



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

Nicola Nadeau, Mayte Ruiz, Patricio Salazar, et al.

Genome Res. published online May 13, 2014

Access the most recent version at doi:[10.1101/gr.169292.113](https://doi.org/10.1101/gr.169292.113)

P<P	Published online May 13, 2014 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see http://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 3.0 Unported), as described at http://creativecommons.org/licenses/by-nc/3.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

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

Nicola J. Nadeau^{1,2}, Mayté Ruiz³, Patricio Salazar^{1,4}, Brian Counterman⁵, Jose Alejandro Medina⁶, Humberto Ortiz-Zuazaga^{6,7}, Anna Morrison¹, W. Owen McMillan⁸, Chris D. Jiggins^{*1}, Riccardo Papa³

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

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

Running title: Population genomics of parallel hybrid zones

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

ABSTRACT

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

INTRODUCTION

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

Butterflies of the Neotropical genus *Heliconius* are extremely diverse in their wing colour patterns and combine within species diversity with convergence among species in wing phenotypes. Their bright wing patterns are used as aposematic warnings to predators, and are under positive frequency dependent selection favouring common colour patterns that predators learn to avoid. This strong selection also maintains narrow hybrid zones between subspecies with different patterns (Benson 1972; Mallet and Barton 1989a; Kapan 2001; Langham 2004). In addition, frequency dependent selection leads to Müllerian mimicry between many distinct species (Müller 1879). For instance, *H. erato* and *H. melpomene* are two distantly related species that diverged around 15-20 million years ago, but have converged on common colour patterns across most of the Neotropics. Divergent races of both species meet in parallel hybrid zones (Figure 1). Evidence suggests that convergent colour patterns in these two species have evolved independently (Supple et al 2013, Hines et al. 2011). It has also been suggested that *H. erato* is more ancient and *H. melpomene* diversified more recently to mimic the *H. erato* forms (Brower 1996; Flanagan et al. 2004; Quek et al. 2010). Nevertheless, it appears that the same handful of genetic loci are responsible for producing most of the colour pattern variation in both species (Joron et al. 2006; Baxter et al. 2008; Martin et al. 2012; Reed et al. 2011). This pattern of parallel adaptive radiation makes *Heliconius* an excellent system in which to address the predictability of the evolutionary process and the extent to which particular genes are re-used when evolving the same phenotypes (Nadeau and Jiggins 2010; Papa, Martin and Reed 2008).

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

Extensive genetic mapping using experimental crosses between different colour pattern forms has identified the chromosomal regions responsible for colour pattern variation (Sheppard et al. 1985; Baxter et al. 2008; Joron et al. 2006; Papa et al. 2013). Three major clusters of loci control most of

the colour pattern variation observed in both species. The tightly linked *B* and *D* loci on chromosome 18 in *H. melpomene* control the red forewing band and the red/orange hindwing rays and proximal “dennis” patches on both wings respectively. These loci are homologous to the *D* locus in *H. erato* (Baxter et al. 2008) and appear to be cis regulatory elements of the *optix* gene (Reed et al. 2011; Supple et al. 2013). The *Ac* and *Sd* loci, in *H. melpomene* and *H. erato* respectively, control the shape of the forewing band via regulation of the *WntA* gene on chromosome 10 (Martin et al. 2012). The presence of most yellow and white elements on the wing is largely controlled by three tightly linked loci, *Yb*, *Sb* and *N*, on chromosome 15 in *H. melpomene* (Ferguson et al. 2010), which are homologous to the *Cr* locus in *H. erato* (Joron et al. 2006). Quantitative trait locus (QTL) mapping has identified other loci of minor effect, including at least 7 additional QTL in *H. erato* (Papa et al. 2013), and QTL in *H. melpomene* on chromosomes 2, 7 and 13 that affect forewing band shape (Baxter et al. 2008). In some cases mapping studies have been followed up by population genetic studies of the mapped intervals across natural hybrid zones, where many generations of backcrossing have led to narrow regions of association, permitting fine scale mapping (Baxter et al. 2010; Counterman et al. 2010; Nadeau et al. 2012; Supple et al. 2013). High-throughput sequencing technologies now provide the feasibility to generate a high density of genomic markers to identify the narrow QTL present in these hybrid zones without the need to perform controlled laboratory crosses (Crawford and Nielsen 2013). Here we test this approach, using a system in which some of the loci responsible for phenotypic differences are known.

The Peru and Ecuador hybrid zones occur across altitudinal gradients (Figure 2A). Therefore, it is possible that traits other than colour pattern may also be differentiated driven by altitudinal selection, for example related to temperature or changes in larval host plants. Such selection on additional regions of the genome could help to stabilise the geographic location of the hybrid zone (Bierne et al. 2011; Barton and Hewitt 1985; Mallet and Barton 1989b; Mallet 2010). Therefore another important question that we will address is whether there are divergent regions of the genome that are not controlling colour pattern. These might be candidates for loci controlling other aspects of ecological adaptation.

