1    **Population genomics of parallel hybrid zones in the mimetic butterflies, *H. melpomene***
2    **and *H. erato***

3    Nicola J. Nadeau[1,2], Mayté Ruiz[3], Patricio Salazar[1,4], Brian Counterman[5], Jose Alejandro Medina[6],
4    Humberto Ortiz-Zuazaga[6,7], Anna Morrison[1], W. Owen McMillan[8], Chris D. Jiggins*[1], Riccardo Papa[3]

5    [1]Department of Zoology, University of Cambridge, UK; [2]Department of Animal and Plant Sciences,
6    University of Sheffield, UK; [3] Department of Biology and Center for Applied Tropical Ecology and
7    Conservation, University of Puerto Rico, Rio Piedras, San Juan, Puerto Rico 00921; [4]Centro de
8    Investigación en Biodiversidad y Cambio Climático (BioCamb), Universidad Tecnológica Indoamérica,
9    Quito, Ecuador; [5]Department of Biology, Mississippi State University, USA; [6]High Performance
10   Computing Facility, University of Puerto Rico, Puerto Rico; [7]Department of Computer Science,
11   University of Puerto Rico Rio Piedras, Puerto Rico; [8]Smithsonian Tropical Research Institute,
12   Panama.

13   *Corresponding author: c.jiggins@zoo.cam.ac.uk; Department of Zoology, University of Cambridge,
14   Downing Street, Cambridge, CB2 3EJ, UK; tel +44 1223769021; fax +44 1223336676

15   Running title: Population genomics of parallel hybrid zones

16   Key words: Hybrid zones, convergent evolution, adaptive divergence, genome-wide association
17   mapping, RAD sequencing

18

19  **ABSTRACT**

20  Hybrid zones can be valuable tools for studying evolution and identifying genomic regions
21  responsible for adaptive divergence and underlying phenotypic variation. Hybrid zones between
22  subspecies of *Heliconius* butterflies can be very narrow and are maintained by strong selection
23  acting on colour pattern. The co-mimetic species *H. erato* and *H. melpomene* have parallel hybrid
24  zones where both species undergo a change from one colour pattern form to another. We use
25  restriction associated DNA sequencing to obtain several thousand genome wide sequence markers
26  and use these to analyse patterns of population divergence across two pairs of parallel hybrid zones
27  in Peru and Ecuador. We compare two approaches for analysis of this type of data; alignment to a
28  reference genome and *de novo* assembly, and find that alignment gives the best results for species
29  both closely (*H. melpomene*) and distantly (*H. erato,* ~15% divergent) related to the reference
30  sequence. Our results confirm that the colour pattern controlling loci account for the majority of
31  divergent regions across the genome, but we also detect other divergent regions apparently
32  unlinked to colour pattern differences. We also use association mapping to identify previously
33  unmapped colour pattern loci, in particular the *Ro* locus. Finally, we identify within our sample a new
34  cryptic population of *H. timareta* in Ecuador, which occurs at relatively low altitude and is mimetic
35  with *H. melpomene malleti*.

36 **INTRODUCTION**

37 Natural hybrid zones occur where divergent forms meet, mate and hybridise. Narrow hybrid zones
38 can be maintained by strong selection that prevents mixing or favours particular forms in particular
39 areas (Barton and Hewitt 1985). Studies of hybrid zones have provided many insights into the origins
40 of diversity and the process of speciation (Mallet et al. 1990; Kawakami and Butlin 2001; Harrison
41 1993). High-throughput sequencing technologies now provide the opportunity for hybrid zones to
42 fully meet their potential as windows into the evolutionary process, by allowing us to move beyond
43 studies of neutral variation at a handful of loci, to describing the genetic basis of phenotypic
44 differences and identifying loci under selection (Crawford and Nielsen 2013; Rieseberg and Buerkle
45 2002; Gompert et al. 2012).

46 Butterflies of the genus *Heliconius* are extremely diverse in their wing colour patterns, and combine
47 extensive intra-specific diversity with perfect inter-specific similarity in wing phenotypes. The bright
48 wing colourations that identify this group of Neotropical butterflies are used as aposematic warnings
49 to predators, and are under positive frequency dependent selection, which favours common colour
50 patterns that predators learn to avoid. This strong selection also maintains narrow hybrid zones
51 between subspecies with different patterns (Benson 1972; Mallet and Barton 1989a; Kapan 2001;
52 Langham 2004). In addition, frequency dependent selection has led to Müllerian mimicry between
53 many distinct species (Müller 1879). For instance, *H. erato* and *H. melpomene* are two distantly
54 related species that diverged around 15-20 million years ago, but have converged on common colour
55 patterns across most of the Neotropics, with parallel hybrid zones between subspecies of both
56 species (Figure 1). Current evidence suggests that the convergent colour patterns in these species
57 have most likely evolved independently (Supple et al 2013, Hines et al. 2011). It has been suggested
58 that *H. erato* is more ancient and *H. melpomene* diversified more recently to mimic the *H. erato*
59 forms (Brower 1996; Flanagan et al. 2004; Quek et al. 2010). Nevertheless, it appears that the same
60 handful of genetic loci are responsible for producing most of the colour pattern variation observed in
61 both species (Joron et al. 2006; Baxter et al. 2008; Martin et al. 2012; Reed et al. 2011). This pattern
62 of parallel adaptive evolution, with multiple instances of both species convergently evolving highly
63 similar phenotypes, makes this an excellent system in which to address questions about the
64 predictability of the evolutionary process and the extent to which particular genes are re-used when
65 evolving the same phenotypes (Nadeau and Jiggins 2010; Papa et al. 2008).

66 In this study, we use high resolution genome scans to investigate patterns of divergence across two
67 pairs of parallel hybrid zones, in Peru and Ecuador. These occur between subspecies with different
68 wing colour patterns in both *H. erato* and *H. melpomene* (Figure 1). In both regions the clines in
69 colour pattern alleles between species are highly coincident (Mallet et al. 1990; Salazar 2012). The
70 two hybrid zones in Peru have been the focus of several previous studies, while those in Ecuador
71 have been less well studied. In Peru, strong natural selection has been shown to maintain colour
72 pattern differences (Mallet and Barton 1989a) and the colour pattern controlling loci show increased
73 divergence compared to other loci in the genome (Baxter et al. 2010; Counterman et al. 2010;
74 Nadeau et al. 2012; Supple et al. 2013; S. H. Martin et al. 2013). However, we are lacking a clear
75 picture of exactly how many loci are divergent between subspecies, and the extent to which the
76 genomic architecture of divergence is the same between mimetic species.

77    Extensive genetic mapping using experimental crosses between different colour pattern forms has
78    identified the chromosomal regions responsible for colour pattern variation in these species
79    (Sheppard et al. 1985; Baxter et al. 2008; Joron et al. 2006; Papa et al. 2013). Three major clusters of
80    loci control most of the colour pattern variation observed in both species. The tightly linked $B$ and $D$
81    loci in $H.$ $melpomene$ control the red forewing band and the red/orange hindwing rays and proximal
82    "dennis" patches on both wings respectively. These loci are homologous to the $D$ locus in $H.$ $erato$
83    (Baxter et al. 2008) and appear to be cis regulatory elements of the $optix$ gene on chromosome 18
84    (Reed et al. 2011; Supple et al. 2013). The $Ac$ and $Sd$ loci, in $H.$ $melpomene$ and $H.$ $erato$ respectively,
85    control the shape of the forewing band via regulation of the $WntA$ gene on chromosome 10 (Martin
86    et al. 2012). The presence of most yellow and white elements on the wing is largely controlled by
87    three tightly linked loci, $Yb,$ $Sb$ and $N,$ on chromosome 15 in $H.$ $melpomene$ (Ferguson et al. 2010),
88    which are homologous to the $Cr$ locus in $H.$ $erato$ (Joron et al. 2006). Quantitative trait locus (QTL)
89    mapping has identified other loci of minor effect, including at least 7 additional QTL in $H.$ $erato$ (Papa
90    et al. 2013), and QTL in $H.$ $melpomene$ on chromosomes 2, 7 and 13 that affect forewing band shape
91    (Baxter et al. 2008). In some cases mapping studies have been followed up by population genetic
92    studies of the mapped intervals across natural hybrid zones where many generations of crossing and
93    backcrossing have led to narrow regions of association, permitting fine scale mapping (Baxter et al.
94    2010; Counterman et al. 2010; Nadeau et al. 2012; Supple et al. 2013). High-throughput sequencing
95    technologies now provide the feasibility to generate a suitably high density of genomic markers to
96    identify the narrow QTL present in these hybrid zones without the necessity of performing
97    controlled laboratory crosses (Crawford and Nielsen 2013). In this study we aim to test this
98    approach, using this system where we already know some of the loci responsible for phenotypic
99    differences.

100   The hybrid zones that we focus on all occur across altitudinal gradients (Figure 2A). Therefore, it is
101   possible that traits other than colour pattern may be differentiated across the hybrid zones driven
102   by altitudinal selection, for example related to temperature or changes in larval host plants. Such
103   selection on additional regions of the genome could help to stabilise the hybrid zone, which, if due
104   to frequency dependent selection acting on colour pattern alone, could be unstable (Bierne et al.
105   2011; Barton and Hewitt 1985; Mallet and Barton 1989b; Mallet 2010). Therefore an important
106   question that we will address is whether there are divergent regions of the genome that are not
107   controlling colour pattern, as these could be candidates for loci controlling other aspects of
108   ecological adaptation.

109   In this study, we use restriction associated DNA (RAD) sequencing (Baird et al. 2008) to determine,
110   for the first time:

111       1)  if association mapping in these hybrid zones can identify known and novel loci underlying
112           phenotypic variation

113       2)  how much of the genome is differentiated and under divergent selection between
114           subspecies

115       3)  how much of this differentiation is due to loci controlling colour pattern variation

116       4)  if the same regions are divergent between co-mimetic species

117    We also investigate the advantages and limitations of alignment and assembly methods when only a
118    single reference genome is available. We compare two widely used approaches: *de novo* assembly of
119    restriction associated reads using the program Stacks (Catchen et al. 2011) versus alignment of
120    paired end reads to the reference *H. melpomene* genome.

121

## RESULTS

### Summary of the data and comparison of alignment and assembly techniques

124    We sequenced a total of 129 individuals of *H. erato* and *H. melpomene* from the four hybrid zones in
125    Peru and Ecuador, and including a small number of additional individuals from across the range of *H.*
126    *erato*. Using restriction associated DNA sequencing (RAD-seq), we obtained a total of 1,496M 150
127    base paired-end reads from the hybrid zone individuals, and an additional 115M 100 base paired-
128    end reads from the other *H. erato* populations and outgroups. We also include in our analyses data
129    from additional *H. melpomene* populations and outgroups already published in a previous study
130    (Nadeau et al. 2013).

131    Our reference genome for *H. melpomene* is highly divergent from *H. erato*. Nonetheless, for both
132    species, alignment of reads to the reference yielded more usable data when compared to *de novo*
133    assembly. The latter produced more bases in assembled contigs (Table 1), but only ~2% of contigs
134    assembled in the *H. erato* populations were present in more than 10 individuals, with the figure
135    being approximately 7% in *H. melpomene*. By comparison when the same data (plus the paired-end
136    reads) were aligned to our reference sequence, approximately 38% of aligned bases were found in
137    more than 10 individuals in *H. erato* and >50% in *H. melpomene*. We hypothesised that high levels of
138    within population variation led to homologous reads being separated into distinct contigs in the *de*
139    *novo* assembly. We could confirm that this was the case for one region of the *H. erato* genome for
140    which a high quality reference sequence is available (Supple 2013).  Across 960kb at the *D* colour
141    pattern locus, we observed that in regions with high sequence divergence between subspecies, our
142    RAD-seq reads were assembled into separate contigs. Overall, we also found a higher frequency of
143    single nucleotide polymorphisms (SNPs) in the reference alignments than the *de novo* assemblies
144    (Table 1). These SNPs were defined as sites that were polymorphic within the sampled populations
145    and so are not inflated by fixed differences from the reference genome.

146    As expected, given that *H. erato* is ~15% divergent from *H. melpomene* in the aligned data, fewer *H.*
147    *erato* reads aligned to the *H. melpomene* genome as compared to those from *H. melpomene*, leading
148    to fewer confidently called bases. Nevertheless, the use of the reference *H. melpomene* genome for
149    aligning the *H. erato* reads resulted in more bases being called across multiple individuals and
150    around 10x more SNPs identified when compared to the *de novo* assembly approach. In addition,
151    the gaps between aligned RAD loci were similar across both species (Table 1), indicating that the
152    reduced number of bases is not due to fewer RAD loci aligning but to fewer confidently called bases
153    at each RAD locus. The power to detect loci under selection or responsible for phenotypic variation
154    should therefore be similar in both species. In summary, it seems that the aligned data should give
155    the most power to detect divergent regions and phenotypic associations for both species. However,
156    we performed outlier and association analyses using the output of both approaches for comparison.

