

1 **Mimicry diversification in *Papilio dardanus* via a genomic inversion in the regulatory**
2 **region of *engrailed-inverted***

3

4

5 Martijn J.T.N. Timmermans^{1,2,3*}, Amrita Srivathsan⁴, Steve Collins⁵, Rudolf Meier^{4,6},
6 Alfried P. Vogler^{1,2}

7

8 1) Department of Life Sciences, Natural History Museum, London, UK

9 2) Department of Life Sciences, Imperial College London, Silwood Park Campus, Ascot,
10 UK

11 3) Department of Natural Sciences, Middlesex University, London, UK

12 4) Department of Biological Sciences, National University of Singapore, Singapore.

13 5) African Butterfly Research Institute, Nairobi, Kenya

14 6) Lee Kong Chian Natural History Museum, National University of Singapore, Singapore

15

16 *) Corresponding author: m.timmermans@mdx.ac.uk

17 Email addresses: Meier, Rudolf (meier@nus.edu.sg); Srivathsan, Amrita

18 (asrivathsan@gmail.com); Vogler, Alfried (a.vogler@imperial.ac.uk); Collins, Steve

19 (collinsabri@gmail.com)

20 **Abstract**

21 Polymorphic Batesian mimics exhibit multiple protective morphs that each mimic different
22 noxious models. Here we study the genomic transitions leading to the evolution of different
23 mimetic wing patterns in the polymorphic Mocker Swallowtail *Papilio dardanus*. We
24 generated a draft genome (231 Mb over 30 chromosomes) and re-sequenced individuals of
25 three morphs. Genome-wide SNP analysis revealed elevated linkage disequilibrium and
26 divergence between morphs in the regulatory region of *engrailed*, a developmental gene
27 previously implicated in the mimicry switch. The diverged region exhibits a discrete
28 chromosomal inversion (of 40 kb) relative to the ancestral orientation that is associated
29 with the *cenea* morph, but not with the bottom-recessive *hippocoonides* morph or with non-
30 mimetic allopatric populations. The functional role of this inversion in the expression of
31 the novel phenotype is currently unknown, but by preventing recombination, it allows the
32 stable inheritance of divergent alleles enabling geographic spread and local co-existence
33 of multiple adaptive morphs.

34

35

36 **Keywords:** Supergene; Batesian mimicry; butterflies; genomic rearrangement;
37 polymorphism

38 **Background**

39 Mimetic butterflies undergo profound evolutionary changes in wing patterns driven by
40 selection for a common signal deterring visual predators [1]. In Batesian mimics, which
41 imitate harmful models but are not chemically defended themselves, the fitness advantage
42 of being mimetic is a function of the predator's encounter frequency of palatable
43 individuals among unpalatable ones. Thus, a rare phenotype has a better chance of survival
44 than a frequent one and the lowered fitness with increasing abundance (negative frequency
45 dependent selection) may favour the evolution of multiple forms that each resemble a
46 different noxious model [2]. In various cases of Batesian mimics several such morphs co-
47 exist as phenotypically discrete, genetically controlled variants within a single population
48 [1, 2]. The African Mocker Swallowtail, *Papilio dardanus*, (Fig. 1) is a widely known
49 example of a polymorphic Batesian mimic. The species has played a central role in the
50 debate about the evolution of phenotypic diversity [3,4,5], starting with Trimen's work in
51 the 1860s [6]. Sometimes referred to as "the most interesting butterfly in the world" [3],
52 well over 100 variants have been named, including geographic races (subspecies) and about
53 a dozen genetically well-defined wing pattern morphs (forms) that may co-occur in
54 populations [7–9]. Females only are mimetic and both sexual dimorphism and female
55 polymorphisms presumably are driven by negative frequency dependent selection from
56 predators [10–12].

57

58 In *P. dardanus*, wing colours and patterns are controlled by a single Mendelian locus, *H*,
59 whose various alleles segregate according to a well-defined hierarchy of dominance [3,13–
60 15]. Phylogenetic analysis of subspecies and closely related species has led to the
61 conclusion that mimicry has arisen fairly recently in *P. dardanus* and that the female
62 mimetic forms are likely to have evolved from a 'male-like', presumed ancestral phenotype
63 that is still found on Madagascar where the species is monomorphic and non-mimetic (Fig.
64 1) [16]. Segregation analysis in pedigree-broods using AFLP [17] and population genetics
65 [18] have shown that the mimicry switch in *P. dardanus* is genetically linked to the
66 *engrailed-invested* locus, a region that codes for two paralogous homeodomain
67 transcription factors involved in anterior-posterior patterning [19].

68 Here, we study the genomic mechanisms that ultimately lead to the evolution of multiple
69 mimetic phenotypes in *P. dardanus*. The simple Mendelian segregation of the wing colour
70 and pattern traits led early geneticists to argue that a novel phenotype arises through a
71 single macromutation [4,20]. However, the idea of achieving perfect mimics in a single
72 step was generally dismissed by proponents of the Modern Synthesis [21,22] who argued
73 that Mendelian inheritance alone was not sufficient to prove an origin through a single
74 mutation. Instead, a two-step mechanism, first proposed by Nicholson [23], became the
75 favoured hypothesis: a new mimetic phenotype originates via an initial large-effect
76 mutation that provides at least moderate resemblance to a new mimicry model, after which
77 genetically linked secondary mutations gradually improve the resemblance [24,25]. A
78 gradual process of mimicry evolution was also favoured by computer simulations of
79 varying recombination frequency and selection strength [26]. Under this hypothesis the
80 initial mutation acts as a ‘genomic sieve’ [27] for closely linked mutations that improve
81 the resemblance to the model; selection against non-mimetic intermediates then leads to
82 the evolution of tighter linkage among genes determining colour and pattern [26,28],
83 potentially producing a ‘supergene’ controlling multiple linked mutations, such that
84 different polymorphic traits show Mendelian co-segregation [28–30].

85

86 A critical aspect of this process is that genetic recombination among functional sites is low,
87 preventing the formation of intermediates with lower fitness. Molecular genetics studies in
88 polymorphic butterflies, beetles and birds have detected associated genomic inversions as
89 a mechanism that increases linkage of co-adapted mutations [31–36]. However, the
90 importance of these inversions in the initial evolution and further diversification of
91 polymorphic forms remains unclear. Mimetic polymorphism may exist with and without
92 genomic inversions, as seen in the closely related Southeast Asian *Papilio polytes* and *P.*
93 *memnon* whose mimicry locus (in the *dsx* genomic region) is contained in an inversion
94 only in *P. polytes* [37].

95

96 To understand the genetic architecture underlying polymorphic mimicry in *P. dardanus* we
97 use comparative genomics of three female ‘forms’ (Fig. 1). Specifically, among the
98 numerous female-limited mimicry types the prevalent morph is the form *hippocoon* (f.