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

- 1) if association mapping in these hybrid zones can identify known and novel loci underlying phenotypic variation
- 2) how much of the genome is differentiated and under divergent selection between subspecies
- 3) how much of this differentiation is due to loci controlling colour pattern variation
- 4) if the same regions are divergent between co-mimetic species

While previous studies have touched on questions 2 and 3 (Martin et al. 2013, Kronforst 2013), here we focus on divergence at the subspecies level where hybridisation is frequent, rather than between occasionally hybridising species. Compared to the study by Martin et al. (2013) we explored additional hybrid zones (Ecuador) and species (*H. erato*) using larger sample sizes, which allowed more robust tests to identify genomic regions under divergent selection. We also investigate the

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

RESULTS

Summary of the data and comparison of alignment and assembly techniques

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

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

As expected, given that *H. erato* is ~15% divergent from *H. melpomene* in the aligned data, fewer *H. erato* reads aligned to the *H. melpomene* genome as compared to those from *H. melpomene*, leading to fewer confidently called bases. Nevertheless, the use of the reference *H. melpomene* genome for aligning the *H. erato* reads resulted in more bases being called across multiple individuals and around 10x more SNPs identified when compared to the *de novo* assembly approach. In addition, the gaps between aligned RAD loci were similar across both species (Table 1), indicating that the reduced number of bases is not due to fewer RAD loci aligning but to fewer confidently called bases at each RAD locus. The power to detect loci under selection or responsible for phenotypic variation should therefore be similar in both species or slightly reduced in *H. erato* due to its larger genome (Tobler et al. 2004). Nevertheless, much of the additional genomic sequence in *H. erato* is likely to be

repetitive DNA (Papa et al. 2008), which would be difficult to align and score variants in, even if a complete reference was available. In summary, it seems that the aligned data should give the most power to detect divergent regions and phenotypic associations for both species. However, we performed outlier and association analyses using the output of both approaches for comparison. It is also possible that the *de novo* assembly might detect divergent regions important in adaptation that could not be aligned to the *H. melpomene* reference.

Phylogenetics and population structure

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

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

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

We conducted further analyses of the genetic structure of each of the hybrid zone populations using the reference aligned data, excluding the *H. timareta* individuals. Overall, these results suggest only very low genetic differentiation between any of the parapatric subspecies. STRUCTURE analyses of each population generally showed very little structure and strongest support for only a single population cluster being present. The only exception was the Peruvian *H. melpomene*, where two population clusters gave the highest posterior probability ($p=1$). However, these clusters did not correspond to the two subspecies. The genetic diversity was partitioned such that most individuals were admixed with about a quarter of their allelic variation from population 2, except for two “hybrid” individuals that had pure population 2 genotypes and two other individuals (one “hybrid” and one *aglaope*) that had almost pure population 1 genotypes (Figure 2B). PCA revealed very similar patterns, with small groups of hybrid phenotype individuals giving the clearest clusters, which in most cases were also identified by STRUCTURE (with $K=2$, Figure 2C). Three of the populations did reveal some separation of the subspecies at one of the first two principal components, but with a gradual change from one genomic “type” to another. The *H. melpomene* subspecies in Ecuador were separated by PC1, which explained 10% of the variation in this population. The two *H. erato* populations both showed some separation by subspecies at PC2, which explained 5.7% and 6.7% of the variation in Peru and Ecuador respectively. We found very similar results with PCA on the *de novo* assembled data (SI figure 3), suggesting that the underlying genetic signal in both data sets is very similar. The lack of strong differentiation between subspecies was also supported by the F_{ST} distributions (calculated by BayeScan), which gave very low F_{ST} values between subspecies at over 99% of the genome, with only a small percentage of SNPs showing high levels of differentiation (Figure 2D, SI figure 3).

Association mapping of loci responsible for phenotypic variation

We performed association mapping to identify genetic regions responsible for the phenotypic variation that segregates across each of the hybrid zones. In general, the expected associations were found at the three major loci known to control colour pattern variation on chromosomes 10, 15 and 18 (Figures 3A/D, 4A/D and Table 2, SI table 2). The majority of SNPs showing significant phenotypic associations fell within or tightly linked to these loci in all populations except in Peruvian *H. erato*, where only 26% were tightly linked to the known loci (SI table 2, SI figures 4-8).

Independent analyses were performed on both the reference alignments and *de novo* assemblies of the data. In all populations, more associated SNPs were identified in the alignments than in the *de novo* assemblies (Table 1, figure 5). We used BLASTN (Altschul et al. 1990) to place *de novo* contigs containing associated SNPs onto the *H. melpomene* genome, and most could be confidently assigned to a unique locus. There was almost no overlap in the particular SNPs detected in the assembled and aligned data sets (Figure 5), although in many cases the SNPs detected were in similar regions (Figures 3 and 4). There was evidence for a higher false positive rate in the *de novo* data, as the majority of the SNPs that were uniquely significantly associated in these data were present in the aligned data but did not reach significance. This, rather than detection of novel regions, appears to be the main cause of the higher proportion of associations found scattered across the genome and away from known colour pattern loci in the *de novo* data.

Red colour pattern elements and the B and D loci

Our mapping of red colour pattern variation was generally consistent with previous studies (Baxter et al. 2008; Counterman et al. 2010; The *Heliconius* Genome Consortium 2012) and in almost all populations was mapped to the expected region of chromosome 18 (Table 2, SI table 2). The only exception was *H. melpomene* in Ecuador, where SNPs in this region did not reach significance (Figure 4A). This is likely to be due to the reduced sample size of this population (22) after removing the *H. timareta* individuals. Colour patterns were scored both as independent elements and also using known patterns of segregation to score the predicted genotype at the *B/D* locus. The red forewing band and red hindwing rays are both dominant traits but are controlled by a single locus (or very tightly linked loci) and inherited in repulsion, meaning that individuals with both traits can be inferred to be heterozygotes (Sheppard et al. 1985, Baxter et al. 2008). This genotypic scoring generally gave stronger associations (Figures 3 and 4), although both methods gave some significant associations for at least one of the traits. In all populations the strongest associations in this region were over 60kb downstream of the *optix* gene that controls red colour pattern (Reed et al. 2011), overlapping the region identified in previous analyses as likely containing the functional regulatory variation (Nadeau et al. 2012; Supple et al. 2013)(Table 2, SI figure 4).

