (samtools depth), Z chromosomal

202    regions were identified from the genome and the sex of offspring was determined. Custom scripts[39]

203    were used to produce genotype likelihoods (called posteriors in Lep-MAP) from the output of

204    SAMtools mpileup.

205    The parental genotypes were inferred with Lep-MAP2 ParentCall module using parameters

206    "ZLimit=2 and ignoreParentOrder=1", first calling Z markers and second calling the parental

207    genotypes by ignoring which way the parents are informative (the parents were not genotyped so

208    we could not separate maternal and paternal markers at this stage). Scripts provided with Lep-

209    MAP2 were used to produce linkage file from the output of ParentCall and all single parent

210    informative markers were converted to paternally informative markers by swapping parents, when

211    necessary. Filtering by segregation distortion was performed using Filtering module.

212    Following this, the SepareteChromosomes module was run on the linkage file and 25 chromosomes

213    were identified using LOD score limit 39. Then JoinSingles module was run twice to add more

214    markers on the chromosomes with LOD score limit of 20. Then SepareteChromosomes was run

215    again but only on markers informative on single parent with LOD limit 10 to separate paternally

216    and maternally informative markers. 51 linkage groups were found and all were ordered using

217    OrderMarkers module. Based on likelihood improvement of marker ordering, paternal and maternal

218    linkage groups were determined. This was possible as there is no recombination in females

219    (achiasmatic meiosis), and thus the order of the markers does not improve likelihood on the female

220    map. The markers on the corresponding maternal linkage groups were converted to maternally

221    informative and OrderMarkers was run on the resulting data twice for each of 25 chromosomes

222    (without allowing recombination in female). The final marker order was obtained as the order with

223    the higher likelihood from the two runs.

224    **Chromosomal assembly**. The 5463 markers that composed the linkage map were mapped to the *P.*

225    *napi* ultrascaffolds using bbmap[40] with sensitivity = slow. Reads that mapped uniquely were used to

226    identify misassemblies in the Ultrascaffolds and arrange those fragments into chromosomal order.

227    54 misassemblies were identified and overall 115 fragments were joined together into 25

228    chromosomes using a series of custom R scripts (supplemental information) and the R package

229    Biostrings[41]. Scaffold joins and misassembly corrections were validated by comparing the number

230    of correctly mapped mate pairs spanning a join between two scaffolds. Mate pair reads from the

231     3kb, 7kb, and 40kb libraries were mapped to their respective assemblies with bbmap (po=t,

232     ambig=toss, kbp=t). SAM output was filtered for quality and a custom script was used to tabulate

233     read spanning counts for each base pair in the assembly.

234     **Synteny Comparisons Between *P. napi, B. mori,* and *H. melpomene*.** A list of 3100 single copy

235     orthologs (SCO) occurring in the Lepidoptera lineage curated by OrthoDB v9.1[42] was used to

236     extract gene names and protein sequences of SCOs in *Bombyx mori* from

237     KaikoBase[43] (Supplemental Note 5) using a custom script. Reciprocal best hits (RBH) between gene

238     sets of *P. napi, P. rapae, H. melpomene, M. cinxia,* and *B. mori* SCOs were identified using

239     BLASTP[44] and custom scripts. Gene sets of *H. melpomene* v2.5  and *M. cinxia* v1 were downloaded

240     from LepBase v4 [45]. Coordinates were converted to chromosomal locations and visualized using

241     Circos[46] and custom R scripts.

242     **Synteny Comparison Within Lepidoptera.** Genome assemblies and annotated protein sets were

243     downloaded for 24 species of Lepidoptera from LepBase v4 [47] and other sources (Supplemental

244     Table 4). Each target species protein set was aligned to its species genome as well as to the *Pieris*

245     *napi* protein set using Diamond v0.9.10[26] with default options. The protein-genome comparison was

246     used to assign each target species gene to one of it's assembled scaffolds, while the  protein-protein

247     comparison was used to identify RBHs  between the protein of each species and its ortholog in *P.*

248     *napi,* and *B. mori* . Using this information we used a custom R script to examine each assembly

249     scaffold for evidence of synteny to either *P. napi* or *B. mori*. First, each scaffold of the target species

250     genome was assigned genes based on the protein-genome blast results, using its own protein set and

251     genome. A gene was assigned to a scaffold if at least 3 HSPs of less than 200bp from a gene aligned

252     with >= 95% identity. Second, if any of these scaffolds then contained genes whose orthologs

253     resided on a single *B. mori* chromosome but two *P. napi* chromosomes, and those same two *P. napi*

254     chromosome segments were also joined in the *B. mori* assembly, that was counted as a 'mori-like

255     join'. Conversely if a target species scaffold contained genes whose orthologs resided on a single *P.*

256  *napi* chromosome but two *B. mori* chromosomes, and those same two *B. mori* chromosome

257  segments were also joined in the *P. napi* assembly, that was counted as a 'napi-like join'.

258  **Pieridae chromosomal evolution.**

259  Chromosomal fusions and fissions were reconstructed across the family Pieridae by placing

260  previously published karyotype studies of haploid chromosomal counts into their evolutionary

261  context. There are approximately 1000 species in the 85 recognized genera of Pieridae and we

262  recently reconstructed a robust fossil-calibrated chronogram for this family at the genus level[48,49].