157  It is also possible that the *de novo* assembly might detect divergent regions important in adaptation
158  that could not be aligned to the *H. melpomene* reference.

159  **Phylogenetics and population structure**

160  Using the reference aligned sequence data, we constructed maximum likelihood phylogenies for the
161  *H. melpomene* and *H. erato* clades, including individuals from additional populations and outgroup
162  taxa (Figure 1). This revealed remarkably similar patterns of divergence between co-occurring, co-
163  mimetic subspecies in both groups. Population divergence in *H. erato* is thought to be deeper than
164  that in *H. melpomene* (Flanagan et al. 2004 but see Cuthill and Charleston 2012), but this was not
165  evident in our tree as branch lengths were very similar between the two species. However this may
166  be due to the poorer alignments for *H. erato*, with the *H. erato* tree based only on about a third as
167  many sites as that for *H. melpomene*. These sites in *H. erato* are likely to be more conserved,
168  resulting in some compression of the tree topology.

169  The most striking finding from the phylogenetic reconstruction was that eight of the presumed *H.*
170  *melpomene* individuals from Ecuador were strongly supported as clustering within the *H. timareta*
171  clade (Figure 1). All of these individuals had a *H. melpomene malleti*-like phenotype with the
172  exception of one individual which had been characterised as a possible hybrid due to a large and
173  rounded yellow forewing band, but was otherwise *H. m. malleti*-like. This finding was surprising
174  because while populations of *H. timareta* mimetic with *H. m. malleti* have previously been described
175  in Colombia (Giraldo et al. 2008) and Northern Peru (Lamas 1997), they are all found in highland
176  areas above ~1000m. Similar populations are not known from lowland sites anywhere in the range.
177  To compare our individuals to these and other populations, we also directly sequenced part of the
178  mitochondrial COI gene that overlaps with the region sequenced in previous studies (Giraldo et al.
179  2008; Mérot et al. 2013). Our phylogeny based on these sequences also robustly supported these
180  eight individuals as being *H. timareta* and placed them closer to the highland *H. timareta timareta* in
181  Ecuador than to *H. timareta florencia* in Colombia that they resemble phenotypically (SI figure 1).

182  The newly identified *H. timareta* subspecies was also clearly evident in a principle components
183  analysis (PCA) of the combined *H. melpomene, H. timareta* and *H. cydno* data. The first principle
184  component separated the Peruvian *H. melpomene* from *H. timareta* and *H. cydno* (which were very
185  similar on this axis, SI figure 2). The grouping of the Ecuadorian samples was consistent with the
186  phylogeny, with the same eight individuals clustering with *H. timareta*. No individuals were
187  intermediate between *H. melpomene* and *H. timareta*, indicating that the level of genetic isolation
188  between the two species is similar to elsewhere in their range. This was also confirmed by a
189  STRUCTURE analysis of the Ecuador *"H. melpomene"* population, where a model with two populations
190  had the best fit to the data (posterior probability=1). Under this model, which allowed for admixture
191  between populations, the *H. timareta* individuals all had 100% of their allelic contribution from
192  population 1, while for *H. melpomene* the maximum contribution of population 1 to any individual's
193  genotype was 1.8% (SI Table 1). In summary, we can conclude that these are distinct species with
194  little gene flow between them.

195  We conducted further analyses of the genetic structure of each of the hybrid zone populations,
196  excluding the *H. timareta* individuals, again using the reference aligned data. Overall, these results
197  suggest only very low genetic differentiation between any of the parapatric subspecies. STRUCTURE
198  analyses of each population generally showed very little structure and strongest support for only a

199     single population cluster being present. The only exception was the Peruvian *H. melpomene*, where
200     two population clusters gave the highest posterior probability (p=1). However, these clusters did not
201     correspond to the two subspecies. The genetic diversity was partitioned such that most individuals
202     were admixed with about a quarter of their allelic variation from population 2, except for two
203     "hybrid" individuals that had pure population 2 genotypes and two other individuals (one "hybrid"
204     and one *aglaope*) that had almost pure population 1 genotypes (Figure 2B). PCA revealed very
205     similar patterns, with small groups of hybrid phenotype individuals giving the clearest clusters, which
206     in most cases were also identified by STRUCTURE (with K=2, Figure 2C). Three of the populations did
207     reveal some separation of the subspecies at one of the first two principle components, but with a
208     gradual change from one genomic "type" to another. The *H. melpomene* subspecies in Ecuador were
209     separated by PC1, which explained 10% of the variation in this population. The two *H. erato*
210     populations both showed some separation by subspecies at PC2, which explained 5.7% and 6.7% of
211     the variation in Peru and Ecuador respectively. We found very similar results with PCA on the *de*
212     *novo* assembled data (SI figure 3), suggesting that the underlying genetic signal in both data sets is
213     very similar. The lack of strong differentiation between subspecies was also supported by the $F_{ST}$
214     distributions (calculated by BayeScan), which gave very low $F_{ST}$ values between subspecies at over
215     99% of the genome, with only a small percentage of SNPs showing high levels of differentiation
216     (Figure 2D, SI figure 3).

217     **Association mapping of loci responsible for phenotypic variation**

218     We performed association mapping to identify genetic regions responsible for the phenotypic
219     variation that segregates across each of the hybrid zones. In general, the expected associations were
220     found at the three major loci known to control colour pattern variation on chromosomes 10, 15 and
221     18 (Figures 3A/D, 4A/D and Table 2). The majority of SNPs showing significant phenotypic
222     associations fell within or tightly linked to these loci in all populations except in Peruvian *H. erato*,
223     where only 26% were tightly linked to the known loci.

224     Independent analyses were performed on both the reference alignments and *de novo* assemblies of
225     the data. In all populations, more associated SNPs were identified in the alignments than in the *de*
226     *novo* assemblies (Table 1, figure 5). We used blastn (Altschul et al. 1990) to place *de novo* contigs
227     containing associated SNPs onto the *H. melpomene* genome, and most could be confidently assigned
228     to a unique locus. There was almost no overlap in the particular SNPs detected in the differently
229     assembled and aligned data sets (Figure 5), although in many cases the SNPs detected were in
230     similar regions (Figures 3 and 4). There was some evidence for a higher false positive rate in the *de*
231     *novo* data, as a higher proportion of associations were found scattered across the genome, away
232     from the known colour pattern loci.

233     *Red colour pattern elements and the B and D loci*

234     Consistent with previous mapping studies based on crosses, red colour pattern variation in almost all
235     populations was mapped to scaffold HE670865 on chromosome 18 (Baxter et al. 2008; Counterman
236     et al. 2010; The *Heliconius* Genome Consortium 2012). The only exception was *H. melpomene* in
237     Ecuador, where SNPs in this region did not reach significance (Figure 4A). This is likely to be due to
238     the reduced sample size of this population (22) after removing the *H. timareta* individuals. Colour
239     patterns were scored both as independent elements and also using known patterns of segregation
240     to score the predicted genotype at the *B/D* locus (individuals with both red forewing and rays are

241    predicted to be heterozygotes, as both traits are dominant). The latter scoring generally gave
242    stronger associations (Figures 3 and 4), although both methods gave some significant associations
243    for at least one of the traits. In all populations the strongest associations in this region were over
244    60kb downstream of the *optix* gene that controls red colour pattern (Reed et al. 2011), but
245    overlapped with the region identified in previous analyses as likely containing the functional
246    regulatory variation (Nadeau et al. 2012; Supple et al. 2013)(Table 2, SI figure 4).

247    In the Peruvian *H. melpomene* population we also found weaker but significant associations with the
248    *B/D* genotype on other linked chromosome 18 scaffolds, particularly HE671488. However, this
249    scaffold was also associated with differences in altitude, which were stronger than the associations
250    with colour (Figure 3A, Table 3, SI figure 4). This could suggest that this *B/D* linked region is
251    responsible for ecological adaptation, although colour and altitude are strongly correlated so we do
252    not have much power in the data set to separate the two. Both the Peruvian and Ecuadorian *H.*
253    *melpomene* populations had a SNP at position 97 on an unmapped scaffold, HE670458, that was
254    highly associated with rays (Table 3). This scaffold appears to consist largely of repetitive elements
255    (BLAST hits match many other regions of the *H. melpomene* genome), suggesting that there may be
256    a copy of a repetitive element that is associated with the presence of rays in both populations. All
257    rayed individuals were heterozygous and all non-rayed individuals were homozygous at this SNP in
258    both *H. melpomene* populations, which would be consistent with the presence of a unique copy of a
259    repetitive element linked to the rayed allele. The existence of such a repetitive element is consistent
260    with previous findings that repetitive elements are present in the region of highest divergence at the
261    *B/D* locus (Nadeau et al. 2012; Papa et al. 2008).

262    Surprisingly, in the Peruvian *H. erato* population the strongest associations with red colour pattern
263    elements were not on chromosome 18, but at two scaffolds (HE670771 and HE670235) on
264    chromosome 2 (Figure 3D, Table 3, SI figure 4). In addition, two SNPs significantly associated with
265    rays and *D* genotype in the *de novo* assembled data of this population could not be confidently
266    assigned to a position in the genome: One with a fairly weak top BLAST hit to the end of
267    chromosome 10 and multiple equally strong hits elsewhere in the genome; the other had a top
268    BLAST hit to an unmapped scaffold but additional good hits to two other unmapped scaffolds and
269    chromosomes 2 and 19.

270    *Yellow colour pattern elements and the Yb, N and Cr loci*

271    In the Peruvian *H. melpomene* population the presence of the yellow hindwing bar and yellow in the
272    forewing band both mapped to chromosome 15, with positions that were consistent with previous
273    work on the *Yb* locus (Ferguson et al. 2010; Nadeau et al. 2012; The *Heliconius* Genome Consortium
274    2012)(Table 2, SI figure 4). Associations with altitude were also found at these associated SNPs but
275    were weaker than the association with colour and so may simply be due to correlations between the
276    altitude of the sampling site and colour pattern.

277    In the Peruvian *H. erato* population we did not recover the expected associations with the yellow
278    hindwing bar, which is known to be controlled by the *Cr* locus on chromosome 15 (Joron et al. 2006;
279    Counterman et al. 2010). Instead, the strongest association with this phenotype in the reference
280    aligned data was found on chromosome 17 (Table 3). Moreover, we also identified significant
281    associations with the yellow hindwing bar on chromosome 10 in both the reference aligned and *de*
282    *novo* assembled data. In particular, one SNP on chromosome 10 was detected in both data sets on

283    scaffold HE668478. These associations can be explained by the presence of *Sd* on this scaffold
284    (Martin et al. 2012) (Table 2), which is known to influence the expression of the yellow hindwing bar,
285    particularly in individuals that are heterozygous at the *Cr* locus (Mallet 1989). The *Cr* locus is not
286    thought to control any aspects of phenotypic variation in Ecuadorian *H. erato* (Salazar 2012) and
287    consistent with this expectation we did not detect any phenotypic associations in this region.

288    Significant associations with yellow in the forewing band were present in the *D* region in Ecuadorian
289    *H. erato* (Figure 4D), consistent with the fact that the *D* locus controls both yellow and red
290    colouration in the forewing band in *H. erato* (Sheppard et al. 1985, Salazar 2012, Papa et al. 2013).
291    No significant associations were found with the presence of yellow colour in the forewing band in
292    either Peruvian *H. erato* or Ecuadorian *H. melpomene*.

293    We scored the Ecuadorian *H. melpomene* individuals for their predicted genotype at the *N* locus
294    (Figure 4C), which is known to control the amount and location of red in the forewing band and the
295    length of the orange hindwing "dennis" bar (Salazar 2012). Despite the epistatic interaction between
296    the *B/D* and the *N* loci, we could still score them independently.  Associations with *N* were found to
297    overlap the *Yb* region (Ferguson et al. 2010; Nadeau et al. 2012) on chromosome 15, scaffold
298    HE667780 (Table 2, SI figure 4), in both the reference aligned and *de novo* assembled data (Figure
299    4A). These are the first genetic mapping results for the *N* locus and confirm previous studies' findings
300    based on segregation patterns that it is tightly linked to *Yb* (Sheppard et al. 1985; Mallet 1989).