In several populations we found additional associations with *B/D* phenotypes on linked chromosome 18 scaffolds (SI table 2). The furthest from the *B/D* locus was HE671488 in Peruvian *H. melpomene*, which is approximately 2Mb away. This scaffold was also associated with differences in altitude in this population, which were stronger than the associations with colour (Figure 3A, Table 3, SI figure 4). This could suggest that this *B/D* linked region is responsible for ecological adaptation, although colour and altitude are strongly correlated so we do not have much power in the data set to separate the two.

Both the Peruvian and Ecuadorian *H. melpomene* populations had a SNP at position 97 on an unmapped scaffold, HE670458, that was highly associated with rays (Table 3, SI table 3). This scaffold appears to consist largely of repetitive elements (BLAST hits match many other regions of the *H. melpomene* genome), suggesting that there may be a copy of a repetitive element that is associated with the presence of rays in both populations. All rayed individuals were heterozygous and all non-rayed individuals were homozygous at this SNP in both *H. melpomene* populations. This would be consistent with multiple alleles aligning to this genomic region but the presence of a unique haplotype sequence linked to the rayed allele. The existence of such a repetitive element is consistent with previous findings that repetitive elements are present in the region of highest divergence at the *B/D* locus (Nadeau et al. 2012; Papa et al. 2008).

Surprisingly, in the Peruvian *H. erato* population the strongest associations with red colour pattern elements were not on chromosome 18, but at two scaffolds (HE670771 and HE670235) on chromosome 2 (Figure 3D, Table 3, SI figure 4, SI table 2). In addition, two SNPs significantly associated with rays and *D* genotype in the *de novo* assembled data of this population could not be confidently assigned to a position in the genome.

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

In the Peruvian *H. melpomene* population the presence of the yellow hindwing bar and yellow in the forewing band both mapped to chromosome 15, with positions that were consistent with previous work on the *Yb* locus (Ferguson et al. 2010; Nadeau et al. 2012; The *Heliconius* Genome Consortium 2012)(Figure 3A, Table 2, SI table 2, SI figure 4). Associations with altitude were also found at these

associated SNPs but were weaker than the association with colour and so may simply be due to correlations between the altitude of the sampling site and colour pattern (SI table 2).

In the Peruvian *H. erato* population we did not recover the expected associations with the yellow hindwing bar, which is known to be controlled by the *Cr* locus on chromosome 15 (Joron et al. 2006; Counterman et al. 2010). Instead, the strongest association with this phenotype in the reference aligned data was found on chromosome 17 (Figure 3D, Table 3, SI table 2). Moreover, we also identified significant associations with the yellow hindwing bar on chromosome 10 in both the reference aligned and *de novo* assembled data. These associations can be explained by the presence of *Sd* on this scaffold (Martin et al. 2012) (Table 2), which is known to influence the expression of the yellow hindwing bar, particularly in individuals that are heterozygous at the *Cr* locus (Mallet 1989). The *Cr* locus is not thought to control any aspects of phenotypic variation in Ecuadorian *H. erato* (Salazar 2012) and consistent with this expectation we did not detect any phenotypic associations in this region.

We scored the Ecuadorian *H. melpomene* individuals for their predicted genotype at the *N* locus (Figure 4C), which controls several wing colour pattern elements such as the yellow forewing band, the amount and location of red in the forewing and the length of the orange hindwing “dennis” bar. Its scoring therefore depends on interactions with the red *B/D* locus (Salazar 2012). Despite the epistatic interaction between the *B/D* and the *N* loci, we could still score them independently. Associations with *N* were found to overlap the *Yb* region (Ferguson et al. 2010; Nadeau et al. 2012) on chromosome 15 (Table 2, SI figure 4), in both the reference aligned and *de novo* assembled data (Figure 4A). Although *N* is known to be tightly linked to *Yb* (Sheppard et al. 1985; Mallet 1989), these are the first genetic mapping results for the *N* locus.

Significant associations with yellow in the forewing band were present in the *D* region in Ecuadorian *H. erato* (Figure 4D, S table 2), consistent with the fact that the *D* locus controls both yellow and red colouration in the forewing band in *H. erato* (Sheppard et al. 1985, Salazar 2012, Papa et al. 2013). No significant associations were found with the presence of yellow colour in the forewing band in either Peruvian *H. erato* or Ecuadorian *H. melpomene*.

