

Research



Cite this article: Timmermans MJTN, Srivathsan A, Collins S, Meier R, Vogler AP. 2020 Mimicry diversification in *Papilio dardanus* via a genomic inversion in the regulatory region of *engrailed–invested*. *Proc. R. Soc. B* **287**: 20200443. <http://dx.doi.org/10.1098/rspb.2020.0443>

Received: 26 February 2020

Accepted: 31 March 2020

Subject Category:

Evolution

Subject Areas:

evolution, genomics, genetics

Keywords:

supergene, Batesian mimicry, butterflies, genomic rearrangement, polymorphism

Author for correspondence:

Martijn J. T. N. Timmermans
e-mail: m.timmermans@mdx.ac.uk

Electronic supplementary material is available online at <https://doi.org/10.6084/m9.figshare.c.4955048>.

Mimicry diversification in *Papilio dardanus* via a genomic inversion in the regulatory region of *engrailed–invested*

Martijn J. T. N. Timmermans^{1,2,3}, Amrita Srivathsan⁴, Steve Collins⁶, Rudolf Meier^{4,5} and Alfried P. Vogler^{1,2}

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

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

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

⁴Department of Biological Sciences, and ⁵Lee Kong Chian Natural History Museum, National University of Singapore, Singapore

⁶African Butterfly Research Institute, Nairobi, Kenya

id MJTNT, 0000-0002-5024-9053; RM, 0000-0002-4452-2885

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

1. Background

Mimetic butterflies undergo profound evolutionary changes in wing patterns driven by selection for a common signal deterring visual predators [1]. In Batesian mimics, which imitate harmful models but are not chemically defended themselves, the fitness advantage of being mimetic is a function of the predator's encounter frequency of palatable individuals among unpalatable ones. Thus, a rare phenotype has a better chance of survival than a frequent one and fitness is lowered with increasing abundance (negative frequency-dependent selection), which may favour the evolution of multiple forms that each resemble a different noxious model [2]. In various cases of Batesian mimics, several such morphs coexist as phenotypically discrete, genetically controlled variants within a single population [1,2]. The African Mocker Swallowtail, *Papilio dardanus*, is a widely known example of a polymorphic Batesian mimic (figure 1). The species has played a central role in the debate about the evolution of phenotypic diversity [4–6], starting with Trimen's work in the 1860s [7]. Sometimes referred to as 'the most interesting butterfly in the world' [4], well over 100 variants have been named, including geographic races (subspecies) and about a dozen genetically well-defined wing pattern morphs (forms) that may co-occur in populations [8–10]. Only the females are mimetic and both sexual dimorphism and female polymorphisms presumably are driven by negative frequency-dependent selection from predators [11–13].

In *P. dardanus*, wing colours and patterns are controlled by a single Mendelian locus, *H*, whose various alleles segregate according to a well-defined hierarchy of dominance [4,10,14,15]. Phylogenetic analysis of subspecies and closely related species has led to the conclusion that mimicry has arisen fairly recently in

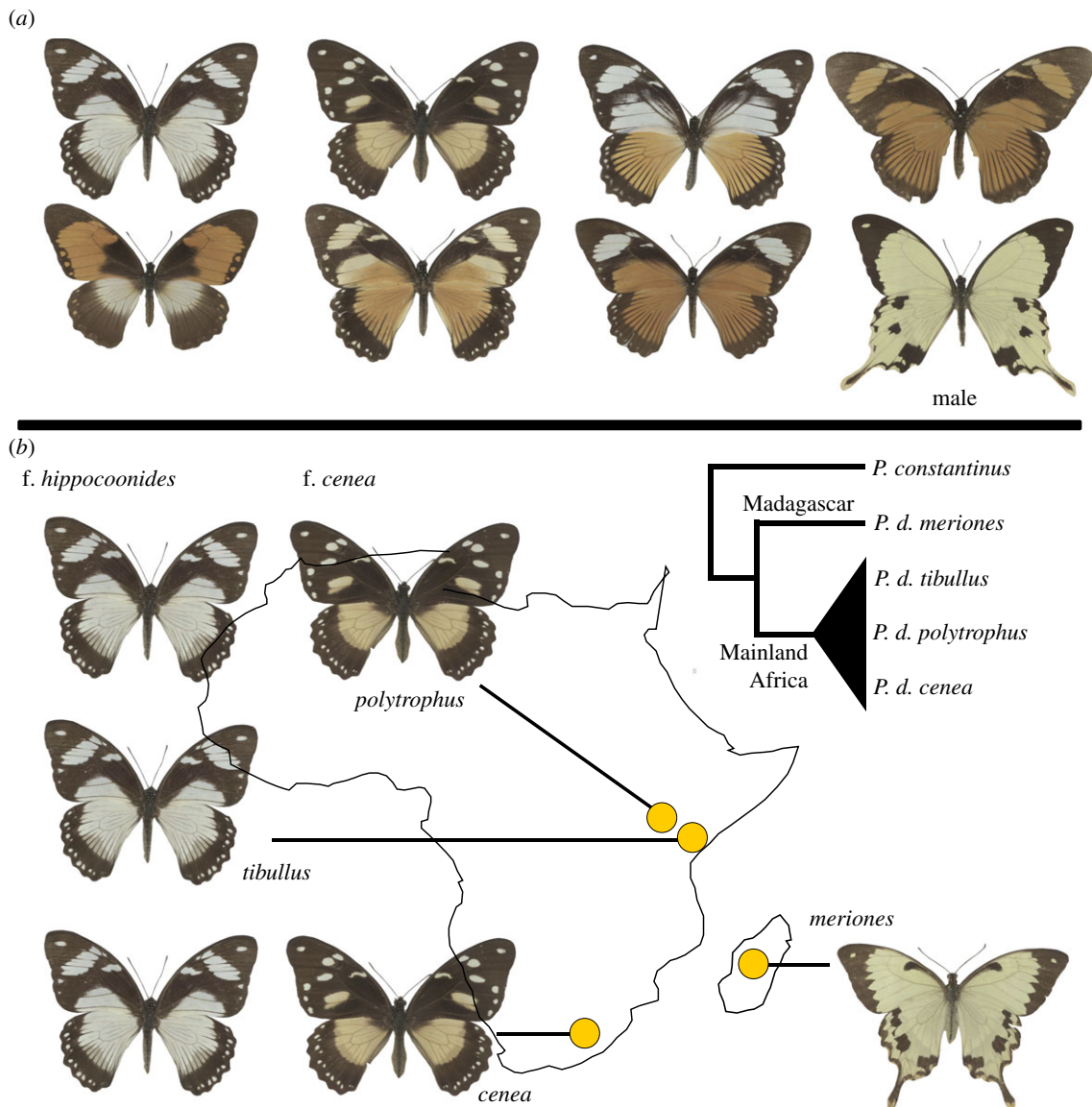


Figure 1. Phenotypic variation in *Papilio dardanus* and samples used. (a) Seven female forms and a male. (b) Origin of samples for sequencing and population genetic analyses, from four subspecies: *P. dardanus polytrophus* (Kenya), *P. dardanus tibullus* (Kenya), *P. dardanus cenea* (South Africa), and *P. dardanus meriones* (Madagascar). The specimen of subspecies *P. dardanus tibullus* was used for the construction of the draft genome sequence. The tree depicts the relationships among these four subspecies and is based on a tree presented in [3]. Three female forms were analysed: *hippocooides*, *cenea*, and 'male-like'.