301    *Forewing band shape and the Ac, Sd and Ro loci*

302    In Peruvian *H. melpomene,* the strongest associations with forewing band shape (cell spot 8 and cell
303    spot 11) were on chromosome 18, within the *B/D* region (Figure 3A). This suggests that the *B/D* locus
304    controls the shape as well as the colour of the forewing band in Peruvian *H. melpomene*. However,
305    we did also find a cluster of 8 SNPs associated with band shape on an unmapped scaffold, HE671554.
306    New mapping analyses suggest that this scaffold is on chromosome 20 (J Davey, Pers. Comm.) and
307    therefore not linked to any previously described colour pattern controlling loci (Table 3).

308    In Peruvian *H. erato*, the SNP in the *Sd* region that was associated with the yellow hindwing bar also
309    showed the expected association with forewing band shape in the *de novo* assembly but not the
310    reference alignment. This SNP was just 5kb upstream of the *WntA* gene (Table 2, Figure 3D, SI figure
311    4). Associations with forewing band shape were also found on chromosome 2 in this and the
312    Ecuadorian *H. erato* populations (Figure 3D, Figure 4D, SI figure 4), in similar regions to those
313    associated with red colour in Peruvian *H. erato* (Table 3).

314    In both species from the Ecuadorian hybrid zone, we found SNPs associated with forewing band
315    shape (cell spot 7/8/11) within introns of the *WntA* gene (Figure 4, Table 2). In Ecuadorian *H. erato,*
316    we also found two tightly linked SNPs on chromosome 13, scaffold HE670984, and three tightly
317    linked SNPs on an unmapped scaffold, HE669551, which were associated with forewing band shape
318    and also rounding of the band (Figure 4D). Rounding of the distal edge of the band in this population
319    has previously been described as being under the control of the unmapped *Ro* locus (Sheppard et al.
320    1985; Salazar 2012). More recent mapping analysis suggests that HE669551 is within 1cM of
321    HE670984 on chromosome 13 (J Davey, Pers. Comm.) so these associations are most likely due to a
322    single locus on this chromosome. This is therefore the most likely genomic location of the *Ro* locus,
323    representing the first time this locus has been mapped.

324     **$F_{ST}$ outlier detection**

325     Outlier detection provides an alternative method for identification of loci under selection that does
326     not depend on phenotypic association. BayeScan detected less than 0.06% of SNPs as outliers in
327     each of the analyses (Table 1).  In the *de novo* assembly, Peruvian hybrid zones showed a greater
328     percentage of SNPs as outliers in both *H. erato* (0.059%) and *H. melpomene* (0.040%), with no
329     outliers detected in Ecuadorian *H. erato* and only five in Ecuadorian *H. melpomene* (0.012%).  The
330     overall proportion of SNPs detected in the reference aligned data was similar. However, unlike the
331     *de novo* assemblies, the proportions of outliers found in the reference alignments were more similar
332     within each species than within each locality.  Reference aligned data from *H. melpomene* contained
333     approximately 0.025% outlier SNPs in both Peru and Ecuador, while reference aligned data from *H.*
334     *erato* had 0.005% outliers in Peru and 0.017% outliers in Ecuador (Table 1). This would be consistent
335     with some of the most rapidly diverging regions being lost in *H. erato* when aligned against the
336     reference *H. melpomene* genome.

337     As suggested by results from the *de novo* assemblies, there do appear to be differences in
338     population structure between the geographic regions that are consistent across both species. This is
339     also reflected in the $F_{ST}$ distributions (from both alignment and assembly approaches), with both *H.*
340     *erato* and *H. melpomene* having higher mean and background levels of $F_{ST}$ in Ecuador as compared to
341     Peru (Table 1, Figure 2, SI figure 3). However, within both regions *H. melpomene* has a lower mean
342     $F_{ST}$ than *H. erato*, which would be consistent with higher dispersal distances in *H. melpomene,* as
343     previously suggested (Mallet et al. 1990). Similar outlier regions were detected by both the
344     alignment and assembly approaches (Figures 3 and 4, B and E), although only Peruvian *H.*
345     *melpomene* gave a good overlap in the specific SNPs detected (Figure 5). Some of the outlier contigs
346     detected in Peruvian *H. erato* could not be positioned on the *H. melpomene* genome with
347     confidence. In particular, one contig containing 4 outlier SNPs, had the highest BLAST homology to
348     an unmapped scaffold (score = 87.7 bits) but similar hits (~86 bits) to chromosomes 1 and 4. In
349     addition, two of the contigs containing 3 outlier SNPs appeared to be linked to the *Sd* and *D* loci but
350     were assigned with very low confidence (BLAST scores < 40 bits).

351     Overall, there was considerable overlap between the genomic regions containing outlier SNPs and
352     those showing phenotypic associations (Figures 3 and 4), and to some extent in the specific SNPs,
353     with the majority of phenotypically associated SNPs also being outliers (Figure 5). The exception to
354     this general trend was the Peruvian *H. erato* population where a large proportion of the phenotypic
355     associations were not strongly divergent between subspecies. In all populations over 50% of outlier
356     SNPs were within 1Mb of a known colour pattern locus (including the newly identified *Ro* region;
357     excluding these, 37.5% of outliers in Ecuadorian *H. erato* were within 1Mb of the *D* and *Sd* loci). The
358     strongest outliers on chromosome 10 in the Ecuadorian populations and Peruvian *H. erato* were
359     within intons of the *WntA* gene and the strongest outliers on the *B/D* scaffold were all 3' of the *optix*
360     gene (Table 2, SI figure 4).

361     In both *H. melpomene* populations there was a second strongly divergent region on chromosome 18
362     about 2Mb from the *B/D* region, which was not divergent in either of the *H. erato* populations
363     (Figure 3B, SI figure 4). This is the same region on scaffold HE671488 that showed associations with
364     colour pattern and altitude in the Peruvian *H. melpomene* population (Table 3). In the Peruvian *H.*
365     *melpomene* population, we detected two clusters of outlier divergent SNPs on chromosome 6, which

366     do not appear to be associated with colour pattern (Figure 3B, Table 3, SI figure 4). Outliers were
367     also detected on chromosome 2 in both *H. erato* populations, some of which were in similar regions
368     to those detected in the association mapping (Table 3, SI figure 4).

369

370     **DISCUSSION**

371     It has long been recognised that convergent and parallel evolution provides a natural experimental
372     system in which to study the predictability of adaptation (Stewart et al. 1987; Wood et al. 2005). This
373     approach has come to the fore with the recent integration of molecular and phenotypic studies of
374     adaptive traits (Stinchcombe and Hoekstra 2007; Nadeau and Jiggins 2010). Here, we have studied
375     parallel divergent clines in two co-mimic species of butterflies, using RAD sequencing to generate an
376     extensive data set covering 1-5% of the entire genome. Previous genomic studies of these species
377     have sampled only a few individuals of divergent wing pattern races (Nadeau et al. 2012; The
378     *Heliconius* Genome Consortium. 2012; Nadeau et al. 2013; Martin et al. 2013; Supple et al. 2013;
379     Kronforst et al. 2013), while previous hybrid zone studies have yet to integrate next-generation
380     sequencing approaches (Mallet and Barton 1989a; Salazar 2012; Baxter et al. 2010; Counterman et
381     al. 2010). Here we have shown that association mapping in the hybrid zones can be used to find
382     known loci, but also identify previously unmapped loci, such as *Ro* in *H. erato* and *N* in *H.*
383     *melpomene*. Moreover, we conducted the first genome-wide scan for divergent loci and identify
384     some that are not wing colour pattern related and so may have a role in other aspects of ecological
385     divergence. With these data, we also identify a cryptic population of *H. timareta* in Ecuador and
386     reveal parallel patterns of divergence between co-mimetic species.

387     **Comparison of de novo assembly and reference alignment of RAD data**

388     Genome-wide association studies (GWAS) are now common in studies of admixed human
389     populations (Visscher et al. 2012). Their use outside of model organisms has mostly been hampered
390     by lack of reference genomes or methods for typing sufficient numbers of markers. However, these
391     limitations are rapidly being eroded as the cost of sequencing decreases and more reference
392     genomes become available. Furthermore, we have shown that alignment of reads to a fairly
393     distantly related reference genome (~15% divergent) can generate meaningful results. In the
394     absence of a reference genome, *de novo* assembly also detects the same loci, but with somewhat
395     reduced efficacy.

396     Alignment of sequence reads to the reference genome produced data for more sites, even in the
397     more distantly related species, *H. erato*. One drawback of the STACKS pipeline that we used for *de*
398     *novo* assembly of the reads is that it does not assemble and call sequence variants in the paired end
399     reads. Hence the available sequence for analysis is almost double in the reference alignments as
400     compared to the de novo assembly. However, it also seems that data was lost in the to *de novo*
401     assembly due to divergent alleles not being assembled together. This may have had a larger
402     influence on the *H. erato* assemblies as this species harbours greater genetic diversity than *H.*
403     *melpomene* (Hines et al. 2011) and so explain why a much lower proportion of the *de novo*
404     assembled contigs were present across multiple individuals in *H. erato* (Table 1). We also found a
405     higher proportion of variable sites in the reference alignments as compared to the *de novo*
406     assemblies. This may again be due to poor assembly of the *de novo* contigs, but it could also

407    represent genetic variability contained in the paired-end reads. It is possible that paired end reads
408    might be located in more variable regions, particularly if restriction-site associated reads were
409    biased towards more conserved regions (The *Heliconius* Genome Consortium 2012).

410    The larger number of SNPs in the reference alignments resulted in larger numbers of outlier and
411    associated SNPs being detected, most of which cluster in the expected genomic regions. Moreover,
412    there appears to be a higher false positive rate in the association mapping using the *de novo*
413    assembled data. The most likely explanation for this result is that the smaller number of SNPs
414    generated from the *de novo* assembly gave less power to correct for underlying population
415    structure. Nevertheless, many of the expected associations and outlier regions were detected in the
416    *de novo* assembled data. The results from assembly and alignment approaches are more concordant
417    in *H. melpomene* than *H. erato* particularly at the level of individual SNPs (Figure 5). This is very likely
418    due to the fact that the *H. melpomene* reference genome was used to generate the sequence
419    alignments in both species. In addition, the lower within-population diversity in *H. melpomene*, may
420    also have led to improved *de novo* assemblies in this species.

421    Overall, our results suggest that detection of loci underlying adaptive change is likely to be more
422    effective where reads can be mapped to a reference genome. The *de novo* approach could, and no
423    doubt will, be improved by developing methods that allow paired end reads to be incorporated into
424    the SNP typing pipeline (Baxter et al. 2011). This would not only allow a higher density of SNPs to be
425    detected but could also improve alignment of divergent alleles. In the mean time, one approach that
426    has been used in other studies is to *de novo* assemble RAD-seq reads to generate a consensus
427    reference that the reads are then mapped back to for SNP calling (Keller et al. 2013).

428    **Association mapping across hybrid zones is a rapid way of detecting loci underlying phenotypic**
429    **differences**

430    We have successfully used association mapping in hybrid zone individuals to identify virtually all of
431    the genomic regions known to control colour pattern in these populations (Reed et. al. 2011, Martin
432    et al. 2012; Nadeau et al. 2012; Supple et al. 2013). It has commonly been supposed that large
433    sample sizes will be necessary in order identify genes in wild populations. Here we have confirmed
434    recent theoretical predictions from simulated data (Crawford and Nielsen 2013), that for large effect
435    adaptive loci, even small sample sizes can be highly effective in identification of narrow genomic
436    regions underlying adaptive traits. We also confirm the prediction that in populations with low
437    background levels of divergence both divergence outlier and association mapping approaches are
438    effective in detecting regions under divergent selection. In our study, association mapping has the
439    added benefit of identifying the phenotypic effects of the selected loci. One anticipated pitfall of this
440    method was that many of the phenotypes covary across the hybrid zone. However, it appears that
441    with just 10 individuals with admixed phenotypes we can disassociate most of the variation and thus
442    find distinct genetic associations for known loci. This therefore gives us some confidence that the
443    novel associations that we have detected are real and not due to covariation with other phenotypes.

444    In Ecuador we intentionally sampled from sites at the edges of the hybrid zone where both pure and
445    hybrid individuals were present, as we anticipated that individuals from these sites would have the
446    highest levels of admixture between selected alleles. This may explain the clearer patterns observed
447    in *H. erato* in Ecuador as compared to Peru (Figures 3D and 4D). In Peruvian *H. erato* we also find
448    several genomic regions showing phenotypic associations that are not divergence outliers, which

449    may suggest that these are likely to be false positives. However, the less clear signal in Peruvian *H.*
450    *erato* could also be due to the reduced sample size in this population (27 individuals). Certainly, the
451    reduced number of Ecuadorian *H. melpomene* individuals (22) seems to have greatly reduced the
452    power of the association mapping to detect significant associations (Figure 4B). The lack of any signal
453    at the *Cr* locus in Peruvian *H. erato* is still surprising, and may be because the *Cr* associated region is
454    very narrow. Linkage disequilibrium breaks down rapidly in *H. erato* (Counterman et al. 2010) and
455    even though there are >700 SNPs in *H. erato* within the region that contains hindwing bar associated
456    SNPs in *H. melpomene*, this may not be enough to identify the loci responsible for this phenotypic
457    variation. Nevertheless, contrary to previous suggestions (Kronforst et al. 2013), the density of RAD
458    markers we have obtained was sufficient to identify many other narrow divergent genomic regions.