Forewing band shape and the Ac, Sd and Ro loci

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

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

In both species from the Ecuadorian hybrid zone, we found SNPs associated with forewing band shape (cell spot 7/8/11) within introns of the *WntA* gene (Figure 4, Table 2, SI table 2). In Ecuadorian *H. erato*, we also found two tightly linked SNPs on chromosome 13 and three tightly linked SNPs on an unmapped scaffold (HE669551) that were associated with forewing band shape and also rounding of the band (Figure 4D, SI table 2). More recent mapping analysis suggests that both these scaffolds are on chromosome 13 and within 1cM of each other (J Davey, Pers. Comm.), so these associations are most likely due to a single locus on this chromosome. Rounding of the distal edge of the band in this population has previously been described as being under the control of the unmapped *Ro* locus (Sheppard et al. 1985; Salazar 2012). We have therefore mapped the *Ro* locus to a region of chromosome 13 (Table 3)

***F_{ST}* outlier detection**

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

As suggested by results from the *de novo* assemblies, there do appear to be differences in population structure between the geographic regions that are consistent across both species. This is also reflected in the *F_{ST}* distributions (from both alignment and assembly approaches), with both *H. erato* and *H. melpomene* having higher mean and background levels of *F_{ST}* in Ecuador as compared to Peru (Table 1, Figure 2, SI figure 3), despite the average distance between sampling locations of “pure” subspecies individuals being similar for both hybrid zones (~56km in Ecuador and 58-60km in Peru). However, the altitudinal range across the hybrid zone in Ecuador is greater than that in Peru (931m versus 318m respectively). Within both regions *H. melpomene* has a lower mean *F_{ST}* than *H. erato*, which would be consistent with higher dispersal distances in *H. melpomene*, as previously suggested (Mallet et al. 1990). Similar outlier regions were detected by both the alignment and assembly approaches (Figures 3 and 4, B and E), although only Peruvian *H. melpomene* gave a good overlap in the specific SNPs detected (Figure 5). Some of the outlier contigs detected in Peruvian *H. erato* could not be positioned on the *H. melpomene* genome with confidence (Figure 3B).

Overall, there was considerable overlap between the genomic regions containing outlier SNPs and those showing phenotypic associations (Figures 3 and 4), and to some extent in the specific SNPs, with the majority of phenotypically associated SNPs also being outliers (Figure 5). The exception to this general trend was the Peruvian *H. erato* population where a large proportion of the phenotypically associated SNPs were not strongly divergent between subspecies. In general, the majority of outlier SNPs were within 1Mb of a known colour pattern locus (including the newly

identified *Ro* region; excluding these, 37.5% of outliers in Ecuadorian *H. erato* were within 1Mb of the *D* and *Sd* loci, SI table 2). The strongest outliers on chromosome 10 in the Ecuadorian populations and Peruvian *H. erato* were within introns of the *WntA* gene and the strongest outliers on the *B/D* scaffold were all 3' of the *optix* gene (Table 2, SI figure 4).

In both *H. melpomene* populations there was a second strongly divergent region on chromosome 18 about 2Mb from the *B/D* region, which was not divergent in either of the *H. erato* populations (Figure 3B, SI figure 4). This is the same region on scaffold HE671488 that showed associations with colour pattern and altitude in the Peruvian *H. melpomene* population (Table 3). In the Peruvian *H. melpomene* population, we detected two clusters of outlier divergent SNPs on chromosome 6, which do not appear to be associated with colour pattern (Figure 3B, Table 3, SI figure 4). Outliers were also detected on chromosome 2 in both *H. erato* populations, some of which were in similar regions to those detected in the association mapping (Table 3, SI figure 4).

DISCUSSION

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

Comparison of de novo assembly and reference alignment of RAD data

Genome-wide association studies (GWAS) are now common in studies of admixed human populations (Visscher et al. 2012). The use of GWAS studies outside of model organisms has mostly been hampered by lack of reference genomes or methods for typing sufficient numbers of markers. However, these limitations are rapidly being eroded as the cost of sequencing decreases and more reference genomes become available. Furthermore, we have shown that alignment of reads to a fairly distantly related reference genome (~15% divergent) can generate meaningful results. In the

absence of a reference genome, *de novo* assembly also detects the same loci, but with somewhat reduced efficacy.

Alignment of sequence reads to the reference genome produced data for more sites, even in the more distantly related species, *H. erato*. One drawback of the Stacks pipeline that we used for *de novo* assembly of the reads is that it does not assemble and call sequence variants in the paired end reads. Hence the available sequence for analysis is almost double in the reference alignments as compared to the *de novo* assembly. However, it also seems that data was lost in the *de novo* assembly due to divergent alleles not being assembled together. This may have had a larger influence on the *H. erato* assemblies as this species harbours greater genetic diversity than *H. melpomene* (Hines et al. 2011) and so explain why a much lower proportion of the *de novo* assembled contigs were present across multiple individuals in *H. erato* (Table 1). We also found a higher proportion of variable sites in the reference alignments as compared to the *de novo* assemblies. This may again be due to poor assembly of the *de novo* contigs, but it could also represent genetic variability contained in the paired-end reads. It is possible that paired end reads might be located in more variable regions, particularly if restriction-site associated reads were biased towards more conserved regions (The *Heliconius* Genome Consortium 2012).

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

Overall, our results suggest that detection of loci underlying adaptive change is likely to be more effective where reads can be mapped to a reference genome. This is perhaps most likely to be the case in populations with high levels of polymorphism, which prevents divergent alleles from assembling. The *de novo* approach could, and no doubt will, be improved by developing methods that allow paired end reads to be incorporated into the SNP typing pipeline (Baxter et al. 2011). This would not only allow a higher density of SNPs to be detected but could also improve alignment of divergent alleles. In the meantime, one approach that has been used in other studies is to first perform *de novo* assembly of RAD-seq reads to generate a consensus reference and to then map reads to this reference for SNP calling (Keller et al. 2013).