263  We then placed the published chromosomal counts for 201 species[9,50] on this time calibrated

264  phylogeny with ancestral chromosomal reconstructions for chromosome count, treated as a

265  continuous character, using the contMap function of the phytools R package[51].

266  **Second Linkage Map for *P. napi*.** A second linkage map was constructed from a different family

267  of *P. napi* in which a female from Abisko, Sweden was crossed with a male from Catalonia, Spain.

268  Genomic DNA libraries were constructed for the mother, father, and four offspring (2 males, 2

269  females). RNA libraries were constructed for an additional 6 female and 6 male offspring. All

270  sequencing was performed on a Illumina HiSeq 2500 platform using High Output mode, with PE

271  2x100bp reads at SciLifeLab (Stockholm, Sweden). Both DNA and RNA reads were mapped to the

272  genome assembly with bbmap. Samtools was used to sort read mappings and merge them into an

273  mpileup file (Supplemental Note 6). Variants were called with BCFtools[52] and filtered with

274  VCFtools[53]. Linkage between SNPs was assessed with PLINK[54]. A custom script was used to assess

275  marker density and determine sex-specific heterozygosity.

276  **Annotation of *P. napi* genome**. Genome annotation was carried out by the Bioinformatics Short-

277  term Support and Infrastructure (BILS, Sweden). BILS was provided with the chromosomal

278  assembly of *P. napi* and 45 RNAseq read sets representing 3 different tissues (head, fat body, and

279  gut) of 7 male and 8 female larva from lab lines were separate from the one used for the initial

280  sequencing. Sequence evidence for the annotation was collected in two complementary ways. First,

281    we queried the Uniprot database[55] for protein sequences belonging to the taxonomic group of

282    Papilionoidea (2,516 proteins). In order to be included, proteins gathered in this way had to be

283    supported on the level of either proteomics or transcriptomics and could not be fragments. In

284    addition, we downloaded the Uniprot-Swissprot reference data set (downloaded on 2014-05-15)

285    (545,388 proteins) for a wider taxonomic coverage with high-confidence proteins. In addition,  493

286    proteins were used that derived from a *P. rapae* expressed sequence tag library that was Sanger

287    sequenced.

288    **Permutation test of syntenic block position within chromosomes**. Syntenic blocks (SBs) were

289    identified as interior vs terminal and the ends of terminal blocks were marked as inward or outward

290    facing. SBs were reshuffled into 25 random chromosomes of 4 SBs in a random orientation and the

291    number of times that a terminal block occurred in a random chromosome with the outward end

292    facing outward was counted. This was repeated 10,000 times to generate a random distribution

293    expectation. The number of terminal outward-facing SBs in *B. mori* that were also terminal and

294    outward facing in *P. napi* was compared to this random distribution to derive the significance of

295    deviation from the expected value. To test the randomness of gene location within chromosomes,

296    orthologs were numbered by their position along each chromosome in both *B. mori* and *P. napi*.

297    10,000 random genomes were generated as above. Distance from the end of the new chromosome

298    and distance from the end of *B. mori* chromosome were calculated for each ortholog and the results

299    were binned. P-values were determined by comparing the number of orthologs in a bin to the

300    expected distribution of genes in a bin from the random genomes. All test were done using a custom

301    R script.

302    **Gene set enrichment analysis of syntenic blocks.** Gene ontology set enrichment was initially

303    tested within syntenic blocks of the *P. napi* genome using topGO[56] with all 13,622 gene models

304    generated from the annotation. For each syntenic block within the genome, each GO term of any

305    level within the hierarchy that had at least 3 genes belonging to it was analyzed for enrichment. If a

306    GO term was overrepresented in a syntenic block compared to the rest of the genome at a p-value of

307    < 0.01 by a Fisher exact test, that block was counted as enriched. 57 of the 99 syntenic blocks in the

308    *P. napi* genome were enriched in this way. Because arbitrarily breaking up a genome and testing for

309    GO enrichment can yield results that are dependent on the distribution of the sizes used, we

310    compared the results of the previous analysis to the enrichment found using the same size genomic

311    regions, randomly selected from the *P. napi* genomes. The size distribution of the 99 syntenic

312    blocks were used to generate fragment sizes into which the genome was randomly assigned. This

313    resulted in a random genome of 99 fragments which in total contained the entire genome but the

314    content of a given fragment was random compared to the syntenic block that defined its size. This

315    random genome was tested for GO enrichment of the fragments in the same way as the syntenic

316    blocks in the original genome, and the number of enriched blocks counted. This was then repeated

317    10,000 times to generate a distribution of expected enrichment in genome fragments of the same

318    size as the *P. napi* syntenic blocks.

319

## Discussion

321        While massive chromosomal fission events are well documented in butterflies  (e.g.

322    *Leptidea* in Pieridae (n=28-103); *Agrodiaetus* in Lycaenidae (n=10-134)), their contribution to

323    Lepidopteran diversity appears to be minimal as all these clades  are very young[57–59]. However, our

324    results challenge this interpretation. Rather,  *P. napi* appears to represent a lineage that has

325    undergone an impressive reconciliation of an earlier series of rampant fission events. Moreover, the

326    subsequent fusion events exhibit a clear bias toward using ancient telomeric ends, as well as

327    returning gene clusters to their relative ancestral position within chromosomes even when the other

328    parts of the newly formed chromosome originated from other sources.  Luckily these initial fission

329    events have been frozen in time as reshuffled syntenic blocks, revealing the potential fitness

330    advantage of maintaining certain functional categories as syntenic blocks.