459    Although we have clearly demonstrated the utility of this approach for association mapping, it
460    should be noted that scoring of some phenotypes was informed by previous crossing experiments.
461    For example, the *N* locus in Ecuadorian *H. melpomene* was scored taking into account the genetic
462    background at *B/D* (Salazar 2012), and the scoring of the predicted genotype at the *B/D* locus
463    yielded stronger associations than scoring of individual colour pattern elements. Nonetheless,
464    scoring based purely on phenotypic variation did successfully identify colour pattern loci in several
465    cases (eg. *Ro, Ac/Sd* and *Yb*). Overall, the prospects for mapping individual phenotypic components
466    and identifying epistatic relationships without prior knowledge are considerable, especially with
467    larger sample sizes.

468    The possibility of using hybrid zones for association mapping has long been recognised (Kocher and
469    Sage 1986) but few studies have successfully applied this technique. Studies in younger hybrid zones,
470    for example *Helianthus* sunflowers, have found that linkage disequilibrium between unlinked
471    genomic regions in early generation hybrids can produce spurious associations (Rieseberg and
472    Buerkle 2002). *Heliconius* hybrid zones seem ideal in this regard because they appear to be fairly
473    ancient and close to linkage equilibrium. However, it seems likely that many other suitable systems
474    do exist for this type of approach (Lexer et al. 2006; Crawford and Nielsen 2013). An additional
475    benefit of the *Heliconius* system is that much of the phenotypic variation is controlled by major
476    effect loci, which can be detected with small sample sizes. Much larger sample sizes would be
477    required in order to detect minor effect loci (Beavis 1997). Nevertheless, adaptive phenotypes
478    appear to frequently arise from major effect loci (Orr 2005; Nadeau and Jiggins 2010) and so there
479    may be many situations in which the use of very large sample sizes is unnecessary. In addition, by
480    incorporating methods that use a probabilistic framework to infer allele frequencies in low coverage
481    sequencing data (Gompert and Buerkle 2011) it should be feasible to sequence large enough
482    samples for analysis of quantitative traits.

483    **Identification of a novel colour pattern locus**

484    Our association mapping results have robustly identified the *H. erato Ro* locus, that controls the
485    shape of the distal edge of the forewing band, as being on chromosome 13 near gene HE669551.
486    This gene has a predicted Gene Ontology (GO) molecular function of microtubule binding and is
487    similar to other insect Radial Spoke Head 3 proteins, which are components of the cilia (Avidor-Reiss
488    et al. 2004). It is therefore not an obvious candidate for control of colour pattern, so may simply be
489    linked to the causative site. Our results are contrary to the suggestion of a recently published QTL
490    study that *Ro* may be linked to *Sd* (Papa et al. 2013). However, that study also identified a major

491    unlinked QTL for forewing band shape, that could not be assigned to a *H. melpomene* chromosome
492    and so may be homologous to the locus we detected here. Furthermore, a QTL for several aspects of
493    forewing band shape and size, including the shape of the distal edge, has previously been identified
494    in *H. melpomene* on chromosome 13 (Baxter et al. 2008). This was located to a fairly broad region
495    but its positioning is consistent with our results for the *Ro* locus in *H. erato*. It therefore seems likely
496    that we have identified a new wing patterning locus that is homologous in *H. melpomene* and *H.*
497    *erato*.

498    **Ecological selection across the hybrid zones**

499    Our results support previous assertions that selection acting on colour pattern is the most important
500    factor in maintaining these hybrid zones (Mallet and Barton 1989a; Baxter et al. 2010; Counterman
501    et al. 2010; Nadeau et al. 2012; Supple et al. 2013). The most divergent genomic regions correspond
502    to colour pattern controlling loci and at least half of all divergence outliers are in these regions.
503    Nevertheless, some divergent regions do not seem to correspond to colour pattern loci, and could
504    be candidates for adaptation to other ecological factors. The best candidates appear to be the
505    regions on chromosome 2 in *H. erato* and chromosome 6 in *H. melpomene*. The regions on
506    chromosome 2 in the Peruvian *H. erato* population are also associated with colour pattern, but such
507    association could be due to the high covariation of colour pattern and sampling location in this
508    population. These regions overlap with predicted genes, including basic metabolic genes and a heat
509    shock protein (Table 3), which could be candidates for adaptation to different temperature regimes.
510    Chemosensory genes were also detected on chromosome 2, and could be candidates for divergent
511    mate preference or host plant adaptation (Briscoe et al. 2013). However, no differences in host plant
512    preference have been observed in Peru where these outliers were detected, and mating within the
513    hybrid zone appears to the random (Mallet and Barton 1989a), although marginal differences in
514    mate preference have been observed in *H. melpomene* (Merrill, Gompert, et al. 2011).

515    There appear to be multiple dispersed divergent regions on chromosome 2 in *H. erato* (SI figure 4).
516    These could be evidence of divergence hitchhiking, whereby new mutations that cause differential
517    fitness are more likely to be fixed by selection if they arise close to other loci already under divergent
518    selection. This could lead to clustering of divergently selected loci in the genome (Via 2012; Feder et
519    al. 2012). The same process could also have led to additional loci under divergent selection to arise
520    in linkage with the colour pattern loci. This could explain the second divergent and altitude
521    associated region on chromosome 18 (linked to the *B/D* locus) in *H. melpomene* (Table 3, SI figure 4).
522    One possibility is that this could be the *B/D* linked mate preference locus that has previously been
523    identified (Merrill, Van Schooten, et al. 2011), although it is not clear if the mate preference locus is
524    an additional linked locus or a pleiotropic effect of the wing colour locus itself. It is also possible that
525    these apparently distinct but linked divergent regions could simply reflect the heterogeneous nature
526    of $F_{ST}$ resulting from strong divergent selection acting on a single locus combined with other
527    background and neutral processes (Charlesworth et al. 1997). A broad region of divergence around
528    the *B/D* locus in *H. melpomene* would fit with other suggestions that it has undergone stronger or
529    more recent selection than other colour pattern loci (Nadeau et al. 2012). The *D* region in *H. erato*
530    does not appear to be extended in the same way as in *H. melpomene* (SI figure 4), suggesting that
531    either the architecture or the selective history of this region is different between these species.

532    **Comparison of the genomic architecture of divergence between convergent species**

533   One interesting question that can be addressed with our results is the extent to which species
534   undergoing parallel divergence will show parallel patterns at the genomic level. In order to address
535   this we first need to know whether the species really have undergone parallel divergence, *i.e.* that
536   both the phenotypic start and end points have been similar. Several previous studies have suggested
537   that this is not the case and that *H. erato* diverged earlier and followed a different trajectory
538   compared to *H. melpomene* (Quek et al. 2010; Brower 1996; Flanagan et al. 2004). However, our
539   phylogenetic results are more consistent with a recent analysis suggesting that the two species do
540   appear to have undergone co-divergence in multiple populations across their range (Cuthill and
541   Charleston 2012). Our results are based on significantly more data than any of the previous analyses
542   (>5Mb in *H. melpomene* and >1Mb in *H. erato*), and should produce a better signal for phylogenetic
543   analysis as compared to AFLPs used previously (Quek et al. 2010). Although the striking similarities in
544   tree topology do seem to support the co-divergence hypothesis, alignment to a reference genome
545   means that the evolutionary rates in our data for *H. erato* and *H. melpomene* are not directly
546   comparable. In addition to the phylogenetic signal, our data also suggested similar patterns of
547   population structure between species in each of the regions, with higher background divergence
548   levels in Ecuador as compared to Peru (Figure 2, Table 1, SI figure 3). This is despite the fact that
549   average distance between sampling locations of "pure" subspecies individuals was similar for both
550   hybrid zones (~56km in Ecuador and 58-60km in Peru).

551   Although some loci show parallel divergence in both species (*B/D* in Peru; *B/D* and *Ac/Sd* in
552   Ecuador), there is surprisingly little similarity in which other loci are divergent between subspecies
553   within each species. This is contrary to the general perception that there are strong genetic parallels
554   in this system (Joron et al. 2006; Baxter et al. 2008; Supple et al. 2013; Papa et al. 2008). Some of
555   these differences were known previously, for example, that in Peru the *Sd/Ac* locus controls band
556   shape variation in *H. erato* but not in *H. melpomene* (Mallet 1989). Our results extend this further
557   through the identification of the *Ro* locus on chromosome 13 in Ecuadorian *H. erato*, which is not
558   divergent in its co-mimic *H. melpomene,* and the identification of divergent regions of chromosome
559   2 in *H. erato* and chromosome 6 in Peruvian *H. melpomene*.

560   In general, it seems that although the same colour pattern loci are present in both species (Joron et
561   al. 2006; Baxter et al. 2008; Martin et al. 2012) they are being used in different ways and
562   combinations in order to produce convergent phenotypes. This is particularly surprising given the
563   pattern of co-divergence observed in the phylogeny, which would appear to suggest that similar
564   colour patterns have arisen at a similar time and from similar ancestral forms in both species.
565   Nonetheless, the apparent pattern of co-divergence could simply reflect more recent patterns of
566   gene flow between geographically proximate populations in both species. This has recently been
567   highlighted by studies showing that patterns of divergence at colour pattern controlling loci can be
568   very different to those found at the rest of the genome (Hines et al. 2011; Pardo-Diaz et al. 2012;
569   The *Heliconius* Genome Consortium 2012; Supple et al. 2013). Therefore, the differences that we
570   observe in the use of particular loci in the two species could reflect different mimetic histories that
571   will only be resolved by studies of the evolutionary history of particular loci.

**Discovery of a new cryptic *H. timareta* population**

573   An unexpected finding of our study was the discovery of a previously undescribed population of *H.*
574   *timareta*, which appears phenotypically virtually indistinguishable from *H. melpomene malleti* in

575   Ecuador but is clearly genetically distinct (Figure 1, SI figures 1 and 2). *H. timareta florencia* is a
576   *malleti* like population that has previously been described in Colombia and also co-occurs with *H.*
577   *melpomene malleti*. In that population the length of the red line on the anterior edge of the ventral
578   forewing was diagnostic (Giraldo et al. 2008). This character was not diagnostic in our genotyped
579   individuals, with overlapping length distributions between the species (data not shown). We noted a
580   tendency towards *H. timareta* having a shorter line on average, but given the small sample sizes in
581   the current study, this remains to be confirmed.

582   A polymorphic high altitude population of *H. timareta* (*H. timareta timareta*) also occurs in this area
583   of Ecuador, overlapping in distribution with *H. melpomene plesseni*. The polymorphism in this
584   population has been somewhat of a puzzle as none of the forms mimic other co-occurring butterflies
585   (Mallet 1999). Our finding of a new *H. timareta* population may help to explain the polymorphism in
586   *H. timareta timareta,* if it is being generated in part by gene flow from this newly identified
587   population.

588   The *H. timareta* radiation has only been recognised in the last 10 years (Giraldo et al. 2008; Jiggins
589   2008, Merot et al., 2013). The *H. timareta* individuals in our study were collected from sites at 824m
590   and 376m. They appear to be fairly common at low altitude as four out of the five individuals
591   sequenced from the site at 376m were *H. timareta*. In a large dataset compiled by Rosser et al.
592   (2012) containing 232 *H. timareta* individuals from all known populations (including *H. tristero*, now
593   thought to be a subspecies of *H. timareta*), the lowest sampling location is around 600m, with 95%
594   of individuals occurring over 800m. Therefore, the population of *H. timareta* that we have
595   discovered occurs below the usual altitudinal range of *H. timareta*. This extends the possible range
596   of this species and suggests that the overlap in distribution of *H. timareta* and *H. melpomene* is
597   greater than previously considered.