P. dardanus and that the female mimetic forms are likely to have evolved from a 'male-like', presumed ancestral phenotype that is still found on Madagascar where the species is monomorphic and non-mimetic (figure 1) [3]. Segregation analysis in pedigree broods using amplified fragment length polymorphism (AFLP) [16] and population genetics [17] have shown that the mimicry switch in *P. dardanus* is genetically linked to the *engrailed–invected* locus, a region that codes for two paralogous homeodomain transcription factors involved in anterior–posterior patterning [18].

Here, we studied the genomic mechanisms that ultimately lead to the evolution of multiple mimetic phenotypes in *P. dardanus*. The simple Mendelian segregation of the wing colour and pattern traits led early geneticists to argue that a novel phenotype arises through a single macromutation [5,19]. However, the idea of achieving perfect mimics in a single step was generally dismissed by proponents of the Modern Synthesis [20,21] who argued that Mendelian inheritance alone was not sufficient to prove an origin through a single mutation. Instead, a two-step mechanism, first proposed by Nicholson [22], became the favoured hypothesis: a new mimetic phenotype originates via an initial large-effect mutation

that provides at least moderate resemblance to a new mimicry model, after which genetically linked secondary mutations gradually improve the resemblance [23,24]. A gradual process of mimicry evolution was also favoured by computer simulations of varying recombination frequency and selection strength [25]. Under this hypothesis, after the initial mutation the surrounding genomic region acts as a 'genomic sieve' [26] for closely linked mutations that improve the resemblance to the model; selection against non-mimetic intermediates then leads to the evolution of tighter linkage among genes determining colour and pattern [25,27], potentially producing a 'supergene' controlling multiple linked mutations, such that different polymorphic traits show Mendelian co-segregation [27–29].

A critical aspect of this process is that genetic recombination among functional sites is low, preventing the formation of intermediates with lower fitness. Molecular genetics studies in polymorphic butterflies, beetles, and birds have detected associated genomic inversions as a mechanism that increases linkage of co-adapted mutations [30–35]. However, the importance of these inversions in the initial evolution and further diversification of polymorphic forms remains unclear. Mimetic

Table 1. Samples used for sequencing. H genotype: $H_c = H_{cenaia}$, $H_h = H_{hippocooides}$. #Paired-end reads: number of raw Illumina reads generated for each specimen. Estimated coverage is calculated via: (number of raw reads * read length) / length of genome assembly. Read length was 125 bp. The actual coverage is expected to be lower due to not all reads passing quality control and the presence of contamination. 'x' in the last two columns indicates whether the specimen carried an allele with the reference orientation and the 40 kb inversion.

voucher number	geographic origin	subspecies	phenotype	H genotype	# paired-end reads	estimated coverage	reference orientation	40 kb inversion
BMNH746848	Kenya	<i>polytrophus</i>	hippocooides	H_h/H_h	38039853	41	x	
BMNH746826	Kenya	<i>polytrophus</i>	hippocooides	H_h/H_h	55548066	60	x	
BMNH847389	Kenya	<i>polytrophus</i>	hippocooides	H_h/H_h	35666867	39	x	
BMNH746846	South Africa	<i>cenea</i>	hippocooides	H_h/H_h	40600315	44	x	
BMNH746453	Kenya	<i>polytrophus</i>	<i>cenea</i>	H_c/H_h	49512794	54	x	x
BMNH746764	Kenya	<i>polytrophus</i>	<i>cenea</i>	H_c/H_h	56085254	61	x	x
Troph-c-02-46	Kenya	<i>polytrophus</i>	<i>cenea</i>	H_c/H_h	42184635	46	x	x
BMNH847353	South Africa	<i>cenea</i>	<i>cenea</i>	$H_c/?$	39434400	43		x
BMNH740167	Madagascar	<i>meriones</i>	<i>meriones</i>		31242775	34	x	

polymorphism may exist with and without genomic inversions, as seen in the closely related Southeast Asian *Papilio polytes* and *P. memnon* whose mimicry locus (in the *dsx* genomic region) is contained in an inversion only in *P. polytes* [36].

To understand the genetic architecture underlying polymorphic mimicry in *P. dardanus*, we used comparative genomics of three female 'forms' (figure 1). Specifically, among the numerous female-limited mimicry types the prevalent morph is the form *hippocoon* (f. *hippocoon*), also referred to as f. *hippocooides* in some parts of its range, which is a black-and-white phenotype mimicking the danaid *Amauris niavivus*. This morph is widely distributed on the African mainland and is recessive to all others. A further widespread phenotype is the black-and-orange form *cenea* (f. *cenea*) present mostly in specific regions of Kenya (subspecies *P. d. polytrophus*) and south-eastern Africa (subspecies *P. d. cenea*). Numerous other mimetic morphs co-occur within populations of these two subspecies at various frequencies throughout sub-Saharan Africa [8], but populations in Madagascar are always monomorphic and have been recognized as a separate subspecies, *P. d. meriones* [15]. Using a newly generated draft genome sequence, we assessed evidence for reduced recombination and genetic divergence in *P. dardanus* and searched for local rearrangements that might control the phenotypic switch. This first genome-wide study of *P. dardanus* allows greater insight into the evolution of multiple mimicry forms and their stable inheritance in populations.

2. Results

(a) Draft genome

A draft genome sequence was constructed using a three-generation laboratory inbred male of subspecies *P. dardanus tibullus*, which was homozygous for the bottom-recessive f. *hippocooides* allele (figure 1). We obtained an assembly of 7365 scaffolds ($N_{50} = 596\,599$; $L_{50} = 99$) with a total length of 231 123 043 bp, which was very similar to a genome size estimate of 232 Mb obtained using k-mer counts (electronic supplementary material, figure S1). We were able to annotate 12 795 potential protein-coding sequences (CDS) and obtained Gene Ontology annotations for 8111 putative protein-coding sequences. The level of completeness was similar to published

draft genomes of three related *Papilio* species (electronic supplementary material, table S1). The entire mimicry locus *H* [16,17] was contained in two scaffolds which were merged into a 2.5 Mb scaffold using information from a publicly available BAC clone sequence from the same morph [17].