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

We have successfully used association mapping in hybrid zone individuals to identify virtually all of the genomic regions known to control colour pattern in these populations (Reed et al. 2011, Martin et al. 2012; Nadeau et al. 2012; Supple et al. 2013). It has commonly been supposed that large sample sizes will be necessary in order identify genes in wild populations. Here we have confirmed

recent theoretical predictions from simulated data (Crawford and Nielsen 2013), that for large effect adaptive loci, even small sample sizes can be highly effective in identification of narrow genomic regions underlying adaptive traits (Table 2, SI figure 4). We also confirm the prediction that in populations with low background levels of divergence both divergence outlier and association mapping approaches are effective in detecting regions under divergent selection. In our study, association mapping has the added benefit of identifying the phenotypic effects of the selected loci. One anticipated pitfall of this method was that many of the phenotypes covary across the hybrid zone. However, it appears that with just 10 individuals with admixed phenotypes we can disassociate most of the variation and thus find distinct genetic associations for known loci. This therefore gives us some confidence that the novel associations that we have detected are real and not due to covariation with other phenotypes.

In Ecuador we intentionally sampled from sites at the edges of the hybrid zone where both pure and hybrid individuals were present, as we anticipated that individuals from these sites would have the highest levels of admixture between selected alleles. This may explain the clearer patterns observed in *H. erato* in Ecuador as compared to Peru (Figures 3D and 4D, SI table 2). In Peruvian *H. erato* we also find several genomic regions showing phenotypic associations that are not divergence outliers, which may suggest that these are false positives. However, the less clear signal in Peruvian *H. erato* could also be due to the reduced sample size in this population (27 individuals). Certainly, the reduced number of Ecuadorian *H. melpomene* individuals (22) seems to have reduced the power of the association mapping (Figure 4B). The lack of any signal at the *Cr* locus in Peruvian *H. erato* is surprising, and may be because the *Cr* associated region is very narrow. There were 789 SNPs present in *H. erato* within the 607kb region that is associated with the yellow hindwing bar in *H. melpomene*. This is only slightly below the genome-wide average for *H. erato* (~1,800 SNPs/Mb) but linkage disequilibrium breaks down rapidly in *H. erato* (Counterterman et al. 2010) and so this may not have been sufficient coverage to identify the *Cr* locus. Nevertheless, contrary to previous suggestions (Kronforst et al. 2013), the density of RAD markers we have obtained was sufficient to identify many narrow divergent genomic regions.

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

The possibility of using hybrid zones for association mapping has long been recognised (Kocher and Sage 1986) but few studies have successfully applied this technique. Studies in younger hybrid zones, for example *Helianthus* sunflowers, have found that linkage disequilibrium between unlinked genomic regions in early generation hybrids can produce spurious associations (Rieseberg and Buerkle 2002). *Heliconius* hybrid zones seem ideal in this regard because they appear to be fairly ancient and close to linkage equilibrium. However, it seems likely that many other suitable systems do exist for this type of approach (Lexer et al. 2006; Crawford and Nielsen 2013). An additional

benefit of the *Heliconius* system is that much of the phenotypic variation is controlled by major effect loci, which can be detected with small sample sizes. Although many adaptive phenotypes appear to involve major effect loci (Orr 2005; Nadeau and Jiggins 2010), in order to move beyond these and detect minor effect loci much larger sample sizes will be required (Beavis 1997). However, by incorporating methods that use a probabilistic framework to infer allele frequencies in low coverage sequencing data (Gompert and Buerkle 2011) it should be feasible to sequence large enough samples for analysis of quantitative traits.

Identification of a novel colour pattern locus

Our association mapping results have robustly identified the *H. erato* *Ro* locus, that controls the shape of the distal edge of the forewing band, as being on chromosome 13 near gene HE669551. This gene has a predicted Gene Ontology (GO) molecular function of microtubule binding and is similar to other insect Radial Spoke Head 3 proteins, which are components of the cilia (Avidor-Reiss et al. 2004). It is therefore not an obvious candidate for control of colour pattern, so may simply be linked to the causative site. Our results are contrary to the suggestion of a recently published QTL study that *Ro* may be linked to *Sd* (Papa et al. 2013). However, that study also identified a major unlinked QTL for forewing band shape, that could not be assigned to a *H. melpomene* chromosome and so may be homologous to the locus we detected here. Furthermore, a QTL for several aspects of forewing band shape and size, including the shape of the distal edge, has previously been identified in *H. melpomene* on chromosome 13 (Baxter et al. 2008). This was located to a fairly broad region but its positioning is consistent with our results for the *Ro* locus in *H. erato*. It therefore seems likely that we have identified a new wing patterning locus that is homologous in *H. melpomene* and *H. erato*.

Ecological selection across the hybrid zones

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

There appear to be multiple dispersed divergent regions on chromosome 2 in *H. erato* (SI figure 4). These could be evidence of divergence hitchhiking, whereby new mutations that cause differential fitness are more likely to be fixed by selection if they arise close to other loci already under divergent

selection. This could lead to clustering of divergently selected loci in the genome (Via 2012; Feder et al. 2012). The same process could also have led to additional loci under divergent selection to arise in linkage with the colour pattern loci. This could explain the second divergent and altitude associated region on chromosome 18 (linked to the *B/D* locus) in *H. melpomene* (Table 3, SI figure 4). One possibility is that this could be the *B/D* linked mate preference locus that has previously been identified (Merrill, Van Schooten, et al. 2011), although it is not clear if the mate preference locus is an additional linked locus or a pleiotropic effect of the wing colour locus itself. It is also possible that these apparently distinct but linked divergent regions could simply reflect the heterogeneous nature of F_{ST} resulting from divergent selection on a single locus combined with other background and neutral processes (Charlesworth et al. 1997). A broad region of divergence around the *B/D* locus in *H. melpomene* would fit with other suggestions that it has undergone stronger or more recent selection than other colour pattern loci (Nadeau et al. 2012). The *D* region in *H. erato* does not appear to be extended in the same way as in *H. melpomene* (SI figure 4), suggesting that either the architecture or the selective history of this region is different between these species.