598   **Conclusions**

599   We have demonstrated that high resolution genome scans using admixed individuals from hybrid
600   zones can be used to identify loci underlying phenotypic variation. Only a small proportion of the
601   genome (about 0.025%) is strongly differentiated between subspecies and most of this can be
602   explained by divergence at loci controlling colour pattern. This is consistent with previous studies
603   based on smaller numbers of markers (Turner 1979; Baxter et al. 2010, Counterman et al. 2010,
604   Nadeau et al. 2012) and suggests that the hybrid zones are ancient or have formed in primary
605   contact, and are maintained by strong selection on colour pattern (Mallet and Barton 1989a, Mallet
606   2010). However, we also find, for the first time, some divergent loci that do not appear to be
607   associated with colour pattern, suggesting that there may be other differences between subspecies.
608   This could explain why several *Heliconius* hybrid zones occur across ecological gradients (Benson
609   1982), if they are coupled with extrinsic selection acting on other loci in the genome (Bierne et al.
610   2011). However, this needs to be confirmed with detailed phenotypic analyses of the subspecies to
611   identify whether differences are present that could be explained by ecological adaptation. In general
612   we find that, although some loci are divergent in all populations, the genomic pattern of divergence
613   between co-mimetic species is not particularly similar, suggesting that the level of parallel genetic
614   evolution between *H. erato* and *H. melpomene* is in fact quite low, despite parallel phylogenetic
615   patterns of divergence. Finally, our analysis shows that alignment to a distantly related reference
616   genome can improve analyses over a *de novo* assembly of the data.

617    **METHODS**

618    **Samples and sequencing**

619    30 *H. erato* and 30 *H. melpomene* individuals were selected from a larger sample taken from the
620    hybrid zone region in Peru. Similarly, 30 *H. erato* and 30 *H. melpomene* were also selected from a
621    larger study of a subspecies hybrid zone in Ecuador (Salazar 2012). Each set of 30 samples comprised
622    10 pure forms of each subspecies and 10 hybrids (based on colour pattern). See Figure 2 and SI table
623    2 for further details of the samples and locations.

624    RAD sequencing libraries were prepared using previously described methodologies (The *Heliconius*
625    Genome Consortium 2012; Baird et al. 2008; Baxter et al. 2011). Briefly, DNA was digested with the
626    restriction enzyme PstI prior to ligation of P1 sequencing adaptors with five-base molecular
627    identifiers (see SI table 2 for MIDs used).  We then pooled samples into groups of 6 before shearing,
628    ligation of P2 adaptors, amplification and fragment size selection (300-600bp). Libraries were then
629    further pooled such that 30 individuals were sequenced on each lane of an Illumina HiSeq2000
630    sequencer to obtain 150 base paired-end sequences. We obtained an average of 374M sequence
631    pairs from each lane.  Following sequencing, three of the *H. erato* individuals from Peru were found
632    to have been incorrectly assigned to this species and were excluded from all further analyses.

633    In order to compare patterns of phylogenetic divergence of the focal subspecies, we also used
634    sequence data from additional subspecies and closely related species in each group.  Two individuals
635    each from 6 additional *H. erato* populations and the closely related *H. himera* were also PstI RAD
636    sequenced with 5 individuals pooled per lane of Illumina GAIIx (100 base paired-end sequencing).
637    These sequences were obtained in the same run as a comparable set of individuals from the *H.*
638    *melpomene* clade, which have been used in previous analyses (The *Heliconius* Genome Consortium
639    2012; Nadeau et al. 2013, European Nucleotide Archive, Accession ERP000991). We also obtained
640    whole-genome shotgun sequence data from an outgroup species, *H. clysonimus*, which was
641    sequenced on a fifth of a HiSeq2000 lane, giving 53.5M 100 base read pairs for this individual.

642    **Alignment to reference genome**

643    We separated paired-end reads by MID using the RADpools script in the RADtools (v1.2.4) package
644    (Baxter et al. 2011), which also filters based on the presence of the restriction enzyme cut site, using
645    the option to allow one mismatch within the MID.  Reads from each individual were then aligned to
646    the *H. melpomene* reference genome (The *Heliconius* Genome Consortium 2012) using Stampy
647    v1.0.17 (Lunter and Goodson 2011), with default parameters except substitution rate, which was set
648    to 0.03 for alignments of *H. melpomene* and 0.10 for alignments of *H. erato*.

649    We then realigned indels and called genotypes using the Genome Analysis Tool Kit (GATK) v1.6.7
650    (DePristo et al. 2011), emitting all confident sites (those with quality ≥30). This was first run on each
651    set of 30 (or 27) individuals from each population group. These genotype calls were used for
652    analyses of genetic variation within each of the groups, including outlier detection, association
653    mapping and analyses of subpopulation structure.  In addition, genotype calling was also performed
654    on a combined dataset of all *H. melpomene* and outgroup taxa (*H. timareta, H. cydno* and *H. hecale*)
655    as well as a combined set of all *H. erato* and its outgroups (*H. himera* and *H. cylsonimus*). These
656    genotype calls were used for the phylogenetic analyses and broader analyses of genetic structure.

657　For all downstream analyses, calls were further filtered to only accept those based on a minimum
658　depth of five reads and minimum genotype and mapping qualities (GQ and MQ) of 30 for *H.*
659　*melpomene* and 20 for *H. erato.*

660　**De novo assembly**

661　We quality-filtered the single-end raw sequence data and separated sequences by MID with the
662　process_radtags program within Stacks (Catchen et al. 2011). This program corrects single errors in
663　the MID or restriction site and then checks quality score using a sliding window across 15% the
664　length of the read.  We discarded sequences with a raw phred score below 10, removed reads with
665　uncalled bases or low quality scores, and trimmed reads to 100 bases to eliminate potential
666　sequencing error occurring at ends of reads. Table 1 shows the mean read numbers per individual
667　obtained after filtering. For each population group, we assembled loci *de novo* using the
668　denovo_map.sh  pipeline in Stacks (Catchen et al. 2011). We set the minimum depth of coverage (m)
669　to 6, allowed 4 mismatches both in creating individual stacks (M) and in secondary reads (N), and
670　removed or separated highly repetitive RadTags. Due to the high level of polymorphism in our
671　dataset, we used these parameters to minimize the exclusion of interesting loci with high variability
672　between populations.  *De novo* assembly was conducted both including (for association mapping)
673　and excluding (for bayescan outlier detection) hybrid individuals in the analysis. Individuals from
674　Ecuador that were identified as being *H. timareta* were excluded.

675　**Phylogenetics and analysis of population structure**

676　Only the reference aligned data were used for phylogenetics and STRUCTURE analyses. We used
677　custom scripts to convert from vcf to Phylip format and to filter sites with a minimum of 95% of
678　individuals with confident calls. Maximum likelihood phylogenies were constructed in PhyML
679　(Guindon and Gascuel 2003) with a GTR model using the resulting 5,737,351 sites (including
680　invariant sites) for the *H. melpomene* group and 1,693,024 sites for the *H. erato* group. Approximate
681　likelihood branch supports were calculated within the program.

682　Population structure within and across each of the hybrid zones was analysed using the program
683　STRUCTURE v2.3 (Pritchard et al. 2000). We prepared input files using custom scripts, and only sites
684　with 100% of individuals present for *H. melpomene* populations or at least 75% of individuals
685　present for *H. erato* populations and with a minor allele frequency of at least 20% were retained.
686　This reduced the number of sampled sites, keeping just the most informative ones, for easier
687　handling by the program. Initial short runs ($10^3$ burn-in, $10^3$ data collection, K=1) were used to
688　estimate the allele frequency distribution parameter λ. We then ran longer clustering runs ($10^4$ burn-
689　in, $10^4$ data collection) with the obtained values of λ for each of the four population groups for K=1-
690　3. For *H. melpomene* in Ecuador the analysis was first run with all individuals included and then
691　excluding the individuals identified as being *H. timareta*.

692　We also performed principle components analysis of the genetic variation in each population group.
693　This was done with the "cmdscale" command in R (R Development Core Team 2011), using genetic
694　distance matrices calculated as 0.5-ibs, where ibs was the identity by sequence matrix calculated in
695　GenABEL (see below). As further confirmation that some of the *H. melpomene* individuals sampled in
696　Ecuador were in fact cryptic *H. timareta*, we also performed principle components analysis on the
697　combined *H. melpomene* and outgroup data set.  We also ran principle components analysis on the

698    *de novo* assembled data for each population group, to test whether both methods were detecting
699    similar underlying patterns of genetic variation.

700    In order to compare our newly identified *H. timareta* individuals to other populations, we Sanger
701    sequenced a 745bp region of mitochondrial COI that overlapped with the regions sequenced in
702    previous studies (Giraldo et al. 2008; Mérot et al. 2013). This was PCR amplified as in Mérot et al.
703    (2013) with primers "Jerry" and "Patlep" and directly sequenced with "Patlep". These sequences
704    were then aligned with those available on Genbank and a maximum likelihood phylogeny was
705    constructed in PhyML (Guindon and Gascuel 2003) with a GTR model and 1000 bootstrap replicates.

706    **Association mapping of loci controlling colour pattern variation**

707    We scored components of phenotypic variation that segregate across each of the hybrid zones. The
708    scored phenotypes are shown in Figure 3 (for Peru) and Figure 4 (for Ecuador) and listed in full in SI
709    Table 3. These were scored mostly as binomial (1,0) traits, but in some cases intermediates were
710    also scored (as 0.5). The width and shape of the forewing band was scored based on whether it
711    extended into each of the wing "cells", demarcated by the major wing veins (as shown in SI figure 5).
712    In Peruvian populations, the size and shape of the forewing band was measured as two components
713    (Figure 3C/F) that extend the band distally (cell spot 8) and proximally (cell spot 11). In Ecuador,
714    three aspects of band shape were scored: cells 8 and 11, which make up the proximal spot in *H. m.*
715    *plesseni* and *H. e. notabilis*, and cell 7, which pushes the band towards the wing margin in *H. m.*
716    *malleti* and *H. e. lativitta* (Figure 4C/F). In our sample of *H. melpomene* the presence of cell spots 8
717    and 11 were perfectly correlated, whereas in *H. erato* the presence of cell spot 7 was perfectly
718    correlated with the absence of cell spot 8. In addition, individuals were also scored for their
719    predicted genotypes at major loci described previously (with predicted heterozygotes scored as 0.5)
720    (Sheppard et al. 1985; J. Mallet 1989) and the altitude at which they were collected was included as
721    a continuous phenotypic trait.

722    We performed association mapping using the R Package GenABEL v 1.7-4 (Aulchenko et al. 2007).
723    This was performed on both the *de novo* assembled and the reference aligned data with a custom
724    script used to convert both from vcf to Illumina SNP format. Individuals identified as being *H.*
725    *timareta* were excluded. Filtering was performed within the program to remove sites with >30%
726    missing data and with a minor allele frequency of <3%.

727    For each population, an analysis of the hind-wing ray phenotype using the reference mapped data
728    was first performed using three methods: a straight score test (qtscore), a score test with the first
729    three principle components of genetic variation (calculated as described above) as covariates, and
730    an EIGENSTRAT analysis (egscore, Price et al. 2006). The presence of genetic stratification and the
731    ability of these methods to correct for this was analysed by comparing the inflation factor, λ. In all
732    cases the analyses incorporating population stratification did not give a reduced value of λ and so
733    were not used for subsequent analyses. As our samples were from hybrid zones with >60% of the
734    samples having extreme values of all scored phenotypes, we would expect similar levels of
735    stratification for all phenotypes, so this test for stratification was not repeated for all phenotypes.

736    We therefore performed score tests for all scored phenotypes across all population groups.
737    Genome-wide significance was determined empirically from 1000 resampling replicates and
738    corrected for population structure using the test specific λ.

739     **BayeScan analysis to identify loci under selection**

740     We used the program BayeScan v2.1 (Foll and Gaggiotti 2008) to look for loci with outlier $F_{ST}$ values
741     between "pure" individuals of each subspecies type (based on wing colour pattern) in each
742     population group. Exclusion of the *H. timareta* individuals meant that only three pure *H. melpomene*
743     *malleti* individuals remained. Therefore, for the purpose of this analysis of *H. melpomene* in Ecuador,
744     the two hybrid individuals closest to the *H. m. malleti* side of the hybrid zone (Figure 2), which also
745     had the most *H. m. malleti* like phenotypes, were included as *H. m. malleti*.

746     The program was run with the prior odds for the neutral model (pr_odds) set to 10 and outlier loci
747     were detected with a false discover rate (FDR) of 0.05. We ran this analysis using both the *de novo*
748     assembled and the reference aligned data. Custom scripts were used to convert these to the correct
749     input format. For both analyses, sites were only kept if at least 75% of individuals were sampled for
750     both subspecies in a given comparison.

751     **DATA ACCESS**

752     DNA sequence reads have been submitted to the European Nucleotide Archive, Accession
753     PRJEB4669. COI sequences have been deposited in EMBL-Bank, accessions HG710096 - HG710125.
754     Custom scripts and wing images are available on request from the authors.