331    Thus, despite the potential for holocentric species to have relaxed constraint upon their

332    chromosomal evolution, we find evidence for selection actively maintaining ancient telomeric ends,

333    as well as gene order within large chromosomal segments. Together these observations suggest that

334    the low chromosome divergence in Lepidoptera over > 100 million generations is at least partially

335    due to purifying selection maintaining an adaptive chromosomal structure.

336

## Bibliography

338

339    1.    Schwander, T., Libbrecht, R. & Keller, L. Supergenes and complex phenotypes. *Curr. Biol.*
340          **24,** R288–R294 (2014).

341    2.    Kunte, K. *et al.* Doublesex Is a Mimicry Supergene. *Nature* **507,** 229–232 (2014).

342    3.    Fishman, L., Stathos, A., Beardsley, P. M., Williams, C. F. & Hill, J. P. Chromosomal
343          rearrangements and the genetics of reproductive barriers in mimulus (monkey flowers).
344          *Evolution (N. Y).* **67,** 2547–2560 (2013).

345    4.    Lamichhaney, S. *et al.* Structural genomic changes underlie alternative reproductive
346          strategies in the ruff (Philomachus pugnax). *Nat. Genet.* **48,** 84–88 (2015).

347    5.    Otto, S. P. & Whitton, J. Polyploid Incidence and Evolution. *Annu. Rev. Genet.* **34,** 401–437
348          (2000).

349    6.    Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of polyploidy.
350          *Nat. Rev. Genet.* **18,** 411–424 (2017).

351    7.    Ahola, V. *et al.* The Glanville fritillary genome retains an ancient karyotype and reveals
352          selective chromosomal fusions in Lepidoptera. *Nat Commun* **5,** 1–9 (2014).

353    8.    Lukhtanov, V. A. Sex chromatin and sex chromosome systems in nonditrysian Lepidoptera
354          (Insecta). *J. Zool. Syst. Evol. Res.* **38,** 73–79 (2000).

355    9.    ROBINSON, R. *Lepidoptera Genetics*. *Lepidoptera Genetics* (1971). doi:10.1016/B978-0-
356          08-006659-2.50017-1

357    10.   Brown, K. S., Von Schoultz, B. & Suomalainen, E. Chromosome evolution in Neotropical
358          Danainae and Ithomiinae (Lepidoptera). *Hereditas* **141,** 216–236 (2004).

359    11.   Kandul, N. P., Lukhtanov, V. A. & Pierce, N. E. Karyotypic diversity and speciation in
360          Agrodiaetus butterflies. *Evolution (N. Y).* **61,** 546–559 (2007).

361    12.   Saura, A., Schoultz, B. Von, Saura, A. O. & Brown, K. S. Chromosome evolution in
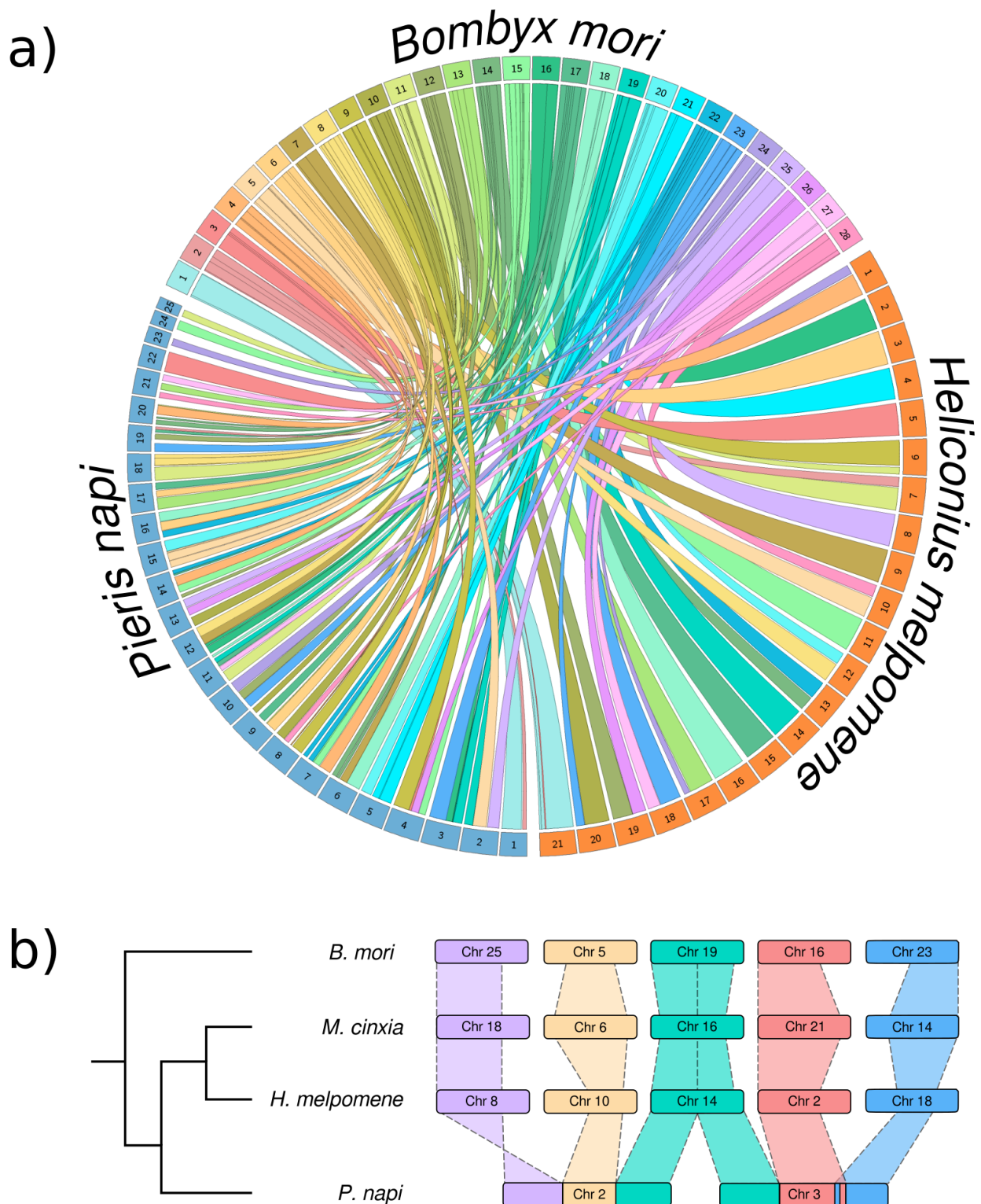362          Neotropical butterflies. *Hereditas* **150,** 26–37 (2013).