99 *hippocoon*), also referred to as f. *hippocoonides* in some parts of its range, which is a black-
100 and-white phenotype mimicking the danaid *Amauris niavius*. This morph is widely
101 distributed on the African mainland and is recessive to all others. A further widespread
102 phenotype is the black-and-orange form *cenea* (f. *cenea*) present mostly in specific regions
103 of Kenya (subspecies *P. d. polytrophus*) and south-eastern Africa (subspecies *P. d. cenea*).
104 Numerous other mimetic morphs co-occur within populations of these two subspecies at
105 various frequencies throughout sub-Saharan Africa [7], but populations in Madagascar are
106 always monomorphic and have been recognised as a separate subspecies, *P. d. meriones*
107 [15]. Using a newly generated draft genome sequence we assess evidence for reduced
108 recombination and genetic divergence in *P. dardanus*, and search for local rearrangements
109 that might control the phenotypic switch. This first genome wide study of *P. dardanus*
110 allows greater insight into the evolution of multiple mimicry forms and their stable
111 inheritance in populations.

112

113 **Results**

114 Draft genome and linkage map

115 A draft genome sequence was constructed using a three-generation laboratory inbred male
116 of subspecies *P. dardanus tibullus*, which was homozygous for the bottom recessive f.
117 *hippocoonides* allele (Fig. 1). We obtained an assembly of 7,365 scaffolds (N50=596,599;
118 L50 = 99) with a total length of 231,123,043 bp, which was very similar to a genome size
119 estimate of 232 Mb obtained using k-mer counts (electronic supplementary
120 material, figure S1). We were able to annotate 12,795 potential protein coding sequences
121 (CDS) and obtained Gene Ontology annotations for 8,111 putative protein coding
122 sequences. The level of completeness was similar to published draft genomes of three
123 related *Papilio* species (electronic supplementary material, table S1). The entire mimicry
124 locus *H* [17,18] was contained in two scaffolds which were merged into a 2.5 Mb scaffold
125 using information from a publicly available BAC clone sequence from the same morph
126 [18].

127

128 The scaffolds were assessed for correct assembly using co-segregation of RADseq
129 polymorphisms generated for two pedigree broods (14 and 33 F1 individuals respectively).

130 For each brood, SNPs were selected that were heterozygous in the female parent and
131 homozygous in the male parent. There is no crossing over in female Lepidoptera [38], and
132 thus all heterozygous positions on a correctly assembled scaffold should show identical
133 inheritance patterns in every offspring of a brood. Of the 7,365 scaffolds, 402 (total length:
134 193,743,404 bp) contained at least two polymorphic RADtags and could be included in
135 this analysis. Using SNP markers within the RADtags that were the furthest apart in the
136 physical maps of the scaffolds, 379 of these 402 scaffolds showed matching SNP patterns
137 in all the progeny, while discrepancies were observed for the remaining 23 scaffolds, whose
138 correct assembly could therefore not be confirmed (electronic supplementary material,
139 figure S2). The RADseq data were further used to merge the scaffolds into 29 unordered
140 bins to represent provisional groups of linked sequences. Of the 12,795 *P. dardanus* CDS,
141 9349 could be associated to one of these chromosome bins. Comparison with the well
142 annotated *Heliconius melpomene* genome largely confirmed the groups (electronic
143 supplementary material, figure S3). The 29 bins are not expected to include sex
144 chromosomes as the analyses only used SNPs that are heterozygous in the female parent
145 (female Lepidoptera are ZW, males ZZ). The data therefore suggests that *P. dardanus*
146 exhibits 30 chromosomes (29 bins plus the sex chromosomes), in accordance with an AFLP
147 study [17] and several related *Papilio* [39].

148

149 Genomics of mimicry morphs

150 Genomic differentiation of morphs was established by shotgun sequencing of specimens
151 of *hippocooides* (n=4), *cenea* (n=4), and an individual of the non-mimetic subspecies *P.*
152 *d. meriones* (Fig. 1; Table 1; electronic supplementary material, table S2). Reads were
153 mapped onto the genomic scaffolds that are longer than 100 kb (n=420). Genome-wide
154 SNP analysis of 5-kb windows (electronic supplementary material, figure S4) detected
155 elevated F_{st} values between the samples of *hippocooides* and *cenea* individuals in various
156 regions throughout the genome, including a region of ~75 kb covering the *engrailed-*
157 *invected* locus. This latter ~75 kb region also showed elevated LD. No such pattern of joint
158 elevated LD and F_{st} was observed in any of the other 420 long contigs (electronic
159 supplementary material, figure S4). These observations support the notion that within this
160 region genetic subdivision is elevated and recombination is rarer than in other regions of

161 the *P. dardanus* genome (Fig. 2). The pinpointed region did not show evidence of elevated
162 nucleotide diversity when analysing sequences from the *hippocoonides* and *cenea* morphs
163 together (Fig. 2). However, sequence divergence (estimated as p-distance) between the
164 *hippocoonides* individuals and the reference genome sequence (derived from a
165 *hippocoonides* individual) was sharply lower in the pinpointed region than for the *cenea*
166 individuals and the more divergent *P. d. meriones* (Fig. 2).

167

168 Closer inspection of the ~75 kb region revealed paired reads that were placed ~40 kb apart
169 and in opposite orientation in all four f. *cenea* individuals (electronic supplementary
170 material, figure S5). Such read-pairs were not observed in the four f. *hippocoonides*
171 samples. This indicates that the genetically diverged region contains a ~40 kb inversion
172 associated to the mimetic f. *cenea*. The inversion was not found in the non-mimetic *P. d.*
173 *meriones* from Madagascar, which indicates that the bottom-recessive mimetic f.
174 *hippocoonides* has the same arrangement as this male-like form, and therefore this specific
175 arrangement is ancestral. The four f. *cenea* specimens represented two distinct subspecies
176 from Kenya (*P. d. polytrophus* f. *cenea*) and South Africa (*P. d. cenea* f. *cenea*). The
177 sequence data furthermore indicated that the Kenyan specimens carried a non-inverted
178 allele too, suggesting they are heterozygous for f. *cenea* and f. *hippocoonides* (H_c/H_h)
179 (Table 1) which is in agreement with breeding experiments (electronic supplementary
180 material, figure S6). The South African f. *cenea* specimen was homozygous for the
181 inversion; while homozygosity for the *cenea* allele has not been confirmed by breeding, it
182 is likely because this morph is very common in this part of the species range [13].

183

184 We validated the inversion for several additional specimens of the two mimetic morphs by
185 PCR amplification with boundary-defining primers (Fig. 3). PCR fragments confirmed the
186 predicted inversion: all 4 additional f. *hippocoonides* individuals retained the arrangement
187 of the draft genome physical map (based on primer pair A-B and C-D; electronic
188 supplementary material, figure S7), consistent with the findings from the sequenced
189 individuals. Four additional f. *cenea* individuals showed the ~40 kb inversion (primer pair
190 A-C and B-D). These *cenea* females also showed the A-B and C-D fragments of the
191 reference map, indicating they are heterozygous (H_c/H_h). Fisher exact tests for association

192 between phenotype and inversion were highly significant ($P < 0.0001$) (electronic
193 supplementary material, table S3).