755

756     **ACKNOWLEDGEMENTS**

769     **FIGURE LEGENDS**

770     **Figure 1.** A) Distribution in South America of the subspecies included in this study. B) Maximum
771     likelihood phylogenies with approximate likelihood branch supports. Co-mimics from outside of the
772     focal hybrid zones are connected with dotted lines. Focal hybrid zone individuals are shown in
773     colour: blue, *H. m. plesseni* and *H. e. notabilis*; purple, Ecuador hybrids; dark red, *H. m. malleti* and *H.*
774     *e. lativitta*; red, *H. m. aglaope* and *H. e. favorinus*; orange, Peru hybrids; yellow, *H. m. amaryllis* and
775     *H. e. emma.* Additional populations are in black. Country abbreviations: Ec, Ecuador; FG, French
776     Guiana; Co, Colombia; Pa, Panama.

777     **Figure 2.** Population structure at each of the hybrid zones using the reference aligned data. A)
778     Sampling locations with altitude (sample number) and pie charts of the proportion of individuals of
779     each type sampled from each site. Colours are the same as in Figure 1 except black indicates *H.*
780     *timareta* in Ecuador. B) STRUCTURE analysis with k=2 (*H. timareta* individuals excluded). Each
781     individual is shown as a horizontal bar with the allelic contribution from population 1 (grey) and
782     population 2 (black) C) Principle components analysis. D) Distribution of $F_{ST}$ values from BayeScan.

783     **Figure 3.** Association mapping (A and D) and outlier analysis (B and E) for *H. melpomene* (A, B, C) and
784     *H. erato* (D, E, F) in Peru. Each phenotype used for the association mapping is shown in a different
785     colour as illustrated in panels C and F. For clarity, only the top 20 associated SNPs are shown for
786     each phenotype. Results from the *de novo* assembled data are shown as crosses (and in orange for
787     the outlier analysis) and positioned based on the top BLAST hit to the *H. melpomene* genome, those
788     with thinner lines were not confidently or uniquely assigned to these positions (eg. those at the end
789     of Chromosome 10 in D). "unmapped" indicates scaffolds of the *H. melpomene* reference genome
790     that were not assigned to chromosomes in v1.1 of the genome assembly.

791     **Figure 4.** Association mapping (A and D) and outlier analysis (B and E) for *H. melpomene* (A, B, C) and
792     *H. erato* (D, E, F) in Ecuador. See Figure 3 legend for further information.

793     **Figure 5.** Venn diagrams of SNPs detected in the *de novo* assembled (orange and red) and reference
794     aligned (blue and green) data by BayeScan outlier detection (red and blue) and association mapping
795     (orange and green), for each of the four populations.

796

797

798     **TABLES**

799     **Table 1. Summary statistics from alignment and assembly approaches**

**_De novo_ assembly with Stacks - single end reads**

| | | n | millions of reads (mean ±SD) | | bases covered (x10$^6$) | bases covered in ≥ 10 inds (x10$^6$) | SNPs used in outlier analysis (x10$^3$) | mean $F_{ST}$ | Outliers | Significant Phenotypic Associations |
|---|---|---|---|---|---|---|---|---|---|---|
| _H. erato_ | Peru | 27 | 8.0 | ±2.2 | 166 | 2.8 | 37 | 0.0280 | 22 | 10 |
| | Ecuador | 30 | 9.2 | ±2.5 | 149 | 3.3 | 31 | 0.0568 | 0 | 2 |
| _H. melpomene_ | Peru | 30 | 7.0 | ±2.4 | 61 | 4.3 | 57 | 0.0145 | 23 | 8 |
| | Ecuador | 22 | 7.8 | ±1.4 | 45 | 3.5 | 43 | 0.0310 | 5 | 4 |

**Aligned to _H. melpomene_ reference - paired end with Stampy**

| | | n | millions of reads (mean ±SD) | | bases covered (x10$^6$) | bases covered in ≥ 10 inds (x10$^6$) | SNPs used in outlier analysis (x10$^3$) | mean $F_{ST}$ | Outliers | Significant Phenotypic Associations | mean gap between RAD loci (kb) | max gap between RAD loci (kb) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| _H. erato_ | Peru | 27 | 11.3 | ±3.2 | 11 | 4.2 | 373 | 0.0142 | 19 | 28 | 9.4 | 116 |
| | Ecuador | 30 | 11.9 | ±3.2 | 13 | 5.1 | 337 | 0.0316 | 56 | 15 | 9.1 | 105 |
| _H. melpomene_ | Peru | 30 | 10.9 | ±3.8 | 28 | 14.4 | 860 | 0.0112 | 235 | 91 | 9.3 | 103 |
| | Ecuador | 22 | 10.7 | ±1.9 | 23 | 15.6 | 788 | 0.0299 | 179 | 14 | 9.5 | 114 |

800

801

802     **Table 2. Positions of the strongest associations and outliers identified in genomic regions known to**
803     **control colour pattern variation, in relation to known genes and functional regions.**

| Colour pattern loci | | | _B/D_ | _Ac/Sd_ | _Yb/N/Cr_ |
|---|---|---|---|---|---|
| | Chromosome | | chr18 | chr10 | chr15 |
| | Scaffold | | HE670865 | HE668478 | HE667780 |
| | Gene | | HMEL001028 (_optix_)[1] | HMEL018100 (_WntA_)[2] | Presently unknown |
| | - Position | | 438,423-439,107 | 450,400-483,854 | |
| | Functional region[3] | | 300,000-400,000 | Presently unknown | 600,000-1,000,000 |
| _H. melpomene_ | Peru | assoc | 161,328-376,651 | N/A | 676,543-697,543 |
| | | outlier | 263,358 | | 676,645 |
| | Ecuador | assoc | none | 454,479-454,496 | 697,118-725,562 |
| | | outlier | 376,651 | 454,404 | 697,118 |
| _H. erato_ | Peru | assoc | 362,793-362,794 | 444,914 | none |
| | | outlier | 362,794 | 478,952 | none |
| | Ecuador | assoc | 282,473-376,342 | 478,897 | N/A |
| | | outlier | 376,250 | 479,220 | |

804     For each population, positions are given for the SNPs showing the strongest phenotypic associations
805     (assoc) and the highest $F_{ST}$ outliers (outlier): N/A, not expected or found; none, not found. [1] from
806     Reed et al. 2011; [2] from Martin et al. 2012; [3] Inferred from population genomics: the _B/D_ region

807    appears to be similar in *H. erato* and *H. melpomene; Yb/N/Cr* region has been localised in *H.*

808    *melpomene* only (Nadeau et al. 2012; Supple et al. 2013).

809

810    **Table 3. Novel genomic regions showing phenotypic associations or divergence outliers**

| chrom. | Scaffold | Position¥ | Comparison* | association (p value)‡ | $F_{ST}$‡ | Closest Gene | Distance† | GO function | putative protein |
|---|---|---|---|---|---|---|---|---|---|
| unmapped (chr13) | HE669551 | 4,467 | *erato Ecuador:* refmap assoc (**Ro**, spot 7/8), refmap outlier | 0.007 | 0.307 | HMEL004352 | 0 (S) | microtubule binding | radial spoke head 3 |
| chr13 | HE670984 | 4,951-4,959 | *erato Ecuador:* refmap assoc (Ro, **spot 11**, spot 7/8), refmap outlier | 0.024 | 0.343 | HMEL009926 | 3,915 | structural constituent of ribosome, RNA binding | ribosomal protein S4 |
| unmapped (chr20) | HE671554 | 82,726-82,851 | *melp Peru:* refmap assoc (spot 8) | 0.017 | | HMEL016146 | 0 (A/S/I) | protein binding, zinc ion binding | MICAL-like |
| chr17 | HE671853 | 190,306 | *erato Ecuador:* refmap assoc (HWY) | 0.001 | | HMEL014236 | 0 (I) | catalytic activity, serine-type endopeptidase activity | serine protease 30 |
| unmapped | HE670458 | 97 | *melp Ecuador:* refmap assoc (rays); *melp Peru:* refmap assoc (alt, **rays, D gen**) | 0.001 | | no genes on this scaffold, in repetitive region | | | |
| chr18 | HE671488 | 228,759-228,922 | *melp Peru:* refmap assoc (**alt**, D gen), refmap outlier, denovo outlier; *melp Ecuador:* refmap outlier | 0.006 | 0.329 | HMEL014920 | 19,171 | | |
| chr2 | HE670771 | 179,067-179,278 | *erato Peru:* refmap assoc (alt, rays, **D gen,** spot 11); *erato Ecuador:* denovo assoc (spot 11) | 0.003 | | HMEL008318 | 0 (I) | catalytic activity, protein binding | fatty acid synthase |
| chr2 | HE670771 | 199,742 | *erato Peru:* denovo assoc (**alt, rays**), denovo outlier | 0.014 | 0.217 | HMEL008322 | 0 (A) | odorant binding | odorant binding protein 7 |
| chr2 | HE670519 | 293,177-293,603 | *erato Peru:* refmap outlier, denovo assoc (spot 11), denovo outlier | 0.021 | 0.158 | HMEL007059 | 0 (I/A) | oxidoreductase activity | 3-dehydroecdysone 3alpha-reductase |
| chr2 | HE670235 | 24,536-24,614 | *erato Peru:* refmap assoc (alt, rays, **D gen**, spot 11) | 0.003 | | HMEL005708 | 56,981 | taste receptor activity | olfactory receptor 4 |

| chr2 | HE671428 | 6,795 | *erato Ecuador:* refmap outlier | 0.143 | HMEL014154 | 0 (S) | choline dehydrogenase activity, oxidoreductase activity, acting on CH-OH group of donors, flavin adenine dinucleotide binding | glucose dehydrogenase |
| chr2 | HE671428 | 99,432 | *erato Ecuador:* refmap outlier | 0.145 | HMEL014163 | 0 (A) | | heat shock protein 70 |
| chr6 | HE671933 | 31,161-31,367 | *melp Peru:* denovo outlier, refmap outlier | 0.175 | HMEL016074 | 7,925 | oxidoreductase activity | amine oxidoreductase |
| chr6 | HE671934 | 15,458 | *melp Peru:* ref map outlier | 0.104 | HMEL016075 | 4,121 | oxidoreductase activity | amine oxidoreductase |

811  ¥The position of the SNP with the strongest result in a given region, unless multiple equally
812  supported SNPs are present in which case a range is given. *Analysis in which SNP is detected: melp,
813  *H. melpomene*; erato, *H. erato*; ref map, reference aligned data; denovo, *de novo* assembled data;
814  outlier, BayeScan outlier analysis; assoc, association analysis (rays, presence of hindwing rays and
815  fore/hindwing dennis patches; Dgen, predicted *B/D* genotype; spot, presence of non-black colour in
816  that wing cell; alt, altitude;; Ro, rounding of distal edge of forewing band). The analysis giving the
817  most significant result is shown in bold. ‡Value for the most significant analysis is given. ₮ If a SNP is
818  within a gene (distance=0) then in parenthesis: A, non-synonymous; S, synonymous; I, within an
819  intron.