363  13.  Davey, J. W. *et al.* Major Improvements to the Heliconius melpomene Genome Assembly
364      Used to Confirm 10 Chromosome Fusion Events in 6 Million Years of Butterfly Evolution.
365      *Genes|Genomes|Genetics* **6,** 695–708 (2016).

366  14.  Yasukochi, Y. A Second-Generation Integrated Map of the Silkworm Reveals Synteny and
367      Conserved Gene Order Between Lepidopteran Insects. *Genetics* **173,** 1319–1328 (2006).

368  15.  Yasukochi, Y. *et al.* A FISH-based chromosome map for the European corn borer yields
369      insights into ancient chromosomal fusions in the silkworm. *Heredity (Edinb).* **116,** 75–83
370      (2016).

371  16.  Beldade, P., Saenko, S. V, Pul, N. & Long, A. D. A gene-based linkage map for Bicyclus
372      anynana butterflies allows for a comprehensive analysis of synteny with the lepidopteran
373      reference genome. *PLoS Genet.* **5,** e1000366 (2009).

374  17.  Šíchová, J. *et al.* Fissions, fusions, and translocations shaped the karyotype and multiple sex
375      chromosome constitution of the northeast-Asian wood white butterfly, Leptidea amurensis.
376      *Biol. J. Linn. Soc.* **118,** 457–471 (2016).

377  18.  Al-Shahrour, F. *et al.* Selection upon genome architecture: Conservation of functional
378      neighborhoods with changing genes. *PLoS Comput. Biol.* **6,** (2010).

379  19.  Gordon, J. L., Byrne, K. P. & Wolfe, K. H. Mechanisms of chromosome number evolution in
380      yeast. *PLoS Genet.* **7,** 0–3 (2011).

381  20.  Coghlan, A., Eichler, E. E., Oliver, S. G., Paterson, A. H. & Stein, L. Chromosome evolution
382      in eukaryotes: A multi-kingdom perspective. *Trends Genet.* **21,** 673–682 (2005).

383  21.  Levis, R. W., Ganesan, R., Houtchens, K., Tolar, L. A. & Sheen, F. miin. Transposons in
384      place of telomeric repeats at a Drosophila telomere. *Cell* **75,** 1083–1093 (1993).

385  22.  Muller, H. J. The remaking of chromosomes. *Collect. net* **13,** 181–198 (1938).

386  23.  Maddox, P. S., Oegema, K., Desai, A. & Cheeseman, I. M. 'Holo'er than thou: Chromosome
387      segregation and kinetochore function in C. elegans. *Chromosom. Res.* **12,** 641–653 (2004).

388  24.  Lorkovic, Z. The genetics and reproductive isolating mechanisms of the Pieris napi -
389      bryoniae group. *J. Lepid. Soc.* **16,** 5–19 (1955).

390  25.  Maeki, K. & Kawazoe, A. On the Hybridization between Two Karyotype Lineages of Pieris
391      napi Linnaeus from Japan. *Cytologia (Tokyo).* **59,** (1994).

392  26.  Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
393      DIAMOND. *Nat. Methods* **12,** 59–60 (2014).

394  27.  Nallu, S. *et al.* The Molecular Genetic Basis of Herbivory between Butterflies and their Host-
395      Plants. *bioRxiv* (2017). doi:10.1101/154799

396  28.  Clark, A. G. *et al.* Evolution of genes and genomes on the Drosophila phylogeny. *Nature* **450,**
397      203–218 (2007).

398  29.  Aljanabi, S. M. & Martinez, I. Universal and rapid salt-extraction of hight quality genomic
399      DNA for PCR-based techniques. *Nucleic Acids Res.* **25,** 4692–4693 (1997).

400   30.   Bentley, D. R. *et al.* Accurate whole human genome sequencing using reversible terminator
401         chemistry. *Nature* **456,** 53–59 (2008).

402   31.   Putnam, N. H. *et al.* Chromosome-scale shotgun assembly using an in vitro method for long-
403         range linkage. *Genome Res.* **26,** 342–350 (2016).