The scaffolds were assessed for correct assembly using co-segregation of RADseq polymorphisms generated for two pedigree broods (14 and 33 F1 individuals, respectively). For each brood, SNPs were selected that were heterozygous in the female parent and homozygous in the male parent. There is no crossing over in female Lepidoptera [37], and thus all heterozygous positions on a correctly assembled scaffold should show identical inheritance patterns in every offspring of a brood. Of the 7365 scaffolds, 402 (total length: 193 743 404 bp) contained at least two polymorphic RADtags and could be included in this analysis. Using SNP markers within the RADtags that were the furthest apart in the physical maps of the scaffolds, 379 of these 402 scaffolds showed matching SNP patterns in all the progeny, while discrepancies were observed for the remaining 23 scaffolds, whose correct assembly could, therefore, not be confirmed (electronic supplementary material, figure S2). The RADseq data were further used to merge the scaffolds into 29 unordered bins to represent provisional groups of linked sequences. Of the 12 795 *P. dardanus* CDS, 9349 could be associated to one of these chromosome bins. Comparison with the well-annotated *Heliconius melpomene* genome largely confirmed the groups (electronic supplementary material, figure S3). The 29 bins are not expected to include sex chromosomes as the analyses only used SNPs that are heterozygous in the female parent (female Lepidoptera are ZW, males ZZ). The data, therefore, suggests that *P. dardanus* exhibits 30 chromosomes (29 bins plus the sex chromosomes), in accordance with an AFLP study [16] and several related *Papilio* [38].

(b) Genomics of mimicry morphs

Genomic differentiation of morphs was established by shotgun sequencing of specimens of f. *hippocooides* ($n = 4$), f. *cenea* ($n = 4$). Data on a previously sequenced individual of the non-mimetic subspecies *P. d. meriones* was also included (figure 1; table 1; electronic supplementary material, table S2). Reads were mapped onto the genomic scaffolds that were

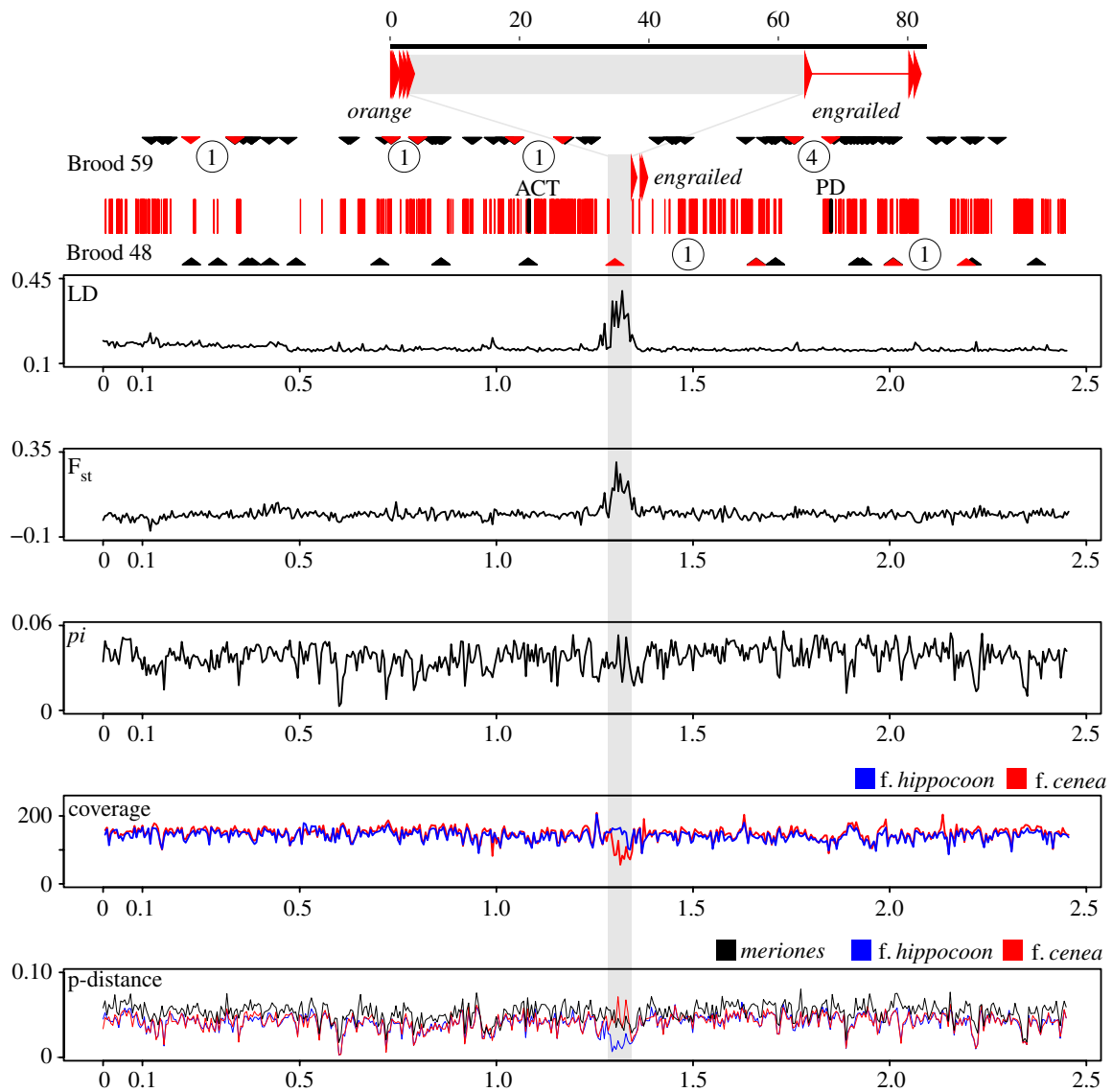


Figure 2. Population genomic analysis of the full *engrailed–invected* containing scaffold. The exons of *orange* and *engrailed* are shown by large arrows and the upstream region marked in grey. Brood 59 and Brood 48: recombination events in pedigree broods. SNPs from RADseq data for two broods are mapped on the *engrailed–invected* scaffold, shown by black triangles. Red triangles mark the intervals with confirmed recombination events, and the number of recombination events within these intervals are circled. The central band in the figure shows the map of the scaffold with exons (thin vertical red lines) and the upstream region of *engrailed* (grey). ACT and PD indicate the position of the AFLP markers of [16]. Linkage disequilibrium (LD; Kelly's ZnS statistic), F_{st} , nucleotide diversity (π), coverage and p-distance (to the reference genome) for the scaffold, calculated for the *f. cenea* and *f. hippocooides* samples in 5 kb windows. Coverage and p-distance were calculated separately for the four *cenea* and for four *hippocooides* specimens. The p-distance to the reference genome is also given for the *P. dardanus meriones* sample. Scales are in million base pairs. (Online version in colour.)