194

195 To test whether the genomic region surrounding *engrailed* and *invected* recombines freely,
196 we used RAD data for two pedigree broods (homozygous f. *hippocoonides*), using 199
197 SNPs at sites with variants in the male but not the female parents in the ~2.5 Mb scaffold
198 containing *engrailed-invected*. We detected seven recombination events in one brood and
199 two in the other (Fig. 2; electronic supplementary material, table S4), for a relatively high
200 recombination frequency of 7.8 cM/Mb (9 recombination events in 47 offspring, or 19.1
201 cM, over a distance of 2,458 Mb of the scaffold). These results were similar to those
202 presented by Clark et al. [17], who analysed a cross between a heterozygous male (H_h/H_c)
203 and a homozygous female (H_h/H_h) and reported 5 recombination events between male-
204 informative AFLP markers ACT and PD (highlighted on Fig. 2), which flank the *engrailed-*
205 *invected* region, after scoring 35 F1 individuals.

206

207 **Discussion**

208 Our genomic analysis revealed a 40 kb inversion in *P. dardanus* at ~6,800 bp upstream of
209 the *engrailed* start codon, which differentiates the haplotype associated with the
210 *hippocoonides* and *cenea* morphs, and coincides with localized peaks in LD and F_{st}
211 between haplotypes of these morphs. Mimicry loci have been postulated to consist of
212 several tightly linked, epistatically interacting loci that in concert determine adaptive
213 phenotypes (i.e. acting as a supergene) [29]. Such interaction of multiple sites requires
214 regions of reduced recombination preventing the segregation of co-adapted loci, which was
215 broadly confirmed in recent work demonstrating inversions in mimicry-linked genomic
216 regions of other mimetic butterflies [31,32,34,37,40]. We have not determined the
217 sequence of the *cenea* (H_c) allele and do not know whether several independent mutations
218 are required for the switch between f. *cenea* and f. *hippocoonides* to happen, but the fact
219 that a recombination suppressing inversion exists suggests a genomic architecture
220 consistent with the supergene hypothesis (although due to the linkage of mutations within
221 the inversion, it will not be possible to uncover the functional sites without functional
222 studies).

223

224 The inversion in *P. dardanus* is small, compared to those associated with the mimicry loci
225 in the Batesian mimic *P. polytes* and the Müllerian mimic *H. numata*, which stretch over
226 130 kb and at least 400 kb, respectively, and in those species result in allelic divergence in
227 several protein coding genes. The *P. dardanus* inversion also differs from those species by
228 the fact that it is found in an extended regulatory region apparently devoid of protein coding
229 sequences. The region contains various enhancer sequences [41,42] that in other species
230 have been shown to exert *cis*-regulatory control of both *engrailed* and *invected* and
231 therefore likely affect unlinked genes determining the colour pattern, as initially envisioned
232 by the ‘regulatory hypothesis’ of Nijhout [8,43]. *Invected* also contains an intronic
233 microRNA (miR-2768) conserved in Lepidoptera (Fig. 3; electronic supplementary
234 material, figure S8), which has been shown to downregulate *cubitus interruptus (ci)*, a gene
235 that determines patterning of the wing primordia via the *hedgehog* signalling pathway in
236 nymphalid butterflies [44].

237

238 In *P. dardanus* the universally recessive *hippocoonides* form, despite being mimetic,
239 apparently retains the presumed ancestral orientation found in the allopatric and genetically
240 divergent (Fig. 2) Madagascan subspecies. This demonstrates that an inversion is not
241 critical for the origin of mimetic forms, as also observed in *P. memnon* [37]. However,
242 when multiple mimetic female forms are found in sympatry chromosomal inversions will
243 assist stable segregation of divergent phenotypes, as has been shown for *P. polytes* [37].
244 Here we show that inversions are associated with multiple sympatric mimicry forms also
245 in *P. dardanus* in mainland Africa. Balanced inversion polymorphisms may be maintained
246 in populations by negative frequency-dependent selection (Type II polymorphisms of
247 [45]). In addition, the spread of an advantageous phenotype is promoted when it is
248 associated with an inversion (e.g. see [46]). The f. *cenea*-linked inversion has spread widely
249 across the African continent and across subspecies boundaries, as evident from the
250 presence of the *cenea* morph in *P. d. polytrophus* from Kenya and *P. d. cenea* from South
251 Africa, geographically separated by at least 3000 km (Fig. 1). The fact that the same
252 inversion is associated with the *cenea* morph in different subspecies, adds support for its
253 role in defining the phenotype.

254

255 It still needs to be confirmed if the regulatory region of *engrailed-invected* plays any
256 functional role in determining the pleiotropic changes of the wing. However, *P. dardanus*
257 would not be unique in having regulatory changes underlying polymorphic mimicry. A
258 recent study on the nymphalid *Hypolimnas misippus*, which displays sex-limited mimicry,
259 revealed a 10-kb intergenic region upstream of the *Sox5/6* gene to be strongly associated
260 to the wing phenotype, suggesting that a cis-regulatory element plays a role in pattern
261 determination [47]. Inversions in an intron of the *pannier* locus determining colour
262 polymorphism in a ladybird beetle have been shown to affect gene expression and to
263 underlie phenotypic differences among colour morphs [48], also supporting *cis*-regulation
264 of the phenotype through inversions of non-coding regions. If the 40-kb inversion in *P.*
265 *dardanus* has *cis*-regulatory effects on the expression of one or more of *engrailed*, *invected*
266 and miR-2768 (and possibly the adjacent gene *orange*), the genetic architecture of the
267 region may be particularly conducive to the evolution of novel phenotypes. Thus, new
268 inversions may provide the hypothesized major-effect shifts through their regulatory
269 function that impacts the mosaic of pattern and colour elements of the wing.

270

271 Other morphs now need to be investigated for chromosomal rearrangements in this region,
272 and may not exclusively involve inversions, given a previously reported duplication of
273 *engrailed-invected* and a few neighbouring genes closely associated with one of the other
274 *P. dardanus* female forms (f. *lamborni*) [18]. Preliminary results also suggest a genomic
275 rearrangement in an individual of f. *planemoides*, which indicates that recombination-
276 suppressing reordering of the *engrailed* region is an integral part of the evolution of new
277 mimicry morphs. Determination of the phenotype likely works in concert with other
278 changes in the *engrailed-invected* region, such as those in the first exon of *engrailed* found
279 in the top-dominant f. *poultoni* and f. *planemoides* that exhibit a statistically significant
280 overrepresentation of non-synonymous substitutions indicative of diversifying selection
281 [49]. These divergent sites are outside of the newly detected inversions, perhaps suggesting
282 that for some morphs a combination of the divergent *engrailed* coding region and the
283 upstream inversion are required for correct specification of the phenotype. The presence of
284 chromosomal rearrangements might suppress the recombination frequency even beyond

285 the inverted region, as already evident from the wider region of high LD and F_{st} extending
286 to ~75 kb (Fig. 3). Accordingly, recombinants producing maladaptive intermediate
287 phenotypes should exist but are rare, and such non-mimetic phenotypes may persist locally.