404   32.   Leggett, R. M., Clavijo, B. J., Clissold, L., Clark, M. D. & Caccamo, M. NextClip: an
405         analysis and read preparation tool for Nextera Long Mate Pair libraries. *Bioinformatics* **30,**
406         566–8 (2014).

407   33.   Gnerre, S. *et al.* High-quality draft assemblies of mammalian genomes from massively
408         parallel sequence data. *Proc. Natl. Acad. Sci.* **108,** 1513–1518 (2011).

409   34.   Parra, G., Bradnam, K. & Korf, I. CEGMA: a pipeline to accurately annotate core genes in
410         eukaryotic genomes. *Bioinformatics* **23,** 1061–1067 (2007).

411   35.   Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-
412         assembled contigs using SSPACE. *Bioinformatics* **27,** 578–579 (2011).

413   36.   Rastas, P., Calboli, F. C. F., Guo, B., Shikano, T. & Merilä, J. Construction of Ultradense
414         Linkage Maps with Lep-MAP2: Stickleback F2 Recombinant Crosses as an Example.
415         *Genome Biol. Evol.* **8,** 78–93 (2015).

416   37.   Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. **0,**
417         1–3 (2013).

418   38.   Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25,** 2078–
419         2079 (2009).

420   39.   Kvist, J. *et al.* Flight-induced changes in gene expression in the Glanville fritillary butterfly.
421         *Mol. Ecol.* **24,** 4886–4900 (2015).

422   40.   Bushnell, B. BBTools. (2017). at <https://jgi.doe.gov/data-and-tools/bbtools/>

423   41.   Pages H, Gentleman R, Aboyoun P,  et al. Biostrings: String objects representing biological
424         sequences, and matching algorithms. *R Packag. version* **2,** 2008 (2008).

425   42.   Zdobnov, E. M. *et al.* OrthoDB v9.1: Cataloging evolutionary and functional annotations for
426         animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Res.* **45,** D744–
427         D749 (2017).

428   43.   Shimomura, M. *et al.* KAIKObase: an integrated silkworm genome database and data mining
429         tool. *BMC Genomics* **10,** 486 (2009).

430   44.   Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST:a new generation of protein database
431         search programs. *Nucleic Acids Res.* **25,** 3389–3402 (1997).

432   45.   Kersey, P. J. *et al.* Ensembl Genomes 2016: More genomes, more complexity. *Nucleic Acids*
433         *Res.* **44,** D574–D580 (2016).

434   46.   Krzywinski, M. *et al.* Circos: An information aesthetic for comparative genomics. *Genome*
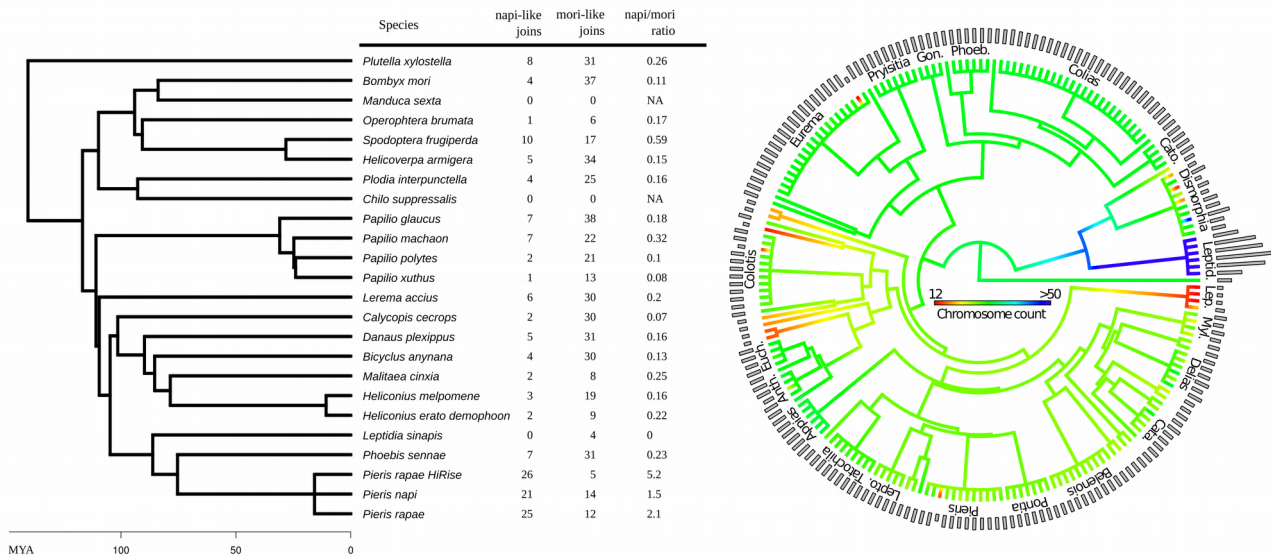435         *Res.* **19,** 1639–1645 (2009).

436  47.  Challis, R. J., Kumar, S., Dasmahapatra, K. K., Jiggins, C. D. & Blaxter, M. Lepbase: The
437       Lepidopteran genome database. *bioRxiv* 56994 (2016). doi:10.1101/056994

438  48.  Wahlberg, N., Rota, J., Braby, M. F., Pierce, N. E. & Wheat, C. W. Revised systematics and
439       higher classification of pierid butterflies (Lepidoptera: Pieridae) based on molecular data.
440       *Zool. Scr.* **43,** 641–650 (2014).