longer than 100 kb ($n = 420$). Genome-wide SNP analysis of 5 kb windows (electronic supplementary material, figure S4) detected elevated F_{st} values between the *f. hippocooides* and *f. cenea* individuals in various regions throughout the genome, including a region of approximately 75 kb covering the *engrailed–invected* locus. This approximately 75 kb region also showed elevated linkage disequilibrium (LD). No such pattern of joint elevated LD and F_{st} was observed in any of the other 420 long contigs (electronic supplementary material, figure S4). These observations support the notion that within this region genetic subdivision is elevated and recombination is rarer than in other regions of the *P. dardanus* genome (figure 2). The pinpointed region did not show evidence of elevated nucleotide diversity when analysing sequences from the *f. hippocooides* and *f. cenea* morphs together (figure 2). However, sequence divergence (estimated as p-distance) between the *f. hippocooides* individuals and the reference genome sequence (derived from an *f. hippocooides* individual) was sharply

lower in the pinpointed region than for the *cenea* individuals and the more divergent *P. d. meriones* (figure 2).

Closer inspection of the approximately 75 kb region revealed paired reads that were placed approximately 40 kb apart and in opposite orientation in all four *f. cenea* individuals (electronic supplementary material, figure S5). Such read-pairs were not observed in the four *f. hippocooides* samples. This indicates that the genetically diverged region contains an approximately 40 kb inversion associated to the mimetic *f. cenea*. The inversion was not found in the non-mimetic *P. d. meriones* from Madagascar, which indicates that the bottom-recessive mimetic *f. hippocooides* has the same arrangement as this male-like form, and therefore this specific arrangement is ancestral. The four *f. cenea* specimens represented two distinct subspecies from Kenya (*P. d. polytrophus f. cenea*) and South Africa (*P. d. cenea f. cenea*). The sequence data furthermore indicated that the Kenyan specimens carried a non-inverted allele too, suggesting they are heterozygous for *f. cenea* and

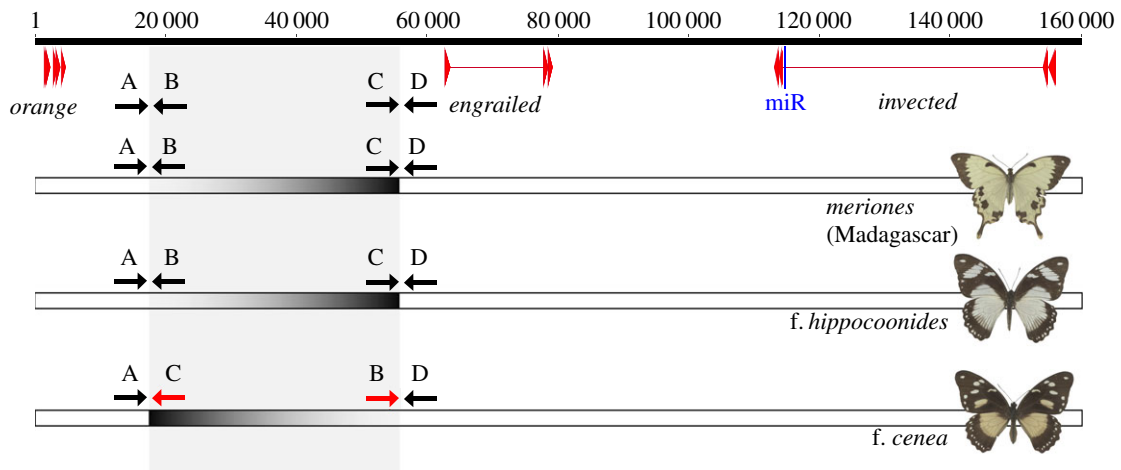


Figure 3. Length and relative position of the inversions in the upstream regulatory region of *engrailed*. At the top, the map of the *engrailed*–*invected* region is shown, with short arrows indicating exons and miR-2768 [39] shown in blue. Below the map is the direction of boundary-defining primers. The grey shading indicates the extent of the 40 kb inversion associated with *f. cenea*. For each of the forms, dark grey–light grey shading is used to indicate directionality of the 40 kb region. Scale is in base pair. (Online version in colour.)

f. hippocooides (H_c/H_h) (table 1) which is in agreement with breeding experiments (electronic supplementary material, figure S6). The South African *f. cenea* specimen was homozygous for the inversion; while homozygosity for the *cenea* allele has not been confirmed by breeding, it is likely because this morph is very common in this part of the species range [14].

We validated the inversion for several additional specimens of the two mimetic morphs by PCR amplification with boundary-defining primers (figure 3). PCR fragments confirmed the predicted inversion: all four additional *f. hippocooides* individuals retained the arrangement of the draft genome physical map (based on primer pair A–B and C–D; electronic supplementary material, figure S7), consistent with the findings from the sequenced individuals. Four additional *f. cenea* individuals showed the approximately 40 kb inversion (primer pair A–C and B–D). These *cenea* females also showed the A–B and C–D fragments of the reference map, indicating they are heterozygous (H_c/H_h). Fisher exact tests for association between phenotype and inversion were highly significant ($p < 0.0001$) (electronic supplementary material, table S3).

To test whether the genomic region surrounding *engrailed* and *invected* recombines freely, we used RAD data for two pedigree broods (homozygous *f. hippocooides*), using 199 SNPs at sites with variants in the male but not the female parents in the approximately 2.5 Mb scaffold containing *engrailed*–*invected*. We detected seven recombination events in one brood and two in the other (figure 2; electronic supplementary material, table S4), for a relatively high recombination frequency of 7.8 cM/Mb (9 recombination events in 47 offspring, or 19.1 cM, over a distance of 2458 Mb of the scaffold). These results were similar to those presented by Clark *et al.* [16], who analysed a cross between a heterozygous male (H_h/H_c) and a homozygous female (H_h/H_h) and reported five recombination events between male-informative AFLP markers ACT and PD (highlighted on figure 2), which flank the *engrailed*–*invected* region, after scoring 35 F1 individuals.