Comparison of the genomic architecture of divergence between convergent species

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

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

In general, it seems that although the same colour pattern loci are present in both species (Joron et al. 2006; Baxter et al. 2008; Martin et al. 2012) they are being used in different ways and

combinations in order to produce convergent phenotypes. This is particularly surprising given the pattern of co-divergence observed in the phylogeny, which would appear to suggest that similar colour patterns have arisen at a similar time and from similar ancestral forms in both species. Nonetheless, the apparent pattern of co-divergence could simply reflect more recent patterns of gene flow between geographically proximate populations in both species. This has recently been highlighted by studies showing that patterns of divergence at colour pattern controlling loci can be very different to those found at the rest of the genome (Hines et al. 2011; Pardo-Diaz et al. 2012; The *Heliconius* Genome Consortium 2012; Supple et al. 2013). Therefore, the differences that we observe in the use of particular loci in the two species could reflect different mimetic histories that will only be resolved by studies of the evolutionary history of particular loci.

Discovery of a new cryptic *H. timareta* population

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

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

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

Conclusions

We have demonstrated that high resolution genome scans using admixed individuals from hybrid zones can be used to identify loci underlying phenotypic variation. Only a small proportion of the genome (about 0.025%) is strongly differentiated between subspecies and most of this can be explained by divergence at loci controlling colour pattern. This is consistent with previous studies

based on smaller numbers of markers (Turner 1979; Baxter et al. 2010, Counterman et al. 2010, Nadeau et al. 2012) and suggests that the hybrid zones are ancient or have formed in primary contact, and are maintained by strong selection on colour pattern (Mallet and Barton 1989a, Mallet 2010). However, we also find, for the first time, some divergent loci that do not appear to be associated with colour pattern, suggesting that there may be other differences between subspecies. This could explain why several *Heliconius* hybrid zones occur across ecological gradients (Benson 1982), if they are coupled with extrinsic selection acting on other loci in the genome (Bierne et al. 2011). However, this needs to be confirmed with detailed phenotypic analyses of the subspecies to identify whether differences are present that could be explained by ecological adaptation. In general we find that, although some loci are divergent in all populations, the genomic pattern of divergence between co-mimetic species is not particularly similar, suggesting that the level of parallel genetic evolution between *H. erato* and *H. melpomene* is in fact quite low, despite parallel phylogenetic patterns of divergence. Finally, our analysis shows that alignment to a distantly related reference genome can improve analyses over a *de novo* assembly of the data.

METHODS

Samples and sequencing

30 *H. erato* and 30 *H. melpomene* individuals were selected from a larger sample taken from the hybrid zone region in Peru. Similarly, 30 *H. erato* and 30 *H. melpomene* were also selected from a larger study of a subspecies hybrid zone in Ecuador (Salazar 2012). Each set of 30 samples comprised 10 pure forms of each subspecies and 10 hybrids (based on colour pattern). See Figure 2 and SI table 4 for further details of the samples and locations.

RAD sequencing libraries were prepared using previously described methodologies (The *Heliconius* Genome Consortium 2012; Baird et al. 2008; Baxter et al. 2011). Briefly, DNA was digested with the restriction enzyme PstI prior to ligation of P1 sequencing adaptors with five-base molecular identifiers (MIDs, SI table 4). We then pooled samples into groups of 6 before shearing, ligation of P2 adaptors, amplification and fragment size selection (300-600bp). Libraries were then further pooled such that 30 individuals were sequenced on each lane of an Illumina HiSeq 2000 sequencer to obtain 150 base paired-end sequences. We obtained an average of 374M sequence pairs from each lane. Following sequencing, three of the *H. erato* individuals from Peru were found to have been incorrectly assigned to this species and were excluded from all further analyses.

In order to compare patterns of phylogenetic divergence of the focal subspecies, we also used sequence data from additional subspecies and closely related species in each group. Two individuals each from 6 additional *H. erato* populations and the closely related *H. himera* were also PstI RAD sequenced with 5 individuals pooled per lane of Illumina GAIIx (100 base paired-end sequencing). These sequences were obtained in the same run as a comparable set of individuals from the *H. melpomene* clade, which have been used in previous analyses and also included *H. cydno*, *H. timareta* and *H. hecale* (The *Heliconius* Genome Consortium 2012; Nadeau et al. 2013, European Nucleotide Archive, Accession ERP000991). We also obtained whole-genome shotgun sequence data from an outgroup species, *H. clysonimus*, which was sequenced on a fifth of a HiSeq 2000 lane, giving 53.5M 100 base read pairs for this individual.