288

289 With each study of polymorphic systems, now including the prototypical *P. dardanus*, the
290 understanding of how discrete adaptive phenotypes evolve and are maintained in natural
291 populations improves: all currently described butterfly mimicry loci show the expected
292 signatures of allelic divergence, indicating that complex phenotypes indeed require
293 multiple sites and probably evolved in smaller steps. However, the mechanisms by which
294 tight linkage is achieved differ, as do the loci that determine the phenotypic switch.
295 Inversions are not necessary, but helpful to promote the capture of alleles under positive
296 selection, because they contribute to maintaining the alleles that would otherwise break up
297 genetically linked sites and lead to poor fitness. They might also contribute to the genetic
298 variation producing novel phenotypes, although for *P. dardanus*, the challenge remains to
299 determine any role of the inversion in gene expression or the regulation of downstream
300 pathways, in order to track the macro- and micro-mutations on the evolutionary trajectory
301 towards stable polymorphisms of mimicry forms.

302

303 **Methods**

304 Genome sequencing, assembly and annotation

305 The draft genome sequence was generated from an inbred male specimen of subspecies *P.*
306 *dardanus tibullus* (electronic supplementary material, table S2). Genomic DNA was used
307 for construction of Illumina TruSeq libraries (insert sizes of 300 bp and 800 bp) and a
308 Nextera mate-pair (MP) library prior to sequencing on Illumina platforms, followed by
309 standard procedures for adapter removal and quality trimming. GenomeScope [50] was
310 used to estimate genome size by obtaining the mean of the k-mer count distribution.
311 Sequencing errors were corrected using QUAKE v0.3.5 [51] using JELLYFISH v1.1.11
312 for k-mer counting [52]. Using an estimated genome size of 200 Mb we used k=17 for error
313 correction and Quake was run using default parameters. Genome assembly was conducted
314 using Platanus v. 1.2.4 [53], using only paired-end data for generating initial contigs, while
315 using mate-pair data for subsequent steps as recommended by the developers (number of

316 links for scaffolding = 10). For improving accuracy of the assembly, removing redundancy
317 and further scaffolding we used HaploMerger2 (Release 20151124) [54]. WindowMasker
318 v1.0.0. was first used to mask repetitive regions and all-against-all whole genome
319 alignments were then obtained using LASTZ and reciprocally-best whole-genome
320 alignments using chainNET to generate an improved haploid assembly.

321

322 The haploid assembly was further scaffolded using SSPACE v3.0 (number of links = 10),
323 using both paired-end and mate-pair libraries. Insert sizes were estimated by using the
324 library *_insFreq.tsv file generated by Platanus. This assembly was further refined by the
325 removal of tandem assembly errors and gaps in the assembly were closed using GapCloser.
326 Lastly, to remove scaffolds that could be from contaminations, we built a custom database
327 consisting of representative bacterial genomes from NCBI RefSeq 6, four reference
328 genomes for *Papilio* sp. (*P. machaon*: GCA_001298355.1; *P. polytes*: GCF_000836215.1;
329 *P. xuthus*: GCF_000836235.1; and *P. glaucus*: GCA_000931545.1) and a reference human
330 genome (GRCh38.p7). All scaffolds were searched against the reference database using
331 BLASTN with e-value of 1E-5. Genome completeness of this draft genome and other
332 *Papilio* genomes was assessed using BUSCO version 3 [55]. The assembly was annotated
333 using MAKER2 [56] with gene predictors trained by AUGUSTUS [57] using the BUSCO
334 ortholog set. Predicted protein and RNA sequences from genome assemblies of other
335 *Papilio* species were used as evidence. For functional annotations, protein sequences were
336 matched to SWISS-PROT [58] using BLASTP (E-value 1e-5) and subject to InterProScan
337 [59] for detection of protein signatures.

338

339 Scaffold clustering and mimicry locus genetic recombination

340 Sets of unordered linked scaffolds (“chromosome bins”) were obtained by SNP segregation
341 in RADseq data generated for two *P. dardanus* broods of 14 and 35 offspring. RAD library
342 construction was performed using *Pst*I restriction digestion and barcoded libraries were
343 sequenced (100 bp single-end reads). Reads were de-multiplexed using the process_radtags
344 script of the package Stacks and subsequently mapped onto the genomic scaffolds using
345 bbmap (sourceforge.net/projects/bbmap/) (setting: ambiguous=toss local=t). The resulting
346 SAM files were sorted and converted to BAM files using SAMtools. Picard-tools-1.117

347 (<http://broadinstitute.github.io/picard>) was used to add read group information and merge
348 the individual files of each brood into a single BAM file (i.e. one merged file per brood).
349 These files were then converted to VCF format using the HaplotypeCaller program of
350 GATK. Positions with 18x coverage or less for at least one of the samples within a brood
351 were removed using SNPsift and the file converted to OneMap format using the
352 `vcf_to_onemap_input` version 1.0 python script and positions heterozygous in the female
353 parent (Onemap notation: 'a,b') and homozygous in the male (Onemap notation: 'a,a')
354 parent (OneMap crosstype: D1.10) were extracted. For each scaffold with at least two
355 segregating RADtags we tested co-segregation of the most distant SNPs to detect
356 inconsistencies in segregation pattern, indicating incorrect assemblies. Co-segregation of
357 SNPs was subsequently used to group scaffolds into linkage groups. CDS from linkage
358 groups were compared to the *Heliconius melpomene* genome (version 2) and the positions
359 of sequence matches on 21 *H. melpomene* chromosomes were recorded. The Perl GD::SVG
360 library was used to visualise the positions of sequence matches. The RAD data was also
361 used to investigate recombination within the scaffold containing *engrailed-inverted*. SNPs
362 homozygous in the female parent and heterozygous in the male parent were extracted and
363 inspected manually for evidence of genomic recombination.

364

365 Population genomics of the *P. dardanus* supergene

366 Genomic data for eight specimens (Table 1) were mapped onto all scaffolds >100 kb using
367 the BWA-MEM algorithm [60], merging the data for *hippocoonides* and *cenea* specimens
368 into two separate files. Mean coverage was calculated for both for 5 kb sliding windows
369 using SAMtools depth function and a custom perl script. To remove repetitive regions,
370 sites with >400x coverage were masked for this analysis. The two files were merged and
371 Kelly's ZnS statistic (the average of the LD measure r^2 calculated between all pairs of
372 SNPs) [61], nucleotide diversity (π), and mean p-distance to the reference genome
373 sequence were calculated using PopBam (sliding window 5 kb) [62]. F_{st} values were
374 calculated using VCFtools 0.1.12 [63] contrasting the *hippocoonides* and *cenea* morphs
375 (window size 5 kb). PCR was used to validate a genomic inversions (Fig. 2) using
376 additional *hippocoonides* and *cenea* specimens (electronic supplementary material, figure
377 S7) and the following primers: A) 5'-gktgtcgattttgctgcta-3', B) 5'-

378 aactaaaactrtyagagacacgcaa-3', C) 5'-tyaacgggtcagacaagttt-3' and D) 5'-
379 amatggcgatgractgmca-3'. Fisher exact tests (two-tailed) were performed to test for
380 association between phenotype and presence of an inversion (taking the dominance
381 hierarchy into account) (electronic supplementary material, table S3).

382

383 **Acknowledgements**

384 The authors thank Martin Thompson for access to DNA from butterfly samples from South
385 Africa. Rebecca Clark kindly provided laboratory cross details on the sequenced *Papilio*
386 *dardanus* f. *cenea* from Kenya. We also would like to thank all reviewers who have
387 commented on various submissions of the manuscript. Sequencing libraries were
388 constructed and sequenced at the NHM London, the Department of Biochemistry
389 (University of Cambridge), and Genepool (University of Edinburgh).