3. Discussion

Our genomic analysis revealed a 40 kb inversion in *P. dardanus* at approximately 6800 bp upstream of the *engrailed* start codon, which differentiates the haplotype associated with the

hippocooides and *cenea* morphs, and coincides with localized peaks in LD and F_{st} between haplotypes of these morphs. Mimicry loci have been postulated to consist of several tightly linked, epistatically interacting loci that in concert determine adaptive phenotypes (i.e. acting as a supergene) [28]. Such interaction of multiple sites requires regions of reduced recombination preventing the segregation of co-adapted loci, which was broadly confirmed in recent work demonstrating inversions in mimicry-linked genomic regions of other mimetic butterflies [30,31,33,36,40]. We have not determined the sequence of the *cenea* (H_c) allele and do not know whether several independent mutations are required for the switch between *f. cenea* and *f. hippocooides* to happen, but the fact that a recombination-suppressing inversion exists suggests a genomic architecture consistent with the supergene hypothesis (although due to the linkage of mutations within the inversion, it will not be possible to uncover the functional sites without functional studies).

The inversion in *P. dardanus* is small, compared to those associated with the mimicry loci in the Batesian mimic *P. polytes* and the Müllerian mimic *H. numata*, which stretch over 130 kb and at least 400 kb, respectively, and in those species result in allelic divergence in several protein-coding genes. The *P. dardanus* inversion also differs from those species by the fact that it is found in an extended regulatory region apparently devoid of protein-coding sequences. The region contains various enhancer sequences [41,42] that in other species have been shown to exert *cis*-regulatory control of both *engrailed* and *invected* and therefore likely affect unlinked genes determining the colour pattern, as initially envisioned by the ‘regulatory hypothesis’ of Nijhout [9,43]. *Invected* also contains an intronic microRNA (miR-2768) conserved in Lepidoptera (figure 3; electronic supplementary material, figure S8), which has been shown to downregulate *cubitus interruptus* (*ci*), a gene that determines patterning of the wing primordia via the *hedgehog* signalling pathway in nymphalid butterflies [39].

In *P. dardanus*, the universally recessive *hippocooides* form, despite being mimetic, apparently retains the presumed ancestral orientation found in the allopatric and genetically divergent (figure 2) Madagascan subspecies. This demonstrates that an inversion is not critical for the origin of mimetic forms, as also observed in *P. memnon* [36]. However, when multiple mimetic female forms are found in sympatry chromosomal

inversions will assist stable segregation of divergent phenotypes, as has been shown for *P. polytes* [36]. Here, we show that inversions are associated with multiple sympatric mimicry forms also in *P. dardanus* in mainland Africa. Balanced inversion polymorphisms may be maintained in populations by negative frequency-dependent selection (Type II polymorphisms of [44]). In addition, the spread of an advantageous phenotype is promoted when it is associated with an inversion (e.g. see [45]). The *f. cenea*-linked inversion has spread widely across the African continent and across subspecies boundaries, as evident from the presence of the *cenea* morph in *P. d. polytrophus* from Kenya and *P. d. cenea* from South Africa, geographically separated by at least 3000 km (figure 1). The fact that the same inversion is associated with the *cenea* morph in different subspecies adds support for its role in defining the phenotype.

It still needs to be confirmed if the regulatory region of *engrailed–invected* plays any functional role in determining the pleiotropic changes of the wing. However, *P. dardanus* would not be unique in having regulatory changes underlying polymorphic mimicry. A recent study on the nymphalid *Hypolimnas misippus*, which displays sex-limited mimicry, revealed a 10 kb intergenic region upstream of the *Sox5/6* gene to be strongly associated with the wing phenotype, suggesting that a *cis*-regulatory element plays a role in pattern determination [46]. Inversions in an intron of the *pannier* locus determining colour polymorphism in a ladybird beetle have been shown to affect gene expression and to underlie phenotypic differences among colour morphs [47], also supporting *cis*-regulation of the phenotype through inversions of non-coding regions. If the 40 kb inversion in *P. dardanus* has *cis*-regulatory effects on the expression of one or more of *engrailed*, *invected* and miR-2768 (and possibly the adjacent gene *orange*), the genetic architecture of the region may be particularly conducive to the evolution of novel phenotypes. Thus, new inversions may provide the hypothesized major-effect shifts through their regulatory function that impacts the mosaic of pattern and colour elements of the wing.

Other morphs now need to be investigated for chromosomal rearrangements in this region and may not exclusively involve inversions, given a previously reported duplication of *engrailed–invected* and a few neighbouring genes closely associated with one of the other *P. dardanus* female forms (*f. lamborni*) [17]. Preliminary results also suggest a genomic rearrangement in an individual of *f. planemoides*, which indicates that recombination-suppressing reordering of the *engrailed* region is an integral part of the evolution of new mimicry morphs. Determination of the phenotype likely works in concert with other changes in the *engrailed–invected* region, such as those in the first exon of *engrailed* found in the top-dominant *f. poultoni* and *f. planemoides* that exhibit a statistically significant overrepresentation of non-synonymous substitutions indicative of diversifying selection [48]. These divergent sites are outside of the newly detected inversions, perhaps suggesting that for some morphs a combination of the divergent *engrailed* coding region and the upstream inversion are required for correct specification of the phenotype. The presence of chromosomal rearrangements might suppress the recombination frequency even beyond the inverted region, as already evident from the wider region of high LD and F_{st} extending to approximately 75 kb (figure 3). Accordingly, recombinants producing maladaptive intermediate phenotypes should exist but be rare, and such non-mimetic phenotypes may persist locally.

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

4. Methods

(a) Genome sequencing, assembly, and annotation

The draft genome sequence was generated from an inbred male specimen of subspecies *P. dardanus tibullus* (electronic supplementary material, table S2). Genomic DNA was used for the construction of Illumina TruSeq libraries (insert sizes of 300 bp and 800 bp) and a Nextera mate-pair (MP) library prior to sequencing on Illumina platforms, followed by standard procedures for adapter removal and quality trimming. GenomeScope [49] was used to estimate genome size by obtaining the mean of the *k*-mer count distribution. Sequencing errors were corrected using QUAKE v0.3.5 [50] using JELLYFISH v1.1.11 for *k*-mer counting [51]. Using an estimated genome size of 200 Mb, we used *k* = 17 for error correction and Quake was run using default parameters. Genome assembly was conducted using Platanus v. 1.2.4 [52], using only paired-end data for generating initial contigs, while using mate-pair data for subsequent steps as recommended by the developers (number of links for scaffolding = 10). For improving accuracy of the assembly, removing redundancy, and further scaffolding, we used HaploMerger2 (Release 20 151 124) [53]. WindowMasker v1.0.0. was first used to mask repetitive regions and all-against-all whole-genome alignments were then obtained using LASTZ and reciprocally best whole-genome alignments using chainNET to generate an improved haploid assembly.