820

821     **REFERENCES**

822     Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic Local Alignment Search Tool.
823             *Journal of Molecular Biology* **215**: 403–410. doi:10.1016/S0022-2836(05)80360-2.
824     Aulchenko YS, Ripke S, Isaacs A, van Duijn CM. 2007. GenABEL: An R Library for Genome-wide
825             Association Analysis. *Bioinformatics* **23**: 1294–1296. doi:10.1093/bioinformatics/btm108.
826     Avidor-Reiss T, Maer AM, Koundakjian E, Polyanovsky A, Keil T, Subramaniam S, Zuker CS. 2004.
827             'Decoding Cilia Function: Defining Specialized Genes Required for Compartmentalized Cilia
828             Biogenesis'. *Cell* 117 (4) (May 14): 527–539. doi:10.1016/S0092-8674(04)00412-X.
829     Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA, Johnson EA.
830             2008. Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers. *PLoS ONE*
831             **3**: e3376. doi:10.1371/journal.pone.0003376.
832     Barton NH, Hewitt GM. 1985. Analysis of Hybrid Zones. *Annual Review of Ecology and Systematics*
833             **16**: 113–148. doi:10.1146/annurev.es.16.110185.000553.
834     Baxter SW, Johnston SE, Jiggins CD. 2008. Butterfly Speciation and the Distribution of Gene Effect
835             Sizes Fixed During Adaptation. *Heredity* **102**: 57–65.
836     Baxter SW, Davey JW, Johnston JS, Shelton AM, Heckel DG, Jiggins CD, Blaxter ML. 2011. Linkage
837             Mapping and Comparative Genomics Using Next-Generation RAD Sequencing of a Non-
838             Model Organism. *PLoS ONE* **6**: e19315. doi:10.1371/journal.pone.0019315.
839     Baxter SW, Nadeau NJ, Maroja LS, Wilkinson P, Counterman BA, Dawson A, Beltran M, Perez-Espona
840             S, Chamberlain N, Ferguson L. et al. 2010. Genomic Hotspots for Adaptation: The Population
841             Genetics of Mullerian Mimicry in the *Heliconius melpomene* Clade. *PLoS Genetics* **6**:
842             e1000794. doi:10.1371/journal.pgen.1000794
843     Baxter SW, Papa R, Chamberlain N, Humphray SJ, Joron M, Morrison C, ffrench-Constant RH,
844             McMillan WO, Jiggins CD. 2008. Convergent Evolution in the Genetic Basis of Müllerian
845             Mimicry in Heliconius Butterflies. *Genetics* **180**: 1567 –1577.
846             doi:10.1534/genetics.107.082982.
847     Beavis WD 1997. QTL Analyses: Power Precision and Accuracy. *Molecular Dissection of Complex*
848             *Traits*. (Paterson AH) 145–162. CRC Press.
849     Bierne N, Welch J, Loire E, Bonhomme F, David P. 2011. The Coupling Hypothesis: Why Genome
850             Scans May Fail to Map Local Adaptation Genes. *Molecular Ecology* **20**: 2044–2072.
851             doi:10.1111/j.1365-294X.2011.05080.x.
852     Benson, WW. 1972. Natural Selection for Miillerian Mimicry in *Heliconius erato* in Costa Rica. *Science*
853             **176**: 936–939. doi:10.1126/science.176.4037.936.
854     Benson, WW. 1982. Alternative models for infrageneric diversification in the humid tropics: tests
855             with passion vine butterflies. *Biological Diversification in the Tropics* (ed. GT Prance), pp.
856             608–640. Columbia Univ. Press, New York.
857     Briscoe AD, Macias-Muñoz A, Kozak KM, Walters JR, Yuan F, Jamie GA, Martin SH, Dasmahapatra KK,
858             Ferguson LC, Mallet J. et al. 2013. Female Behaviour Drives Expression and Evolution of
859             Gustatory Receptors in Butterflies. *PLoS Genet* **9**: e1003620.
860             doi:10.1371/journal.pgen.1003620.
861     Brower AV. 1996. Parallel Race Formation and the Evolution of Mimicry in *Heliconius* Butterflies: A
862             Phylogenetic Hypothesis from Mitochondrial DNA Sequences. *Evolution* **50**: 195–221.
863             doi:10.2307/2410794.
864     Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH. 2011. Stacks: Building and
865             Genotyping Loci De Novo From Short-Read Sequences. *G3: Genes, Genomes, Genetics* **1**:
866             171–182. doi:10.1534/g3.111.000240.
867     Charlesworth B, Nordborg M, Charlesworth D. 1997. The Effects of Local Selection, Balanced
868             Polymorphism and Background Selection on Equilibrium Patterns of Genetic Diversity in
869             Subdivided Populations. *Genetical Research* **70**: 155–174.
870     Counterman BA, Araujo-Perez F, Hines HM, Baxter SW, Morrison CM, Lindstrom DP, Papa R,
871             Ferguson L, Joron M, ffrench-Constant RH, et al. 2010. Genomic Hotspots for Adaptation:

872     The Population Genetics of Müllerian Mimicry in *Heliconius erato*. *PLoS Genetics* **6**:
873         e1000796. doi:10.1371/journal.pgen.1000796.
874 Crawford JE, Nielsen R. 2013. Detecting Adaptive Trait Loci in Non-model Systems: Divergence or
875         Admixture Mapping? *Molecular Ecology* in press. doi:10.1111/mec.12562.
876 Cuthill JH, Charleston M. 2012. Phylogenetic Codivergence Supports Coevolution of Mimetic
877         Heliconius Butterflies. *PloS One* **7**: e36464. doi:10.1371/journal.pone.0036464.
878 DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del Angel G, Rivas
879         MA, Hanna M, et al. 2011. A Framework for Variation Discovery and Genotyping Using Next-
880         generation DNA Sequencing Data. *Nat Genet* **43**: 491–498. doi:10.1038/ng.806.
881 Feder JL, Gejji R, Yeaman S, Nosil P. 2012. Establishment of New Mutations Under Divergence and
882         Genome Hitchhiking. *Phil. Trans. Roy. Soc. B* **367**: 461–474. doi:10.1098/rstb.2011.0256.
883 Ferguson L, Lee SF, Chamberlain N, Nadeau NJ, Joron M, Baxter S, Wilkinson P, Papanicolaou A,
884         Kumar S, Kee T-J, et al. 2010. Characterization of a Hotspot for Mimicry: Assembly of a
885         Butterfly Wing Transcriptome to Genomic Sequence at the *HmYb/Sb* Locus. *Mol. Ecol.* **19**:
886         240–254. doi:10.1111/j.1365-294X.2009.04475.x.
887 Flanagan NS, Tobler, Davison AA, Pybus OG, Kapan DD, Planas S, Linares M, Heckel D, McMillan WO.
888         2004. Historical Demography of Müllerian Mimicry in the Neotropical Heliconius Butterflies.
889         *PNAS* **101**: 9704–9709. doi:10.1073/pnas.0306243101.
890 Foll M, Gaggiotti O. 2008. A Genome-Scan Method to Identify Selected Loci Appropriate for Both
891         Dominant and Codominant Markers: A Bayesian Perspective. *Genetics* **180**: 977–993.
892         doi:10.1534/genetics.108.092221.
893 Giraldo N., Salazar C, iggins C, Bermingham E, and Linares M. 2008. Two Sisters in the Same Dress:
894         Heliconius Cryptic Species. *BMC Evolutionary Biology* **8**: 324. doi:10.1186/1471-2148-8-324.
895 Gompert Z, Buerkle CA. 2011. A Hierarchical Bayesian Model for Next-Generation Population
896         Genomics. *Genetics* **187**: 903–917. doi: 10.1534/genetics.110.124693.
897 Gompert, Z, Lucas LK, Nice CC, Fordyce JA, Forister ML, Buerkle AC. 2012. Genomic regions with a
898         history of divergent selection affect fitness of hybrids between two butterfly species.
899         *Evolution* **66**: 2167–2181. doi:10.1111/j.1558-5646.2012.01587.x.
900 Guindon S, Gascuel O. 2003. A Simple, Fast, and Accurate Algorithm to Estimate Large Phylogenies
901         by Maximum Likelihood. *Systematic Biology* **52**: 696–704. doi:10.1080/10635150390235520.
902 Harrison RG. 1993. *Hybrid Zones and the Evolutionary Process*. Oxford University Press.
903 Hines HM, Counterman BA, Papa R, Albuquerque de Moura P, Cardoso MZ, Linares M, Mallet J, Reed
904         RD, Jiggins CD, Kronforst MR, et al. 2011. Wing Patterning Gene Redefines the Mimetic
905         History of Heliconius Butterflies. *PNAS* **108**: 19666–19671. doi:10.1073/pnas.1110096108.
906 Jiggins, CD. 2008. Ecological Speciation in Mimetic Butterflies. *BioScience* **58**: 541–548.
907 Joron M, Papa R, Beltrán M, Chamberlain N, Mavárez J, Baxter S, Abanto M, Bermingham E,
908         Humphray SJ, Rogers J, et al. 2006. A Conserved Supergene Locus Controls Colour Pattern
909         Diversity in Heliconius Butterflies. *PLoS Biology* **4**: e303. doi:10.1371/journal.pbio.0040303.
910 Kapan DD. 2001. Three-butterfly System Provides a Field Test of Mullerian Mimicry. *Nature* **409**:
911         338–340. doi:10.1038/35053066.
912 Kawakami T, Butlin RK. 2001. Hybrid Zones. *eLS*. John Wiley & Sons, Ltd.
913         http://onlinelibrary.wiley.com/doi/10.1002/9780470015902.a0001752.pub2/abstract.
914 Keller I, Wagner CE, Greuter L, Mwaiko S, Selz OM, Sivasundar A, Wittwer S, Seehausen O. 2013.
915         Population Genomic Signatures of Divergent Adaptation, Gene Flow and Hybrid Speciation in
916         the Rapid Radiation of Lake Victoria Cichlid Fishes. *Molecular Ecology* **22**: 2848–2863.
917         doi:10.1111/mec.12083.
918 Kocher TD, Sage RD. 1986. Further Genetic Analyses of a Hybrid Zone Between Leopard Frogs (Rana
919         Pipiens Complex) in Central Texas. *Evolution* **40**: 21–33. doi:10.2307/2408600.
920 Kronforst MR, Hansen MEB, Crawford NG, Gallant JR, Zhang W, Kulathinal RJ, Kapan DD, Mullen SP.
921         2013. Hybridization Reveals the Evolving Genomic Architecture of Speciation. *Cell Reports*.
922         doi:10.1016/j.celrep.2013.09.042.

923  Lamas G. 1997. Comentarios Taxonomicos y Nomenclaturales Sobre Heliconiini Neotropicales, Con
924        Designacion de Lectotipos y Descripcion de Cuatro Subespecies Nuevas (Lepidoptera:
925        Nymphalidae: Heliconiinae). *Rev. Per. Ent.* **40**: 111–125.
926  Langham GM. 2004. Specialized Avian Predators Repeatedly Attack Novel Color Morphs of
927        Heliconius Butterflies. *Evolution* **58**: 2783–2787. doi:10.1111/j.0014-3820.2004.tb01629.x.
928  Lexer C, Buerkle CA, Joseph JA, Heinze B, Fay MF. 2006. Admixture in European *Populus* Hybrid
929        Zones Makes Feasible the Mapping of Loci That Contribute to Reproductive Isolation and
930        Trait Differences. *Heredity* **98**: 74–84. doi:10.1038/sj.hdy.6800898.
931  Lunter G, Goodson M. 2011. Stampy: A Statistical Algorithm for Sensitive and Fast Mapping of
932        Illumina Sequence Reads. *Genome Research* **21**: 936–939. doi:10.1101/gr.111120.110.
933  Mallet J. 1989. The Genetics of Warning Colour in Peruvian Hybrid Zones of *Heliconius erato* and *H.*
934        *melpomene*. *Proc. Roy. Soc. B* **236**: 163–185.
935  Mallet J. 1999. Causes and Consequences of a Lack of Coevolution in Müllerian Mimicry.
936        *Evolutionary Ecology* **13**: 777–806. doi:10.1023/A:1011060330515.
937  Mallet J. 2010. Shift Happens! Shifting Balance and the Evolution of Diversity in Warning Colour and
938        Mimicry. *Ecological Entomology* **35**: 90–104. doi:10.1111/j.1365-2311.2009.01137.x.
939  Mallet J, Barton N. 1989a. Strong Natural Selection in a Warning-Color Hybrid Zone. *Evolution* **43**:
940        421–431. doi:10.2307/2409217.
941  Mallet J, Barton N. 1989b. Inference from Clines Stabilized by Frequency-dependent Selection.
942        *Genetics* **122**: 967–976.
943  Mallet, J, Barton N, Lamas G, Santisteban J, Muedas M, and Eeley H. 1990. Estimates of Selection and
944        Gene Flow from Measures of Cline Width and Linkage Disequilibrium in *Heliconius* Hybrid
945        Zones. *Genetics* **124**: 921–936.
946  Martin, A, Papa R, Nadeau NJ, Hill RI, Counterman BA, Halder G, Jiggins CD, Kronforst MR, Long AD,
947        McMillan WO, et al. 2012. Diversification of Complex Butterfly Wing Patterns by Repeated
948        Regulatory Evolution of a Wnt Ligand. *PNAS* **109**: 12632–12637.
949        doi:10.1073/pnas.1204800109.
950  Martin SH, Dasmahapatra KK, Nadeau NJ, Salazar C, Walters JR, Simpson F, Blaxter M, Manica A,
951        James Mallet, Jiggins CD. 2013. Genome-wide Evidence for Speciation with Gene Flow in
952        Heliconius Butterflies. *Genome Research* doi:10.1101/gr.159426.113.
953  Mérot C., Mavárez J, Evin A, Dasmahapatra KK, Mallet J, Lamas G, and Joron M. 2013. Genetic
954        Differentiation Without Mimicry Shift in a Pair of Hybridizing Heliconius Species
955        (Lepidoptera: Nymphalidae). *Biological Journal of the Linnean Society* doi:10.1111/bij.12091.
956  Merrill RM, Gompert Z, Dembeck LM, Kronforst MR, McMillan WO, Jiggins CD. 2011. Mate
957        Preference Across the Speciation Continuum in a Clade of Mimetic Butterflies. *Evolution* **65**:
958        1489–1500. doi:10.1111/j.1558-5646.2010.01216.x.
959  Merrill RM, Schooten BV, Scott JA, and Jiggins CD. 2011. Pervasive Genetic Associations Between
960        Traits Causing Reproductive Isolation in Heliconius Butterflies. *Proc. Roy Soc B* **278**: 511–518.
961        doi:10.1098/rspb.2010.1493.
962  Müller F. 1879. Ituna and Thyridia; a Remarkable Case of Mimicry in Butterflies. *Trans. Entomol. Soc.*
963        *Lond.* 1879: xx–xxix.
964  Nadeau NJ, Jiggins CD. 2010. A Golden Age for Evolutionary Genetics? Genomic Studies of
965        Adaptation in Natural Populations. *Trends in Genetics* **26**: 484–492. doi:16/j.tig.2010.08.004.
966  Nadeau NJ., Martin SH, Kozak KM, Salazar C, Dasmahapatra KK, Davey JW, Baxter SW, Blaxter ML,
967        Mallet J, Jiggins CD. 2013. Genome-wide Patterns of Divergence and Gene Flow Across a
968        Butterfly Radiation. *Molecular Ecology* **22**: 814–826. doi:10.1111/j.1365-294X.2012.05730.x.
969  Nadeau NJ, Whibley A, Jones RT, Davey JW, Dasmahapatra KK, Baxter SW, Quail MA, Joron M,
970        ffrench-Constant RH, Blaxter, M, et al. 2012. Genomic Islands of Divergence in Hybridizing
971        Heliconius Butterflies Identified by Large-scale Targeted Sequencing. *Phil. Trans. Roy. Soc. B*
972        **367**: 343–353. doi:10.1098/rstb.2011.0198.
973  Orr HA. 2005. The Genetic Theory of Adaptation: a Brief History. *Nat Rev Genet* **6**: 119–127.