Alignment to reference genome

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

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

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

De novo assembly

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

Phylogenetics and analysis of population structure

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

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

We also performed principal components analysis of the genetic variation in each population group. This was done with the “`cmdscale`” command in R (R Development Core Team 2011), using genetic distance matrices calculated as 0.5-ibs, where ibs was the identity by sequence matrix calculated in GenABEL (see below). As further confirmation that some of the *H. melpomene* individuals sampled in Ecuador were in fact cryptic *H. timareta*, we also performed principal components analysis on the

combined *H. melpomene* and outgroup data set. We also ran principal components analysis on the *de novo* assembled data for each population group, to test whether both methods were detecting similar underlying patterns of genetic variation.

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

Association mapping of loci controlling colour pattern variation

We scored components of phenotypic variation that segregate across each of the hybrid zones. The scored phenotypes are shown in Figure 3 (for Peru) and Figure 4 (for Ecuador) and listed in full in SI table 5. These were scored mostly as binomial (1,0) traits, but in some cases intermediates were also scored (as 0.5). The width and shape of the forewing band was scored based on whether it extended into each of the wing “cells”, demarcated by the major wing veins (as shown in SI figure 9). In Peruvian populations, the size and shape of the forewing band was measured as two components (Figure 3C/F) that extend the band distally (cell spot 8) and proximally (cell spot 11). In Ecuador, three aspects of band shape were scored: cells 8 and 11, which make up the proximal spot in *H. m. plesseni* and *H. e. notabilis*, and cell 7, which pushes the band towards the wing margin in *H. m. malleti* and *H. e. lativitta* (Figure 4C/F). In our sample of *H. melpomene* the presence of cell spots 8 and 11 were perfectly correlated, whereas in *H. erato* the presence of cell spot 7 was perfectly correlated with the absence of cell spot 8. In addition, individuals were also scored for their predicted genotypes at major loci described previously (with predicted heterozygotes scored as 0.5) (Sheppard et al. 1985; J. Mallet 1989) and the altitude at which they were collected was included as a continuous phenotypic trait.

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

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

We therefore performed score tests for all scored phenotypes across all population groups. Genome-wide significance was determined empirically from 1000 resampling replicates and corrected for population structure using the test specific λ (SI table 5).

BayeScan analysis to identify loci under selection

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

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

DATA ACCESS

DNA sequence reads from this study have been submitted to the European Nucleotide Archive (ENA; <http://www.ebi.ac.uk/ena/>) under accession number ERP003980. COI sequences have been submitted to the EMBL Nucleotide Sequence Database (EMBL-Bank) at ENA under accession numbers HG710096 - HG710125. Custom scripts and wing images are available from Data Dryad with DOI: doi:10.5061/dryad.1nc50.

ACKNOWLEDGEMENTS

We thank Simon Baxter, Doug Turnbull and William Cresko for their help and advice with RAD library preparation and sequencing. Sequencing was performed at the University of Oregon, Genomics Core Facility and The Gene Pool genomics facility in the University of Edinburgh. We would also like to thank Julian Catchen for his help with Stacks. We thank the governments of Peru and Ecuador for their permission to collect and export specimens. Santiago Villamarín from the Museo Ecuatoriano de Ciencias Naturales provided institutional support in Ecuador. We also thank Ismael Aldás, Carlos Robalino and Patricia Salazar for their assistance with fieldwork. Joanna Riley assisted with DNA extractions. John Davey gave us access to his unpublished mapping results. We thank 3 anonymous reviewers for their comments. This project was supported by research grants BBSRC H01439X/1, NSF-CREST #0206200, NSF-DEB-1257839 and NSF-IOS 1305686. NJN was funded by a Leverhulme Trust award to CDJ, MR was funded through the Ford Foundation Postdoctoral Fellowship Program administered by the National Academies, HOZ and JAM were partially supported by NIH-NIGMS INBRE award P20GM103475 and NSF-EPSCoR award 1002410.

FIGURE LEGENDS

Figure 1. A) Distribution in South America of the subspecies included in this study. B) Maximum likelihood phylogenies with approximate likelihood branch supports. Co-mimics from outside of the focal hybrid zones are connected with dotted lines. Focal hybrid zone individuals are shown in

colour: blue, *H. m. plesseni* and *H. e. notabilis*; purple, Ecuador hybrids; dark red, *H. m. malleti* and *H. e. lativitta*; red, *H. m. aglaope* and *H. e. emma*; orange, Peru hybrids; yellow, *H. m. amaryllis* and *H. e. favorinus*. Additional populations are in black. Country abbreviations: Ec, Ecuador; FG, French Guiana; Co, Colombia; Pa, Panama.

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

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

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

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

TABLES

Table 1. Summary statistics from alignment and assembly approaches

<i>De novo</i> assembly with Stacks - single end reads										
		n	millions of reads (mean \pm SD)	bases covered ($\times 10^6$)	bases covered in ≥ 10 inds ($\times 10^6$)	SNPs used in outlier analysis ($\times 10^3$)	mean F_{ST}	Outliers	Significant Phenotypic Associations	
<i>H. erato</i>	Peru	27	8.0 \pm 2.2	166	2.8	37	0.0280	22	10	
	Ecuador	30	9.2 \pm 2.5	149	3.3	31	0.0568	0	2	
<i>H. melpomene</i>	Peru	30	7.0 \pm 2.4	61	4.3	57	0.0145	23	8	
	Ecuador	22	7.8 \pm 1.4	45	3.5	43	0.0310	5	4	

Aligned to <i>H. melpomene</i> reference - paired end with Stampy											
		n	millions of reads (mean \pm SD)	bases covered ($\times 10^6$)	bases covered in ≥ 10 inds ($\times 10^6$)	SNPs used in outlier analysis ($\times 10^3$)	mean F_{ST}	Outliers	Significant Phenotypic Associations	mean gap between RAD loci (kb)	max gap between RAD loci (kb)
<i>H. erato</i>	Peru	27	11.3 \pm 3.2	11	4.2	373	0.0142	19	28	9.4	116
	Ecuador	30	11.9 \pm 3.2	13	5.1	337	0.0316	56	15	9.1	105
<i>H. melpomene</i>	Peru	30	10.9 \pm 3.8	28	14.4	860	0.0112	235	91	9.3	103
	Ecuador	22	10.7 \pm 1.9	23	15.6	788	0.0299	179	14	9.5	114

Table 2. Accuracy of identification of genomic regions known to control colour pattern variation.