390

391 **Data accessibility**

392 Sequence data that support the findings of this study have been deposited in GenBank with
393 the accession codes PRJNA451133, PRJNA600400, PRJNA600373 and SAMN05819004.

394

395 **Author' contributions**

396 MJTNT participated in the design of the study, carried out the molecular lab work, analysed
397 data, and drafted the manuscript; AS carried out bioinformatics analyses and drafted the
398 manuscript; SC participated in the design of the study and provided specimens; RM
399 provided bioinformatics resources and critically revised the manuscript; APV participated
400 in the design of the study and drafted the manuscript. All authors gave final approval for
401 publication and agree to be held accountable for the work performed therein.

402

403 **Competing interests**

404 There are no competing interests.

405

406 **Funding**

407 This study was funded by NERC Postdoctoral Fellowship NE/I021578/1 (to MJTNT) and
408 NERC NE/F006225/1 (to APV). AS was supported by SEABIG (R-154-000-648-646 and
409 R-154-000-648-733).

410 **References**

- 411 1. Kunte K. 2009 The diversity and evolution of Batesian mimicry in *Papilio swallowtail*
412 butterflies. *Evolution* **63**, 2707–2716.
- 413 2. Joron M, Mallet JLB. 1998 Diversity in mimicry: paradox or paradigm? *Trends Ecol. Evol.* **13**,
414 461–466. (doi:10.1016/S0169-5347(98)01483-9)
- 415 3. Poulton EB. 1924 *Papilio dardanus*. The most interesting butterfly in the world. *IJ E Afr*
416 *Uganda Nat Hist Soci* **20**, 4–22.
- 417 4. Davis FR. 2009 *Papilio dardanus*: The natural animal from the experimentalist's point of
418 view. In *Descended from Darwin: Insights into the history of evolutionary studies, 1900-1970*
419 (eds J Cain, M Ruse), pp. 221–242. Philadelphia: American Philosophical Society.
- 420 5. Ford EB. 1936 The genetics of *Papilio dardanus* Brown (Lep.). *Trans. R. Entomol. Soc. Lond.*
421 **85**, 435–466.
- 422 6. Trimén R. 1869 On some remarkable mimetic analogies among African butterflies. *Trans.*
423 *Linn. Soc. Lond.* **26**, 497–522.
- 424 7. Thompson MJ, Timmermans MJTN. 2014 Characterising the phenotypic diversity of *Papilio*
425 *dardanus* wing patterns using an extensive museum collection. *PLoS ONE* **9**, e96815.
- 426 8. Nijhout HF. 2003 Polymorphic mimicry in *Papilio dardanus*: mosaic dominance, big effects,
427 and origins. *Evol. Dev.* **5**, 579–592.
- 428 9. Clarke CA, P. M. Sheppard. 1960 The genetics of *Papilio dardanus* Brown. II. Races *dardanus*,
429 *polytrophus*, *meseres*, and *tibullus*. *Genetics* **45**, 439–456.
- 430 10. Cook SE. 1994 Mate Choice in the Polymorphic African Swallowtail Butterfly, *Papilio*
431 *dardanus*: Male-like Females May Avoid Sexual Harassment. *Anim. Behav.* **47**, 389–397.
- 432 11. Turner JRG. 1978 Why male butterflies are non-mimetic: natural selection, group selection,
433 modification and sieving. *Biol. J. Linn. Soc.* **10**, 385–432.
- 434 12. O'Donald P. 1969 The selective coefficients that keep modifying genes in a population.
435 *Genetics* **62**, 435–444.
- 436 13. Clarke CA, P. M. Sheppard. 1959 The genetics of *Papilio dardanus* Brown. I. Race *cenea* from
437 South Africa. *Genetics* **44**, 1347–1358.
- 438 14. Clarke CA, P. M. Sheppard. 1960 The genetics of *Papilio dardanus* Brown. II. Races *dardanus*,
439 *polytrophus*, *meseres*, and *tibullus*. *Genetics* **45**, 439–456.

- 440 15. Clarke CA, Sheppard PM. 1960 The genetics of *Papilio dardanus* Brown. III. Race *antinatorii*
441 from Abyssinia and race *meriones* from Madagascar. *Genetics* **45**, 683–698.
- 442 16. Timmermans MJTN, Thompson MJ, Collins S, Vogler AP. 2017 Independent evolution of
443 sexual dimorphism and female-limited mimicry in swallowtail butterflies (*Papilio dardanus*
444 and *Papilio phorcas*). *Mol. Ecol.* **26**, 1273–1284. (doi:10.1111/mec.14012)
- 445 17. Clark R, S. M. Brown, S. C. Collins, C. D. Jiggins, D. G. Heckel, A. P. Vogler. 2008 Colour
446 pattern specification in the Mocker Swallowtail *Papilio dardanus*: the transcription factor
447 *invected* is a candidate for the mimicry locus *H*. *Proc. R. Soc. B* **275**, 1181–1188.
- 448 18. Timmermans MJTN *et al.* 2014 Comparative genomics of the mimicry switch in *Papilio*
449 *dardanus*. *Proc. R. Soc. B Biol. Sci.* **281**, 20140465–20140465. (doi:10.1098/rspb.2014.0465)
- 450 19. Peel AD, Telford MJ, Akam M. 2006 The evolution of hexapod *engrailed*-family genes:
451 evidence for conservation and concerted evolution. *Proc. R. Soc. B-Biol. Sci.* **273**, 1733–
452 1742.
- 453 20. Punnett RC. 1915 *Mimicry in Butterflies*. London and Edinburgh: Cambridge University Press.
- 454 21. Fisher RA. 1927 On some objections to mimicry theory; statistical and genetic. *Trans. R.*
455 *Entomol. Soc.* **75**, 269–274.
- 456 22. Ford EB. 1975 *Ecological Genetics, Fourth ed.* London: Chapman and Hall.
- 457 23. Nicholson AJ. 1927 A new theory of mimicry in insects. *Aust. J. Zool.* **5**, 10–104.
- 458 24. Baxter SW, Johnston SE, Jiggins CD. 2009 Butterfly speciation and the distribution of gene
459 effect sizes fixed during adaptation. *Heredity* **102**, 57–65.
- 460 25. Joron M. 2003 Mimicry. In *Encyclopedia of insects* (eds RT Carde, VH Resh), pp. 714–726.
461 New York: Academic Press.
- 462 26. Charlesworth D, Charlesworth B. 1975 Theoretical genetics of Batesian mimicry. 2. Evolution
463 of supergenes. *J. Theor. Biol.* **55**, 305–324.
- 464 27. Turner JRG. 1977 Butterfly mimicry: the genetical evolution of an adaptation. *Evol. Biol.* **10**,
465 163–206.
- 466 28. Clarke CA, P. M. Sheppard. 1960 Super-genes and mimicry. *Heredity* **14**, 175–185.
- 467 29. Thompson MJ, Jiggins CD. 2014 Supergenes and their role in evolution. *Heredity* **113**, 1–8.
468 (doi:10.1038/hdy.2014.20)
- 469 30. Charlesworth D. 2016 The status of supergenes in the 21st century: recombination
470 suppression in Batesian mimicry and sex chromosomes and other complex adaptations.
471 *Evol. Appl.* **9**, 74–90. (doi:10.1111/eva.12291)