The haploid assembly was further scaffolded using SSPACE v3.0 (number of links = 10), using both paired-end and mate-pair libraries. Insert sizes were estimated by using the library *_insFreq.tsv file generated by Platanus. This assembly was further refined by the removal of tandem assembly errors and gaps in the assembly were closed using GapCloser. Lastly, to remove scaffolds that could be from contaminations, we built a custom database consisting of representative bacterial genomes from NCBI RefSeq 6, four reference genomes for *Papilio* sp. (*P. machaon*: GCA_001298355.1; *P. polytes*: GCF_000836215.1; *P. xuthus*: GCF_000836235.1; and *P. glaucus*: GCA_000931545.1) and a reference human genome (GRCh38.p7). All scaffolds were searched against the reference database using BLASTN with an *e*-value of 1E-5. Genome completeness of this draft genome and other *Papilio* genomes was assessed using BUSCO version 3 [54]. The assembly was annotated using MAKER2 [55] with gene predictors trained by AUGUSTUS [56] using the BUSCO ortholog set. Predicted protein and RNA sequences from genome assemblies of other *Papilio* species were used as evidence. For functional annotations, protein sequences were matched to SWISS-PROT [57] using BLASTP

(E-value 1×10^{-5}) and subject to InterProScan [58] for detection of protein signatures.

(b) Scaffold clustering and mimicry locus genetic recombination

Sets of unordered linked scaffolds (chromosome bins) were obtained by SNP segregation in RADseq data generated for two *P. dardanus* broods of 14 and 33 offspring. RAD library construction was performed using *Pst*I restriction digestion and barcoded libraries were sequenced (100 bp single-end reads). Reads were de-multiplexed using the `process_radtags` script of the package Stacks [59], and subsequently mapped onto the genomic scaffolds using `bbmap` (sourceforge.net/projects/bbmap/) (setting: `ambiguous = toss local = t`). The resulting SAM files were sorted and converted to BAM files using SAMtools [60]. Picard-tools-1.117 (<http://broadinstitute.github.io/picard>) was used to add read group information and merge the individual files of each brood into a single BAM file (i.e. one merged file per brood). These files were then converted to VCF format using the HaplotypeCaller program of GATK [61]. Positions with $18\times$ coverage or less for at least one of the samples within a brood were removed using SNPsift [62] and the file converted to OneMap format [63], using the `vcf_to_onemap_input` version 1.0 python script (<https://github.com/UT-Python-Fall-2013/Class-Projects>) and positions heterozygous in the female parent (Onemap notation: 'a,b') and homozygous in the male (Onemap notation: 'a,a') parent (OneMap crosstype: D1.10) were extracted. For each scaffold with at least two segregating RADtags, we tested co-segregation of the most distant SNPs to detect inconsistencies in segregation pattern, indicating incorrect assemblies. Co-segregation of SNPs was subsequently used to group scaffolds into linkage groups. CDS from linkage groups were compared to the *Heliconius melpomene* genome (version 2) and the positions of sequence matches on 21 *H. melpomene* chromosomes were recorded. The Perl GD::SVG library was used to visualize the positions of sequence matches. The RAD data was also used to investigate recombination within the scaffold containing *engrailed-inverted*. SNPs homozygous in the female parent and heterozygous in the male parent were extracted and inspected manually for evidence of genomic recombination.

(c) Population genomics of the *P. dardanus* supergene

Genomic data for nine specimens (table 1) were mapped onto all scaffolds greater than 100 kb using the BWA-MEM algorithm [64], merging the data for *f. hippocooides* and *f. cenea* specimens

into two separate files. Mean coverage was calculated for both for 5 kb sliding windows using SAMtools depth function and a custom perl script. To remove repetitive regions, sites with greater than $400\times$ coverage were masked for this analysis. The two files were merged, and Kelly's ZnS statistic (the average of the LD measure r^2 calculated between all pairs of SNPs) [65], nucleotide diversity (π), and mean p-distance to the reference genome sequence were calculated using PopBam (sliding window 5 kb) [66]. F_{st} values were calculated using VCFtools 0.1.12 [67] contrasting the *hippocooides* and *cenea* morphs (window size 5 kb). PCR was used to validate a genomic inversion (figure 2) using additional *hippocooides* and *cenea* specimens (electronic supplementary material, figure S7) and the following primers: A) 5'-GKTGTGCGATTTTGGCGCTA-3', B) 5'-AACTAAACTRTRTYAGAGACACGCAA-3', C) 5'-TYAACCGGGTCAGACAAGTTT-3', and D) 5'-AMATGGCGATGRACCTGMCGA-3'. Fisher exact tests (two-tailed) were performed to test for association between phenotype and the presence of an inversion (taking the dominance hierarchy into account) (electronic supplementary material, table S3).

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

Authors' contributions. M.J.T.N.T. participated in the design of the study, carried out the molecular laboratory work, analysed data, and drafted the manuscript; A.S. carried out bioinformatics analyses and drafted the manuscript; S.C. participated in the design of the study and provided specimens; R.M. provided bioinformatics resources and critically revised the manuscript; A.P.V. participated in the design of the study and drafted the manuscript. All authors gave final approval for publication and agree to be held accountable for the work performed therein.

Competing interests. There are no competing interests.

Funding. This study was funded by NERC Postdoctoral Fellowship NE/I021578/1 (to M.J.T.N.T.) and NERC NE/F006225/1 (to APV). A.S. was supported by SEABIG (R-154-000-648-646 and R-154-000-648-733).

Acknowledgements. The authors thank Martin Thompson for access to DNA from butterfly samples from South Africa. Rebecca Clark kindly provided details on laboratory crosses from previous studies. We also would like to thank all reviewers who have commented on various submissions of the manuscript. Sequencing libraries were constructed and sequenced at the NHM London, the Department of Biochemistry (University of Cambridge), and Genepool (University of Edinburgh).