441  49.  Edger, P. P. *et al.* The butterfly plant arms-race escalated by gene and genome duplications.
442       *Proc. Natl. Acad. Sci.* **112,** 8362–8366 (2015).

443  50.  Lukhtanov, V. A. Karyotype evolution and systematics of higher taxa of Pieridae
444       (Lepidoptera) of the World. *Ent. Obozr. 70 619-636, 3 figs* (1991).

445  51.  Revell, L. J. phytools: An R package for phylogenetic comparative biology (and other
446       things). *Methods Ecol. Evol.* **3,** 217–223 (2012).

447  52.  Narasimhan, V. *et al.* BCFtools/RoH: A hidden Markov model approach for detecting
448       autozygosity from next-generation sequencing data. *Bioinformatics* **32,** 1749–1751 (2016).

449  53.  Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27,** 2156–2158
450       (2011).

451  54.  Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based
452       Linkage Analyses. *Am. J. Hum. Genet.* **81,** 559–575 (2007).

453  55.  The UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.* **43,**
454       D204-12 (2015).

455  56.  Alexa A and Rahnenfuhrer J. topGO: Enrichment Analysis for Gene Ontology. *R package
456       version 2.26.0.* (2016). at <http://bioconductor.org/packages/release/bioc/html/topGO.html>

457  57.  Dincǎ, V., Lukhtanov, V. A., Talavera, G. & Vila, R. Unexpected layers of cryptic diversity in
458       wood white Leptidea butterflies. *Nat. Commun.* **2,** (2011).

459  58.  Vila, R. *et al.* Phylogeny and palaeoecology of Polyommatus blue butterflies show Beringia
460       was a climate-regulated gateway to the New World. *Proc. R. Soc. B Biol. Sci.* **278,** 2737–
461       2744 (2011).

462  59.   Lukhtanov, V. The blue butterfly Polyommatus (Plebicula) atlanticus (Lepidoptera,
463  Lycaenidae) holds the record of the highest number of chromosomes in the non-polyploid
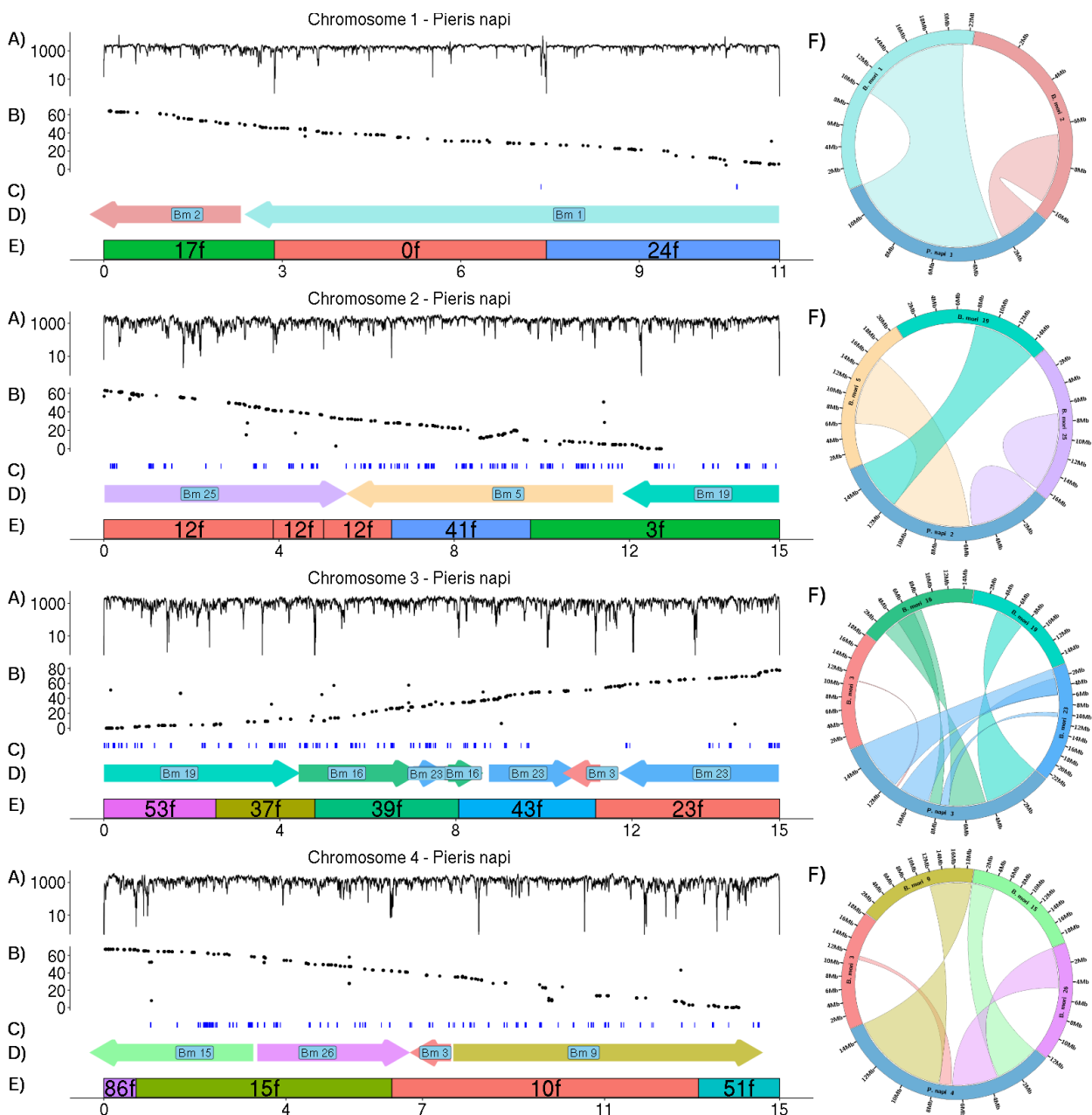464  eukaryotic organisms. *Comp. Cytogenet.* **9,** 683–690 (2015).