- 472 31. Joron M *et al.* 2011 Chromosomal rearrangements maintain a polymorphic supergene
473 controlling butterfly mimicry. *Nature* **477**, 203-U102.
- 474 32. Kunte K, Zhang W, Tenger-Trolander A, Palmer DH, Martin A, Reed RD, Mullen SP, Kronforst
475 MR. 2014 *doublesex* is a mimicry supergene. *Nature* **507**, 229–232.
476 (doi:10.1038/nature13112)
- 477 33. Küpper C *et al.* 2016 A supergene determines highly divergent male reproductive morphs in
478 the ruff. *Nat. Genet.* **48**, 79–83. (doi:10.1038/ng.3443)
- 479 34. Nishikawa H *et al.* 2015 A genetic mechanism for female-limited Batesian mimicry in *Papilio*
480 butterfly. *Nat. Genet.* **47**, 405–409. (doi:10.1038/ng.3241)
- 481 35. Zhang W, Westerman E, Nitzany E, Palmer S, Kronforst MR. 2017 Tracing the origin and
482 evolution of supergene mimicry in butterflies. *Nat. Commun.* **8**. (doi:10.1038/s41467-017-
483 01370-1)
- 484 36. Tuttle EM *et al.* 2016 Divergence and functional degradation of a sex chromosome-like
485 supergene. *Curr. Biol.* **26**, 344–350. (doi:10.1016/j.cub.2015.11.069)
- 486 37. Iijima T, Kajitani R, Komata S, Lin C-P, Sota T, Itoh T, Fujiwara H. 2018 Parallel evolution of
487 Batesian mimicry supergene in two *Papilio* butterflies, *P. polytes* and *P. memnon*. *Sci. Adv.* **4**,
488 eaao5416. (doi:10.1126/sciadv.aao5416)
- 489 38. Turner JRG, P. M. Sheppard. 1975 Absence of crossing-over in female butterflies
490 (*Heliconius*). *Journal Hered.* **34**, 265–269.
- 491 39. Meenu Sadhotra. 2016 A chromosomal investigation of three species of *Papilio*
492 (*Papilionidae*:*Lepidoptera*). *Asian J. Anim. Sci.* **11**, 135–139.
- 493 40. Jay P, Whibley A, Frézal L, Rodríguez de Cara MÁ, Nowell RW, Mallet J, Dasmahapatra KK,
494 Joron M. 2018 Supergene evolution triggered by the introgression of a chromosomal
495 inversion. *Curr. Biol.* **28**, 1839-1845.e3. (doi:10.1016/j.cub.2018.04.072)
- 496 41. Cheng Y, Brunner AL, Kremer S, DeVido SK, Stefaniuk CM, Kassis JA. 2014 Co-regulation of
497 *invected* and *engrailed* by a complex array of regulatory sequences in *Drosophila*. *Dev. Biol.*
498 **395**, 131–143. (doi:10.1016/j.ydbio.2014.08.021)
- 499 42. Gustavson E, Goldsborough AS, Ali Z, Kornberg TB. 1996 The *Drosophila engrailed* and
500 *invected* Genes: Partners in Regulation, Expression and Function. *Genetics* **142**, 893–906.
- 501 43. Nijhout HF. 1991 *The development and evolution of butterfly wing patterns*. Washington:
502 Smithsonian Institution Press.
- 503 44. Quah S, Hui JHL, Holland PWH. 2015 A Burst of miRNA Innovation in the Early Evolution of
504 Butterflies and Moths. *Mol. Biol. Evol.* **32**, 1161–1174. (doi:10.1093/molbev/msv004)
- 505 45. Faria R, Johannesson K, Butlin RK, Westram AM. 2019 Evolving Inversions. *Trends Ecol. Evol.*
506 **34**, 239–248. (doi:10.1016/j.tree.2018.12.005)

- 507 46. Kirkpatrick M, Barton N. 2006 Chromosome inversions, local adaptation and speciation.
508 *Genetics* **173**, 419–434.
- 509 47. VanKuren NW, Massardo D, Nallu S, Kronforst MR. 2019 Butterfly Mimicry Polymorphisms
510 Highlight Phylogenetic Limits of Gene Reuse in the Evolution of Diverse Adaptations. *Mol.*
511 *Biol. Evol.* **36**, 2842–2853. (doi:10.1093/molbev/msz194)
- 512 48. Ando T *et al.* 2018 Repeated inversions within a pannier intron drive diversification of
513 intraspecific colour patterns of ladybird beetles. *Nat. Commun.* **9**, 3843.
514 (doi:10.1038/s41467-018-06116-1)
- 515 49. Thompson MJ, Timmermans MJ, Jiggins CD, Vogler AP. 2014 The evolutionary genetics of
516 highly divergent alleles of the mimicry locus in *Papilio dardanus*. *BMC Evol. Biol.* **14**, 140.
517 (doi:10.1186/1471-2148-14-140)
- 518 50. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC.
519 2017 GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics*
520 **33**, 2202–2204. (doi:10.1093/bioinformatics/btx153)
- 521 51. Kelley DR, Schatz MC, Salzberg SL. 2010 Quake: quality-aware detection and correction of
522 sequencing errors. *Genome Biol.* **11**, R116. (doi:10.1186/gb-2010-11-11-r116)
- 523 52. Marçais G, Kingsford C. 2011 A fast, lock-free approach for efficient parallel counting of
524 occurrences of k-mers. *Bioinformatics* **27**, 764–770. (doi:10.1093/bioinformatics/btr011)
- 525 53. Kajitani R *et al.* 2014 Efficient de novo assembly of highly heterozygous genomes from
526 whole-genome shotgun short reads. *Genome Res.* **24**, 1384–1395.
527 (doi:10.1101/gr.170720.113)
- 528 54. Huang S, Kang M, Xu A. 2017 HaploMerger2: rebuilding both haploid sub-assemblies from
529 high-heterozygosity diploid genome assembly. *Bioinformatics* **33**, 2577–2579.
530 (doi:10.1093/bioinformatics/btx220)
- 531 55. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015 BUSCO:
532 assessing genome assembly and annotation completeness with single-copy orthologs.
533 *Bioinformatics* **31**, 3210–3212. (doi:10.1093/bioinformatics/btv351)
- 534 56. Holt C, Yandell M. 2011 MAKER2: an annotation pipeline and genome-database
535 management tool for second-generation genome projects. *BMC Bioinformatics* **12**, 491.
536 (doi:10.1186/1471-2105-12-491)
- 537 57. Stanke M, Morgenstern B. 2005 AUGUSTUS: a web server for gene prediction in eukaryotes
538 that allows user-defined constraints. *Nucleic Acids Res.* **33**, W465–W467.
539 (doi:10.1093/nar/gki458)
- 540 58. UniProt Consortium T. 2018 UniProt: the universal protein knowledgebase. *Nucleic Acids*
541 *Res.* **46**, 2699–2699. (doi:10.1093/nar/gky092)