974    Papa R, Kapan DD, Counterman BA, Maldonado K, Lindstrom DP, Reed RD, Nijhout HF, Hrbek T,
975        McMillan WO. 2013. Multi-Allelic Major Effect Genes Interact with Minor Effect QTLs to
976        Control Adaptive Color Pattern Variation in *Heliconius erato*. *PLoS ONE* **8**: e57033.
977        doi:10.1371/journal.pone.0057033.
978    Papa R, Martin A, and Reed RD. 2008. Genomic Hotspots of Adaptation in Butterfly Wing Pattern
979        Evolution. *Current Opinion in Genetics & Development* **18**: 559–564.
980        doi:10.1016/j.gde.2008.11.007.
981    Papa R, Morrison CM, Walters JR, Counterman BA, Chen R, Halder G, Ferguson L, Chamberlain N,
982        ffrench-Constant R, Kapan DD et al. 2008. Highly Conserved Gene Order and Numerous
983        Novel Repetitive Elements in Genomic Regions Linked to Wing Pattern Variation in
984        *Heliconius* Butterflies. *BMC Genomics* **9**: 345. doi:10.1186/1471-2164-9-345.
985    Pardo-Diaz C, Salazar C, Baxter SW, Merot C, Figueiredo-Ready W, Joron M, McMillan WO, Jiggins
986        CD. 2012. Adaptive Introgression Across Species Boundaries in Heliconius Butterflies. *PLoS*
987        *Genet* **8**: e1002752. doi:10.1371/journal.pgen.1002752.
988    Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. 2006. Principal Components
989        Analysis Corrects for Stratification in Genome-wide Association Studies. *Nature Genetics* **38**:
990        904–909. doi:10.1038/ng1847.
991    Pritchard JK, Stephens M, Donnelly P. 2000. 'Inference of Population Structure Using Multilocus
992        Genotype Data. *Genetics* **155**: 945–959.
993    Quek S-P, Counterman BA, Albuquerque de Moura P, Cardoso MZ, Marshall C, McMillan WO,
994        Kronforst MR. 2010. Dissecting Comimetic Radiations in Heliconius Reveals Divergent
995        Histories of Convergent Butterflies. *PNAS* **107**: 7365–7370. doi:10.1073/pnas.0911572107.
996    R Development Core Team. 2011. *R: A Language and Environment for  Statistical Computing*
997        (version 2.14). Vienna, Austria: R Foundation for Statistical Computing. http://www.R-
998        project.org/.
999    Reed RD, Papa R, Martin A, Hines HM, Counterman BA, Pardo-Diaz C, Jiggins CD, Chamberlain Nl,
1000        Kronforst MR, Chen R, et al. 2011. Optix Drives the Repeated Convergent Evolution of
1001        Butterfly Wing Pattern Mimicry. *Science* **333**: 1137 –1141. doi:10.1126/science.1208227.
1002    Rieseberg L, Buerkle CA. 2002. Genetic Mapping in Hybrid Zones. *The American Naturalist* **159**: S36–
1003        S50. doi:10.1086/338371.
1004    Rosser N, Phillimore AB, Huertas B, Willmott KR, Mallet J. 2012. Testing Historical Explanations for
1005        Gradients in Species Richness in Heliconiine Butterflies of Tropical America. *Biological*
1006        *Journal of the Linnean Society* **105**: 479–497. doi:10.1111/j.1095-8312.2011.01814.x.
1007    Salazar PA. 2012. Hybridization and the Genetics of Wing Colour-pattern Diversity in Heliconius
1008        Butterflies. PhD, Cambridge, UK: University of Cambridge.
1009    Sheppard PM, Turner JRG, Brown KS, Benson WW, Singer MC. 1985. Genetics and the Evolution of
1010        Muellerian Mimicry in Heliconius Butterflies. *Phil. Trans. Roy. Soc. B* **308**: 433–610.
1011        doi:10.2307/2398716.
1012    Stewart C-B, Schilling JW, Wilson AC. 1987. Adaptive Evolution in the Stomach Lysozymes of Foregut
1013        Fermenters. *Nature* **330** : 401–404. doi:10.1038/330401a0.
1014    Stinchcombe JR, Hoekstra HE. 2007. Combining Population Genomics and Quantitative Genetics:
1015        Finding the Genes Underlying Ecologically Important Traits. *Heredity* **100**: 158–170.
1016    Supple MA, Hines HM, Dasmahapatra KK, Lewis JJ, Nielsen DM, Lavoie C, Ray DA, Salazar C, McMillan
1017        WO, Counterman BA. 2013. Genomic Architecture of Adaptive Color Pattern Divergence and
1018        Convergence in *Heliconius* Butterflies. *Genome Research* **23**: 1248-1257
1019        doi:10.1101/gr.150615.112.
1020    The *Heliconius* Genome Consortium. 2012. Butterfly Genome Reveals Promiscuous Exchange of
1021        Mimicry Adaptations Among Species. *Nature* **487**: 94–98. doi:10.1038/nature11041.
1022    Turner JR, Johnson MS, Eanes WF. 1979. Contrasted Modes of Evolution in the Same Genome:
1023        Allozymes and Adaptive Change in *Heliconius*. *Proceedings of the National Academy of*
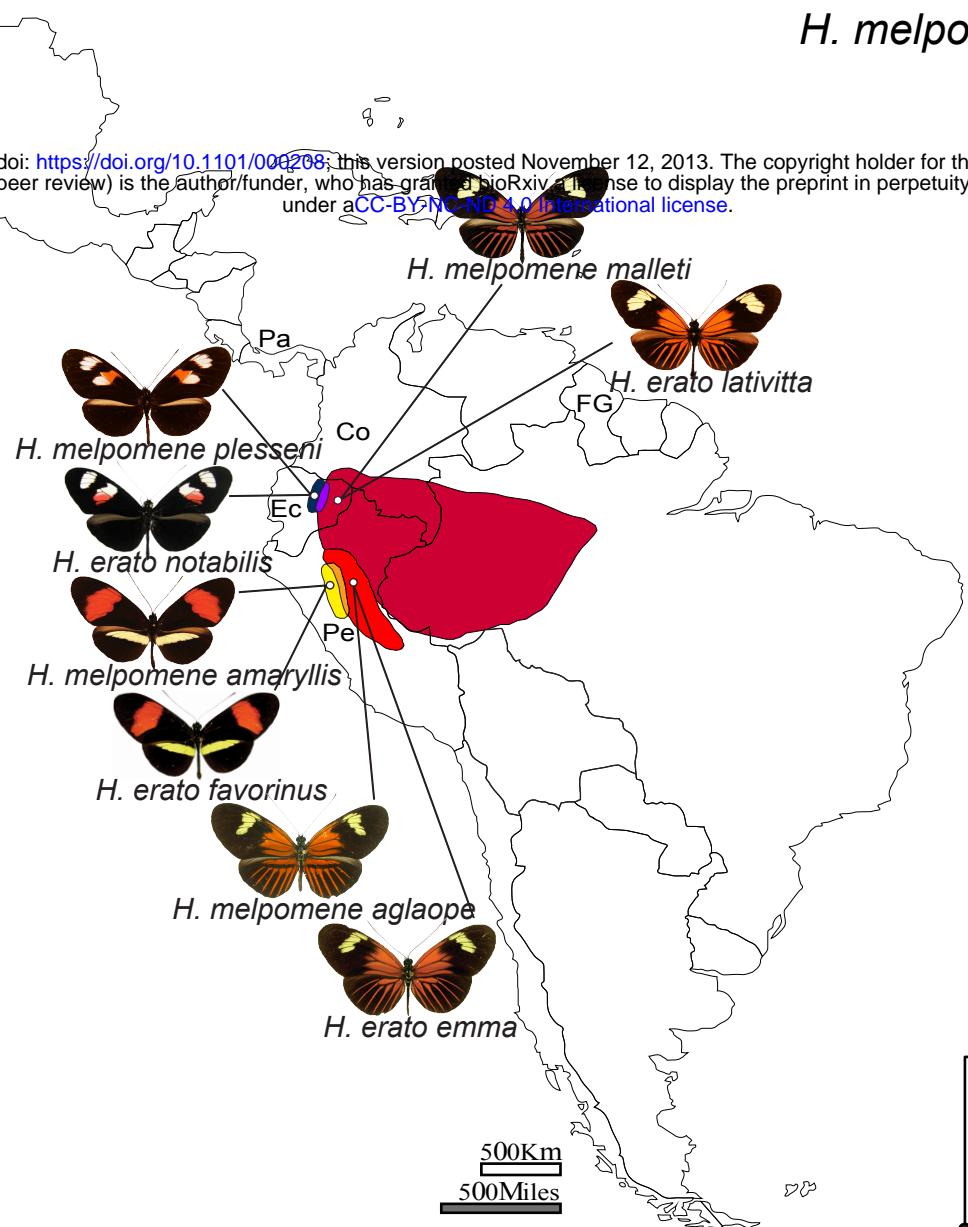1024        *Sciences* **76**: 1924–1928.

1025    Via S. 2012. Divergence Hitchhiking and the Spread of Genomic Isolation During Ecological
1026            Speciation-with-gene-flow. *Phil. Trans. Roy. Soc. B* **367**: 451–460.
1027            doi:10.1098/rstb.2011.0260.
1028    Visscher PM, Brown MA, McCarthy MI, Yang J. 2012. Five Years of GWAS Discovery. *The American*
1029            *Journal of Human Genetics* **90**: 7–24. doi:10.1016/j.ajhg.2011.11.029.
1030    Wood TE, Burke TM, Rieseberg LH. 2005. Parallel Genotypic Adaptation: When Evolution Repeats
1031            Itself. *Genetica* **123**: 157–170. doi:10.1007/s10709-003-2738-9.
1032

A.

*H. melpomene malleti*

*H. erato lativitta*

Pa

Co

*H. melpomene plesseni*

FG

Ec

*H. erato notabilis*

*H. melpomene amaryllis*

Pe

*H. erato favorinus*

*H. melpomene aglaope*

*H. erato emma*

500Km

500Miles

B.

*H. melpomene* clade

*H. erato* clade

*H. timareta*

*H. cydno*

Approx. likelihood branch support

● 1

● <1 ≥0.95

○ <0.95 ≥0.9

0.01

0.01

## A. H. melpomene Peru

## B. Structure K=2

## C. PCA

## D. Cumulative Distribution of between-subspecies $F_{ST}$

**Reference alignment $F_{ST}$ outliers**

*De novo* assembly $F_{ST}$ outliers

**Reference alignment phenotypic associations**

*De novo* assembly phenotypic associations

*H. melpomene* Peru

141
72
12
3
7
1
5
12
2

*H. melpomene* Ecuador

169
1
8
2
2
6

*H. erato* Peru

13
17
5
4
1
5
22

*H. erato* Ecuador

44
12
2
3