465



**Figure 1 a)** Chromosomal mapping between the moth *Bombyx mori* (Bombycoidea) and the butterflies *Pieris napi* (Pieridae) and *Heliconius melpomene* (Nymphalidae). These species last shared a common ancestor > 100 million generations ago[49]. Depicted are the reciprocal best hit orthologs identified between *B. mori* and *P. napi* (n=2354) and between *B. mori* and *H. melpomene* (n=2771). Chromosome 1 is the Z chromosome in *B. mori* and *P. napi* and 21 is the Z chromosome in *H. melpomene*. Chromosomes 2-25 in *P. napi* are ordered in size from largest to smallest. Links between orthologs originate from the *B. mori* chromosome and are colored by their chromosome of origin, while *P. napi* chromosomes are colored blue and *H. melpomene* chromosomes are colored

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

481



| Species | napi-like joins | mori-like joins | napi/mori ratio |
|---|---|---|---|
| Plutella xylostella | 8 | 31 | 0.26 |
| Bombyx mori | 4 | 37 | 0.11 |
| Manduca sexta | 0 | 0 | NA |
| Operophtera brumata | 1 | 6 | 0.17 |
| Spodoptera frugiperda | 10 | 17 | 0.59 |
| Helicoverpa armigera | 5 | 34 | 0.15 |
| Plodia interpunctella | 4 | 25 | 0.16 |
| Chilo suppressalis | 0 | 0 | NA |
| Papilio glaucus | 7 | 38 | 0.18 |
| Papilio machaon | 7 | 22 | 0.32 |
| Papilio polytes | 2 | 21 | 0.1 |
| Papilio xuthus | 1 | 13 | 0.08 |
| Lerema accius | 6 | 30 | 0.2 |
| Calycopis cecrops | 2 | 30 | 0.07 |
| Danaus plexippus | 5 | 31 | 0.16 |
| Bicyclus anynana | 4 | 30 | 0.13 |
| Malitaea cinxia | 2 | 8 | 0.25 |
| Heliconius melpomene | 3 | 19 | 0.16 |
| Heliconius erato demophoon | 2 | 9 | 0.22 |
| Leptidia sinapis | 0 | 4 | 0 |
| Phoebis sennae | 7 | 31 | 0.23 |
| Pieris rapae HiRise | 26 | 5 | 5.2 |
| Pieris napi | 21 | 14 | 1.5 |
| Pieris rapae | 25 | 12 | 2.1 |

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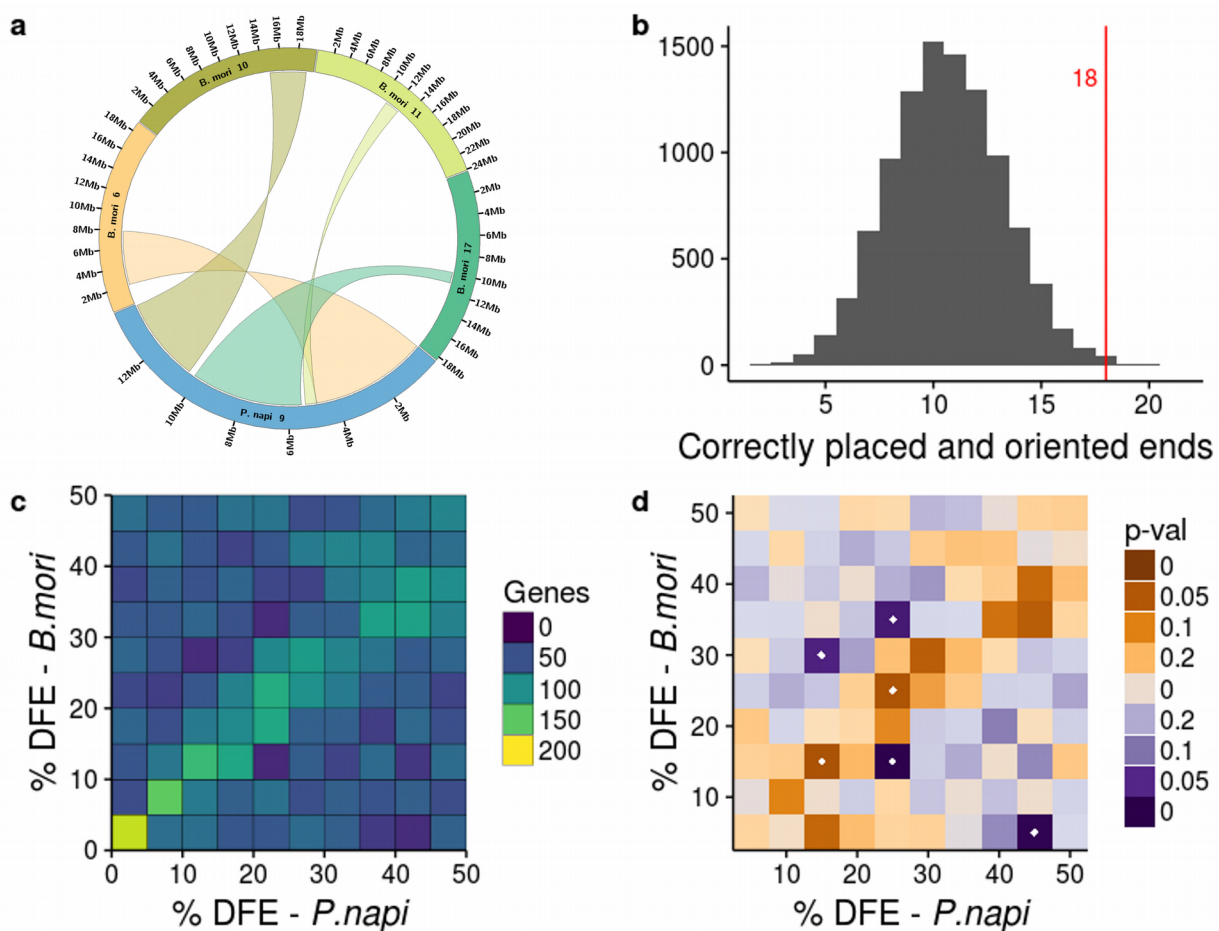493 Material).