- 542 59. Jones P *et al.* 2014 InterProScan 5: genome-scale protein function classification.
543 *Bioinformatics* **30**, 1236–1240. (doi:10.1093/bioinformatics/btu031)
- 544 60. Li H, Durbin R. 2009 Fast and accurate short read alignment with Burrows-Wheeler
545 transform. *Bioinformatics* **25**, 1754–1760. (doi:10.1093/bioinformatics/btp324)
- 546 61. Kelly JK. 1997 A test of neutrality based on interlocus associations. *Genetics* **146**, 1197–
547 1206.
- 548 62. Garrigan D. 2013 POPBAM: Tools for evolutionary analysis of short read sequence
549 alignments. *Evol. Bioinforma.* **9**, EBO.S12751. (doi:10.4137/EBO.S12751)
- 550 63. Danecek P *et al.* 2011 The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158.
551 (doi:10.1093/bioinformatics/btr330)

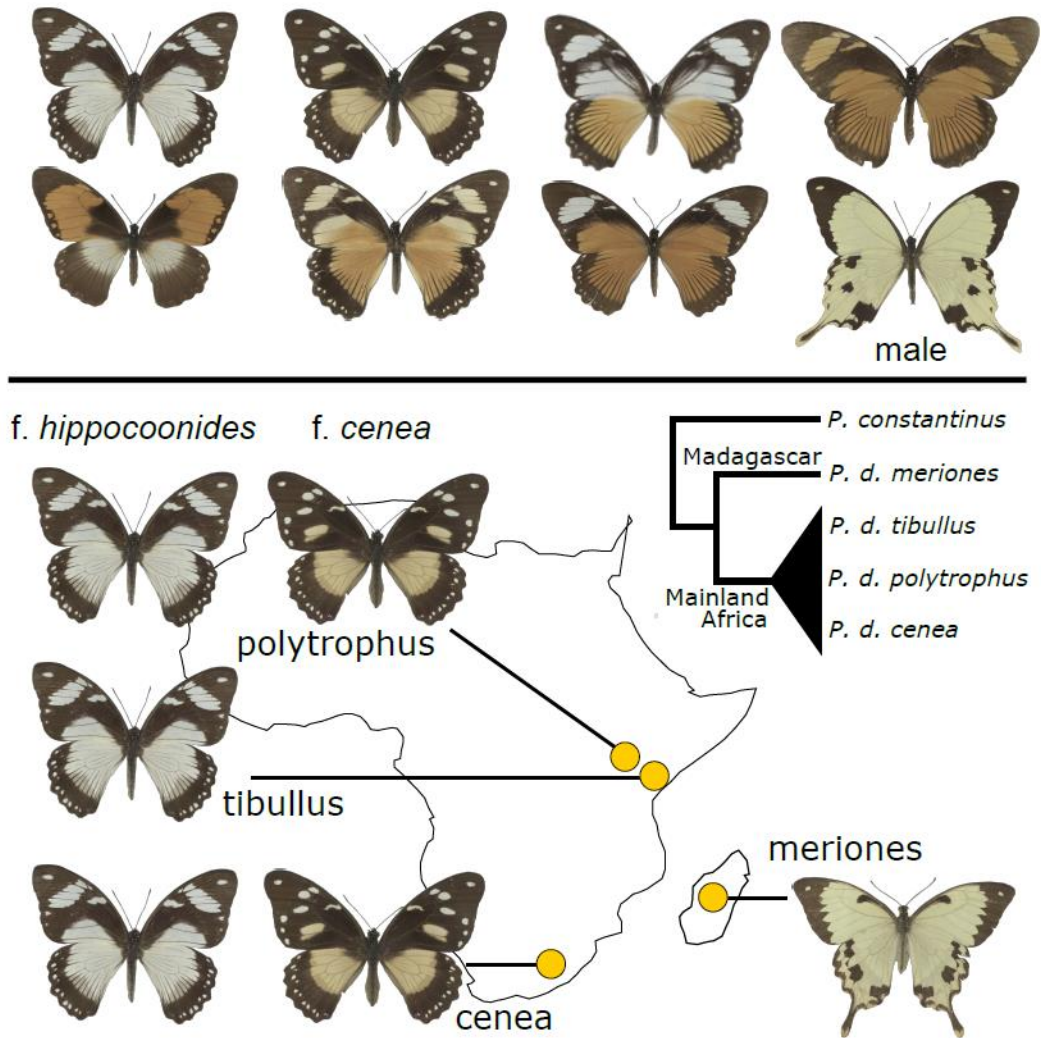
552
|
553

554 **Table 1:** Samples used for sequencing.