References

- Kunte K. 2009 The diversity and evolution of Batesian mimicry in *Papilio* swallowtail butterflies. *Evolution* **63**, 2707–2716. (doi:10.1111/j.1558-5646.2009.00752.x)
- Joron M, Mallet JLB. 1998 Diversity in mimicry: paradox or paradigm? *Trends Ecol. Evol.* **13**, 461–466. (doi:10.1016/S0169-5347(98)01483-9)
- Timmermans MJTN, Thompson MJ, Collins S, Vogler AP. 2017 Independent evolution of sexual dimorphism and female-limited mimicry in swallowtail butterflies (*Papilio dardanus* and *Papilio phorcas*). *Mol. Ecol.* **26**, 1273–1284. (doi:10.1111/mec.14012)
- Poulton EB. 1924 *Papilio dardanus*. The most interesting butterfly in the world. *IJ E Afr. Uganda Nat. Hist. Soc.* **20**, 4–22.
- Davis FR. 2009 *Papilio dardanus*: The natural animal from the experimentalist's point of view. In *Descended from Darwin: insights into the history of evolutionary studies, 1900–1970* (eds J Cain, M Ruse), pp. 221–242. Philadelphia, PA: American Philosophical Society.
- Ford EB. 1936 The genetics of *Papilio dardanus* Brown (Lep). *Trans. R. Entomol. Soc. Lond.* **85**, 435–466. (doi:10.1111/j.1365-2311.1936.tb00238.x)
- Trimen R. 1869 On some remarkable mimetic analogies among African butterflies. *Trans. Linn. Soc. Lond.* **26**, 497–522. (doi:10.1111/j.1096-3642.1869.tb00538.x)
- Thompson MJ, Timmermans MJTN. 2014 Characterising the phenotypic diversity of *Papilio dardanus* wing patterns using an extensive museum collection. *PLoS ONE* **9**, e96815. (doi:10.1371/journal.pone.0096815)
- Nijhout HF. 2003 Polymorphic mimicry in *Papilio dardanus*: mosaic dominance, big effects, and origins. *Evol. Dev.* **5**, 579–592. (doi:10.1046/j.1525-142X.2003.03063.x)
- Clarke CA, Sheppard PM. 1960 The genetics of *Papilio dardanus* Brown. II. Races *dardanus*, *polytrophus*, *meseres*, and *tibullus*. *Genetics* **45**, 439–456.
- Cook SE. 1994 Mate choice in the polymorphic African swallowtail butterfly, *Papilio dardanus*: male-like females may avoid sexual harassment. *Anim. Behav.* **47**, 389–397. (doi:10.1006/anbe.1994.1053)
- Turner JRG. 1978 Why male butterflies are non-mimetic: natural selection, group selection,