**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.
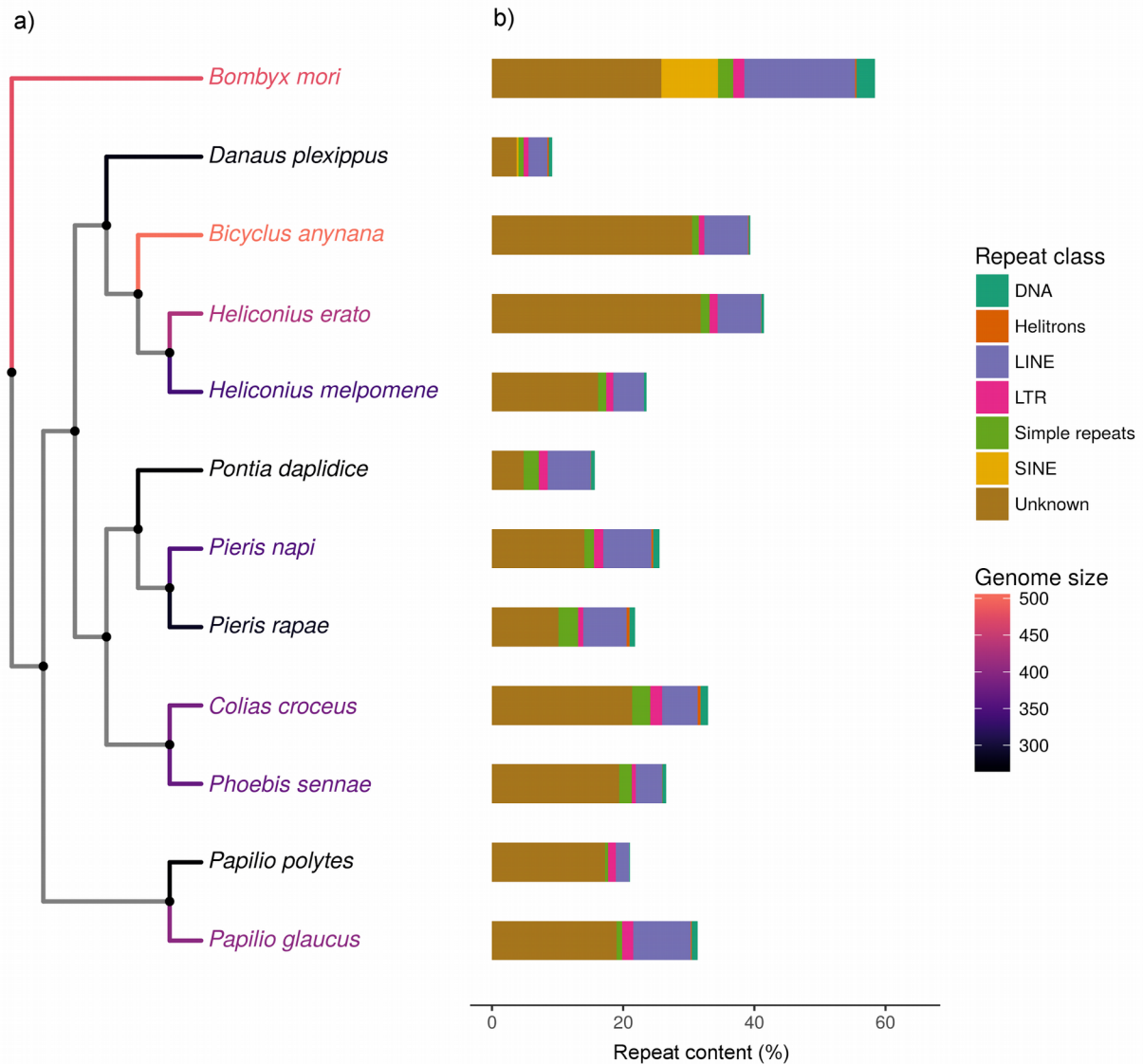
512
513
514
515
516

**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 $p < 0.05$).

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