555

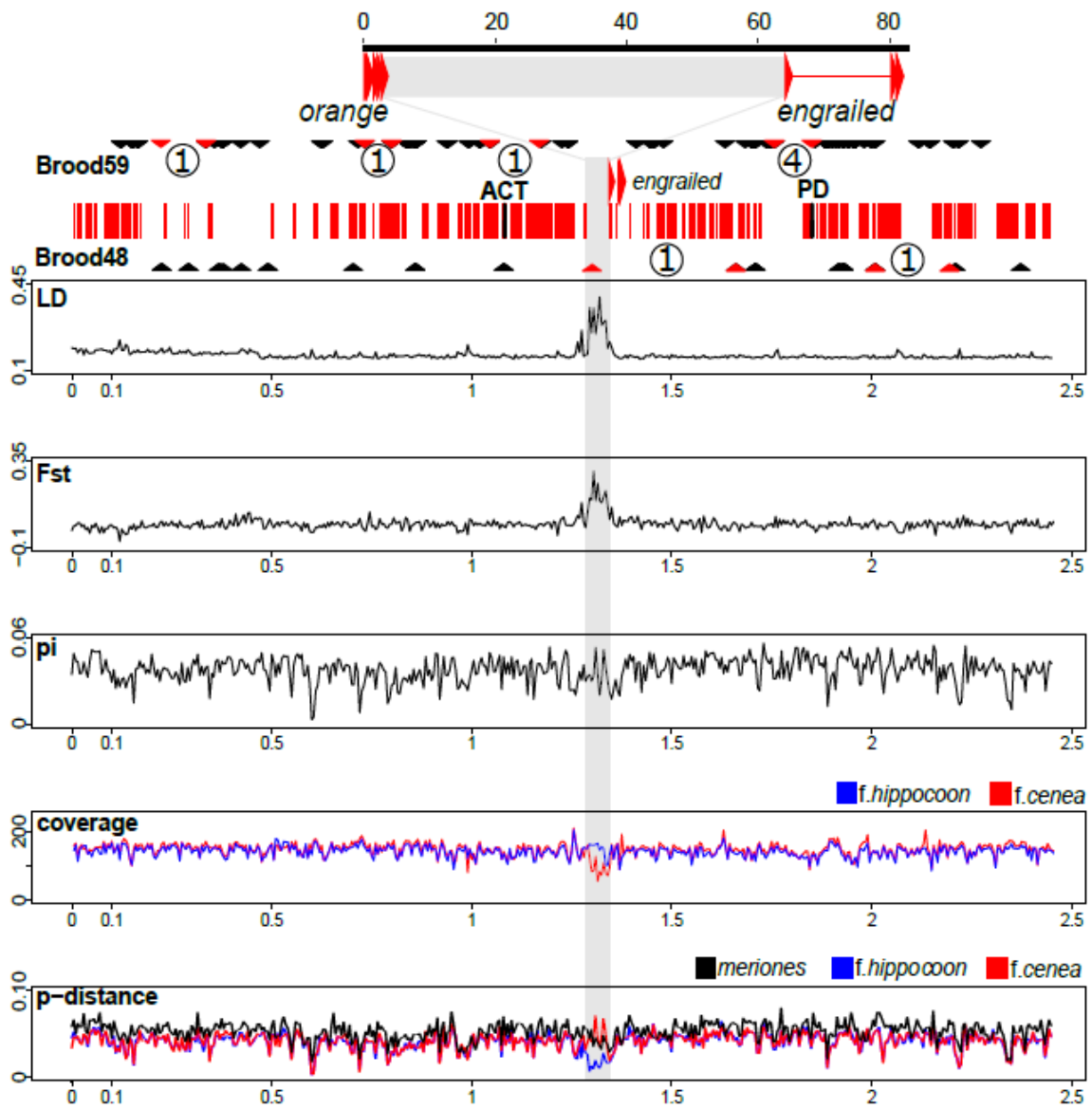
Voucher number	Geographic origin	Subspecies	Phenotype	H genotype	# Paired-End Reads	Estimated coverage	Reference orientation	40 kb inversion
BMNH746848	Kenya	polytrophus	hippocoonides	H _h / H _h	38039853	41	x	
BMNH746826	Kenya	polytrophus	hippocoonides	H _h / H _h	55548066	60	x	
BMNH847389	Kenya	polytrophus	hippocoonides	H _h / H _h	35666867	39	x	
BMNH746846	South Africa	cenea	hippocoonides	H _h / H _h	40600315	44	x	
BMNH746453	Kenya	polytrophus	cenea	H _c /H _h	49512794	54	x	x
BMNH746764	Kenya	polytrophus	cenea	H _c /H _h	56085254	61	x	x
Troph-c-02-46	Kenya	polytrophus	cenea	H _c /H _h	42184635	46	x	x
BMNH847353	South Africa	cenea	cenea	H _c /?	39434400	43		x
BMNH740167	Madagascar	meriones	meriones		31242775	34	x	

556 H genotype: H_c = H_{cenea} , H_h = H_{hippocoonides}. #Paired End Reads: number of raw reads generated for each specimen. Estimated coverage
 557 is calculated via: (number of raw reads * read length) / length of genome assembly. Read length was 125 bp. The actual coverage is
 558 expected to be lower due to not all reads passing quality control and the presence of contamination. “x” in the last three columns indicates
 559 whether the specimen carried an allele with the reference orientation and the 40 kb inversion.



561

562 **Figure 1:** Phenotypic variation in *Papilio dardanus* and samples used. Top: Seven female
 563 forms and a male. Bottom: Origin of samples for sequencing and population genetic
 564 analyses, from four subspecies: *P. dardanus polytrophus* (Kenya), *P. dardanus tibullus*
 565 (Kenya), *P. dardanus cenea* (South-Africa), and *P. dardanus meriones* (Madagascar). The
 566 specimen of subspecies *P. dardanus tibullus* was used for the construction of the draft
 567 genome sequence. The tree depicts the relationships among these four subspecies and is
 568 based on a tree presented in [16]. Three female forms were analysed: *hippocoonides*, *cenea*,
 569 and ‘male-like’.



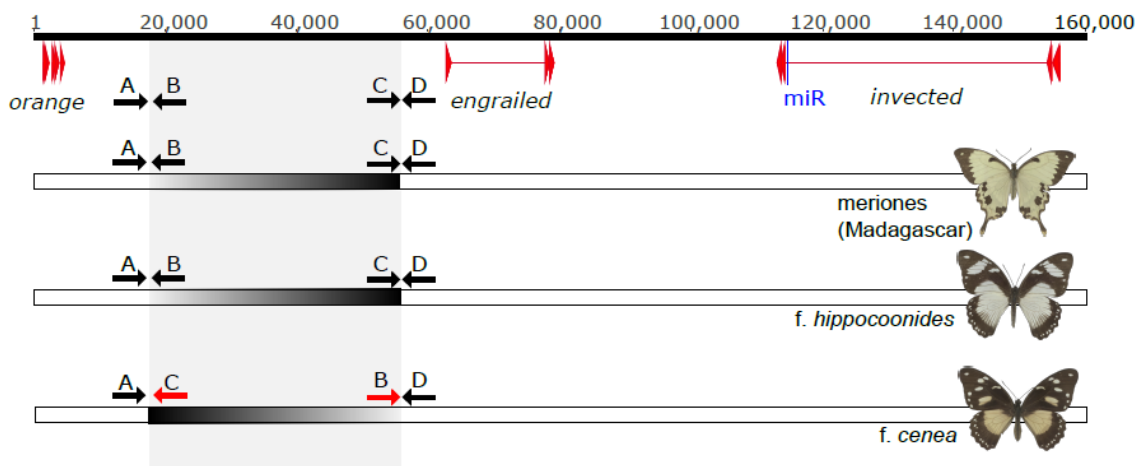
570

571 **Figure 2:** Population genomic analysis of the full *engrailed-inverted* containing scaffold.
 572 Thin vertical red lines: exons of various genes in the region. Note the exons of *engrailed*
 573 shown by large arrows and the upstream region marked in grey. Brood 59 and Brood 48:
 574 Recombination events in pedigree broods. SNPs from RADseq data for two broods are
 575 mapped on the *engrailed-inverted* scaffold, shown by black triangles. Red triangles mark
 576 the intervals with confirmed recombination events, and the number of recombination
 577 events within these intervals are circled. The central band in the figure shows the map of
 578 the scaffold with exons (red arrows) and the upstream region of *engrailed* (grey). ACT and
 579 PD indicate the position of the AFLP markers of [17]. Linkage disequilibrium (LD; Kelly's

580 ZnS statistic), F_{st} , nucleotide diversity (π), coverage and p-distance (to the reference
 581 genome) for the scaffold, calculated for the *f. cenea* and *f. hippocooides* samples in 5 kb
 582 windows. Coverage and p-distance was calculated separately for the four *cenea* and for
 583 four *hippocooides* specimens. The p-distance to the reference genome is also given for
 584 the *P. dardanus meriones* sample. Scales are in million base pairs.

585

586



587

588 **Figure 3:** Length and relative position of the inversions in the upstream regulatory region
 589 of *engrailed*. At the top, the map of the *engrailed-invested* region is shown, with short
 590 arrows indicating exons and the miR-2768 [44] shown in blue. Below the map is the
 591 direction of boundary-defining primers. The grey shading indicates the extent of the 40 kb
 592 inversion associated with *f. cenea*. For each of the forms, dark grey – light grey shading is
 593 used to indicate directionality of the 40 kb region. Scale is in base pair.

594

595