- modification and sieving. *Biol. J. Linn. Soc.* **10**, 385–432. (doi:10.1111/j.1095-8312.1978.tb00023.x)
13. O'Donald P. 1969 The selective coefficients that keep modifying genes in a population. *Genetics* **62**, 435–444.
 14. Clarke CA, Sheppard PM. 1959 The genetics of *Papilio dardanus* Brown. I. Race *cenea* from South Africa. *Genetics* **44**, 1347–1358.
 15. Clarke CA, Sheppard PM. 1960 The genetics of *Papilio dardanus* Brown. III. Race *antinorii* from Abyssinia and race *meriones* from Madagascar. *Genetics* **45**, 683–698.
 16. Clark R, Brown SM, Collins SC, Jiggins CD, Heckel DG, Vogler AP. 2008 Colour pattern specification in the Mocker Swallowtail *Papilio dardanus*: the transcription factor *invected* is a candidate for the mimicry locus *H*. *Proc. R. Soc. B* **275**, 1181–1188. (doi:10.1098/rspb.2007.1762)
 17. Timmermans MJTN *et al.* 2014 Comparative genomics of the mimicry switch in *Papilio dardanus*. *Proc. R. Soc. B* **281**, 20140465. (doi:10.1098/rspb.2014.0465)
 18. Peel AD, Telford MJ, Akam M. 2006 The evolution of hexapod *engrailed*-family genes: evidence for conservation and concerted evolution. *Proc. R. Soc. B* **273**, 1733–1742. (doi:10.1098/rspb.2006.3497)
 19. Punnett RC. 1915 *Mimicry in butterflies*. London and Edinburgh, UK: Cambridge University Press.
 20. Fisher RA. 1927 On some objections to mimicry theory; statistical and genetic. *Trans. R. Entomol. Soc.* **75**, 269–274. (doi:10.1111/j.1365-2311.1927.tb00074.x)
 21. Ford EB. 1975 *Ecological genetics, fourth ed.* London, UK: Chapman and Hall.
 22. Nicholson AJ. 1927 A new theory of mimicry in insects. *Aust. J. Zool.* **5**, 10–104.
 23. Baxter SW, Johnston SE, Jiggins CD. 2009 Butterfly speciation and the distribution of gene effect sizes fixed during adaptation. *Heredity* **102**, 57–65. (doi:10.1038/hdy.2008.109)
 24. Joron M. 2003 Mimicry. In *Encyclopedia of insects* (eds RT Carde, VH Resh), pp. 714–726. New York, NY: Academic Press.
 25. Charlesworth D, Charlesworth B. 1975 Theoretical genetics of Batesian mimicry. 2. Evolution of supergenes. *J. Theor. Biol.* **55**, 305–324. (doi:10.1016/S0022-5193(75)80082-8)
 26. Turner JRG. 1977 Butterfly mimicry: the genetical evolution of an adaptation. *Evol. Biol.* **10**, 163–206.
 27. Clarke CA, Sheppard PM. 1960 Super-genes and mimicry. *Heredity* **14**, 175–185. (doi:10.1038/hdy.1960.15)
 28. Thompson MJ, Jiggins CD. 2014 Supergenes and their role in evolution. *Heredity* **113**, 1–8. (doi:10.1038/hdy.2014.20)
 29. Charlesworth D. 2016 The status of supergenes in the 21st century: recombination suppression in Batesian mimicry and sex chromosomes and other complex adaptations. *Evol. Appl.* **9**, 74–90. (doi:10.1111/eva.12291)
 30. Joron M *et al.* 2011 Chromosomal rearrangements maintain a polymorphic supergene controlling butterfly mimicry. *Nature* **477**, 203–206. (doi:10.1038/nature10341)
 31. Kunte K, Zhang W, Tenger-Trolander A, Palmer DH, Martin A, Reed RD, Mullen SP, Kronforst MR. 2014 *doublesex* is a mimicry supergene. *Nature* **507**, 229–232. (doi:10.1038/nature13112)
 32. Küpper C *et al.* 2016 A supergene determines highly divergent male reproductive morphs in the ruff. *Nat. Genet.* **48**, 79–83. (doi:10.1038/ng.3443)
 33. Nishikawa H *et al.* 2015 A genetic mechanism for female-limited Batesian mimicry in *Papilio* butterfly. *Nat. Genet.* **47**, 405–409. (doi:10.1038/ng.3241)
 34. Zhang W, Westerman E, Nitzany E, Palmer S, Kronforst MR. 2017 Tracing the origin and evolution of supergene mimicry in butterflies. *Nat. Commun.* **8**, 1. (doi:10.1038/s41467-017-01370-1)
 35. Tuttle EM *et al.* 2016 Divergence and functional degradation of a sex chromosome-like supergene. *Curr. Biol.* **26**, 344–350. (doi:10.1016/j.cub.2015.11.069)
 36. Iijima T, Kajitani R, Komata S, Lin C-P, Sota T, Itoh T, Fujiwara H. 2018 Parallel evolution of Batesian mimicry supergene in two *Papilio* butterflies, *P. polytes* and *P. memnon*. *Sci. Adv.* **4**, eaao5416. (doi:10.1126/sciadv.aao5416)
 37. Turner JRG, Sheppard PM. 1975 Absence of crossing-over in female butterflies (*Heliconius*). *J. Hered.* **34**, 265–269. (doi:10.1038/hdy.1975.29)
 38. Sadhotra M. 2016 A chromosomal investigation of three species of *Papilio* (Papilionidae:Lepidoptera). *Asian J. Anim. Sci.* **11**, 135–139. (doi:10.15740/HAS/TAJAS/11.2/135-139)
 39. Quah S, Hui JHL, Holland PWH. 2015 A burst of miRNA innovation in the early evolution of butterflies and moths. *Mol. Biol. Evol.* **32**, 1161–1174. (doi:10.1093/molbev/msv004)
 40. Jay P, Whibley A, Frézal L, de Cara MÂR, Nowell RW, Mallet J, Dasmahapatra KK, Joron M. 2018 Supergene evolution triggered by the introgression of a chromosomal inversion. *Curr. Biol.* **28**, 1839–1845. e3.
 41. Cheng Y, Brunner AL, Kremer S, DeVido SK, Stefaniuk CM, Kassis JA. 2014 Co-regulation of *invected* and *engrailed* by a complex array of regulatory sequences in *Drosophila*. *Dev. Biol.* **395**, 131–143. (doi:10.1016/j.ydbio.2014.08.021)
 42. Gustavson E, Goldsborough AS, Ali Z, Kornberg TB. 1996 The *Drosophila engrailed* and *invected* genes: partners in regulation, expression and function. *Genetics* **142**, 893–906.
 43. Nijhout HF. 1991 *The development and evolution of butterfly wing patterns*. Washington, DC: Smithsonian Institution Press.
 44. Faria R, Johannesson K, Butlin RK, Westram AM. 2019 Evolving Inversions. *Trends Ecol. Evol.* **34**, 239–248. (doi:10.1016/j.tree.2018.12.005)
 45. Kirkpatrick M, Barton N. 2006 Chromosome inversions, local adaptation and speciation. *Genetics* **173**, 419–434. (doi:10.1534/genetics.105.047985)
 46. VanKuren NW, Massardo D, Nallu S, Kronforst MR. 2019 Butterfly mimicry polymorphisms highlight phylogenetic limits of gene reuse in the evolution of diverse adaptations. *Mol. Biol. Evol.* **36**, 2842–2853. (doi:10.1093/molbev/msz194)
 47. Ando T *et al.* 2018 Repeated inversions within a pannier intron drive diversification of intraspecific colour patterns of ladybird beetles. *Nat. Commun.* **9**, 3843. (doi:10.1038/s41467-018-06116-1)
 48. Thompson MJ, Timmermans MJ, Jiggins CD, Vogler AP. 2014 The evolutionary genetics of highly divergent alleles of the mimicry locus in *Papilio dardanus*. *BMC Evol. Biol.* **14**, 140. (doi:10.1186/1471-2148-14-140)
 49. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC. 2017 GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics* **33**, 2202–2204. (doi:10.1093/bioinformatics/btx153)
 50. Kelley DR, Schatz MC, Salzberg SL. 2010 Quake: quality-aware detection and correction of sequencing errors. *Genome Biol.* **11**, R116. (doi:10.1186/gb-2010-11-11-r116)
 51. Marçais G, Kingsford C. 2011 A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* **27**, 764–770. (doi:10.1093/bioinformatics/btr011)
 52. Kajitani R *et al.* 2014 Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res.* **24**, 1384–1395. (doi:10.1101/gr.170720.113)
 53. Huang S, Kang M, Xu A. 2017 HaploMerger2: rebuilding both haploid sub-assemblies from high-heterozygosity diploid genome assembly. *Bioinformatics* **33**, 2577–2579. (doi:10.1093/bioinformatics/btx220)
 54. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015 BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* **31**, 3210–3212. (doi:10.1093/bioinformatics/btv351)
 55. Holt C, Yandell M. 2011 MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinf.* **12**, 491. (doi:10.1186/1471-2105-12-491)
 56. Stanke M, Morgenstern B. 2005 AUGUST US: a web server for gene prediction in eukaryotes that allows user-defined constraints. *Nucleic Acids Res.* **33**, W465–W467. (doi:10.1093/nar/gki458)
 57. Consortium TU. 2018 UniProt: the universal protein knowledgebase. *Nucleic Acids Res.* **46**, 2699. (doi:10.1093/nar/gky092)
 58. Jones P *et al.* 2014 InterProScan 5: genome-scale protein function classification. *Bioinformatics* **30**, 1236–1240. (doi:10.1093/bioinformatics/btu031)
 59. Catchen J, Hohenlohe PA, Bassham S, Amores A, Cresko WA. 2013 Stacks: an analysis tool set for population genomics. *Mol. Ecol.* **22**, 3124–3140. (doi:10.1111/mec.12354)
 60. Li H *et al.* 2009 The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079. (doi:10.1093/bioinformatics/btp352)
 61. McKenna A *et al.* 2010 The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research* **20**, 1297–1303. (doi:10.1101/gr.107524.110)
 62. Cingolani P, Platts A, Wang Le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. Fly (Austin). 2012 A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w¹¹¹⁸;

- iso-2; iso-3. 2012 Apr-Jun; 6(2): 80–92. PMID: 22728672. (doi:10.4161/fly.19695)
63. Margarido GRA, Souza AP, Garcia AAF. 2007 OneMap: software for genetic mapping in outcrossing species: OneMap. *Hereditas* **144**, 78–79. (doi:10.1111/j.2007.0018-0661.02000.x)
64. Li H, Durbin R. 2009 Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760. (doi:10.1093/bioinformatics/btp324)
65. Kelly JK. 1997 A test of neutrality based on interlocus associations. *Genetics* **146**, 1197–1206.
66. Garrigan D. 2013 POPBAM: tools for evolutionary analysis of short read sequence alignments. *Evol. Bioinforma.* **9**, EBO.S12751. (doi:10.4137/EBO.S12751)
67. Danecek P *et al.* 2011 The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158. (doi:10.1093/bioinformatics/btr330)