Colour pattern loci		<i>B/D</i>	<i>Yb/N/Cr</i>	<i>Ac/Sd</i>
Chromosome		chr18	chr15	chr10
Scaffold		HE670865	HE667780	HE668478
Gene		HMEL001028 (<i>optix</i>) ¹	Presently unknown	HMEL018100 (<i>WntA</i>) ²
- Position		438,423-439,107		450,400-483,854
Functional region ³		300,000-400,000	600,000-1,000,000	Presently unknown
<i>H. melpomene</i>	Peru	assoc	161,328-376,651	676,543-697,543
		outlier	263,358	676,645
	Ecuador	assoc	none	697,118-725,562
		outlier	376,651	697,118
<i>H. erato</i>	Peru	assoc	362,793-362,794	none
		outlier	362,794	none
	Ecuador	assoc	282,473-376,342	N/A
		outlier	376,250	479,220

For each population, positions are given for the SNPs showing the strongest phenotypic associations (assoc) and the highest F_{ST} outliers (outlier) on the given scaffold: N/A = not expected or found; none = not found. ¹ from Reed et al. 2011; ² from Martin et al. 2012; ³ Inferred from population genomics: the *B/D* region appears to be similar in *H. erato* and *H. melpomene*; *Yb/N/Cr* region has been localised in *H. melpomene* only (Nadeau et al. 2012; Supple et al. 2013).

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

Chrom.	Scaffold	Comparison*	Closest Gene	Distance†	GO function	putative protein
chr18	HE671488	<i>melp Peru</i> : assoc alt (D gen); outlier <i>melp Ecuador</i> : outlier	MEL014920	19,171		
unmapped	HE670458	<i>melp Ecuador</i> : assoc rays <i>melp Peru</i> : assoc rays, D gen (alt)	no genes on this scaffold, in repetitive region			
chr2	HE670771	<i>erato Peru</i> : assoc D gen (alt, rays, spot 11) <i>erato Ecuador</i> : assoc spot 11	HMEL008318	0 (I)	catalytic activity, protein binding	fatty acid synthase
chr2	HE670771	<i>erato Peru</i> : assoc alt, rays; outlier	HMEL008322	0 (A)	odorant binding	odorant binding protein 7
chr2	HE670519	<i>erato Peru</i> : assoc spot 11; outlier	HMEL007059	0 (I/A)	oxidoreductase activity	3-dehydroecdysone 3alpha-reductase
chr2	HE670235	<i>erato Peru</i> : assoc D gen (alt, rays, spot 11)	HMEL005708	56,981	taste receptor activity	olfactory receptor 4
chr2	HE671428	<i>erato Ecuador</i> : outlier	HMEL014154	0 (S)	choline dehydrogenase activity, oxidoreductase activity, acting on CH-OH group of donors, flavin adenine dinucleotide binding	glucose dehydrogenase
chr2	HE671428	<i>erato Ecuador</i> : outlier	HMEL014163	0 (A)		heat shock protein 70
chr17	HE671853	<i>erato Ecuador</i> : assoc HWY	HMEL014236	0 (I)	catalytic activity, serine-type endopeptidase activity	serine protease 30
unmapped (chr20)	HE671554	<i>melp Peru</i> : assoc spot 8	HMEL016146	0 (A/S/I)	protein binding, zinc ion binding	MICAL-like
unmapped (chr13)	HE669551	<i>erato Ecuador</i> : assoc Ro (spot 7/8); outlier	HMEL004352	0 (S)	microtubule binding	radial spoke head 3
chr13	HE670984	<i>erato Ecuador</i> : assoc spot 11 (Ro, spot 7/8); outlier	HMEL009926	3,915	structural constituent of ribosome, RNA binding	ribosomal protein S4
chr6	HE671933	<i>melp Peru</i> : outlier	HMEL016074	7,925	oxidoreductase activity	amine oxidoreductase
chr6	HE671934	<i>melp Peru</i> : outlier	HMEL016075	4,121	oxidoreductase activity	amine oxidoreductase

*Analysis in which SNP is detected: *melp*, *H. melpomene*; *erato*, *H. erato*; outlier, BayeScan F_{ST} outlier analysis; assoc, association analysis with the strongest associated phenotype and additional phenotypes in parentheses (rays, presence of hindwing rays and fore/hindwing dennis patches; Dgen, predicted *B/D* genotype; spot, presence of non-black colour in that wing cell; alt, altitude; Ro,

rounding of distal edge of forewing band). † If a SNP is within a gene (distance=0) then in parentheses: A, non-synonymous; S, synonymous; I, within an intron. Further information is given in SI table 3.

REFERENCES

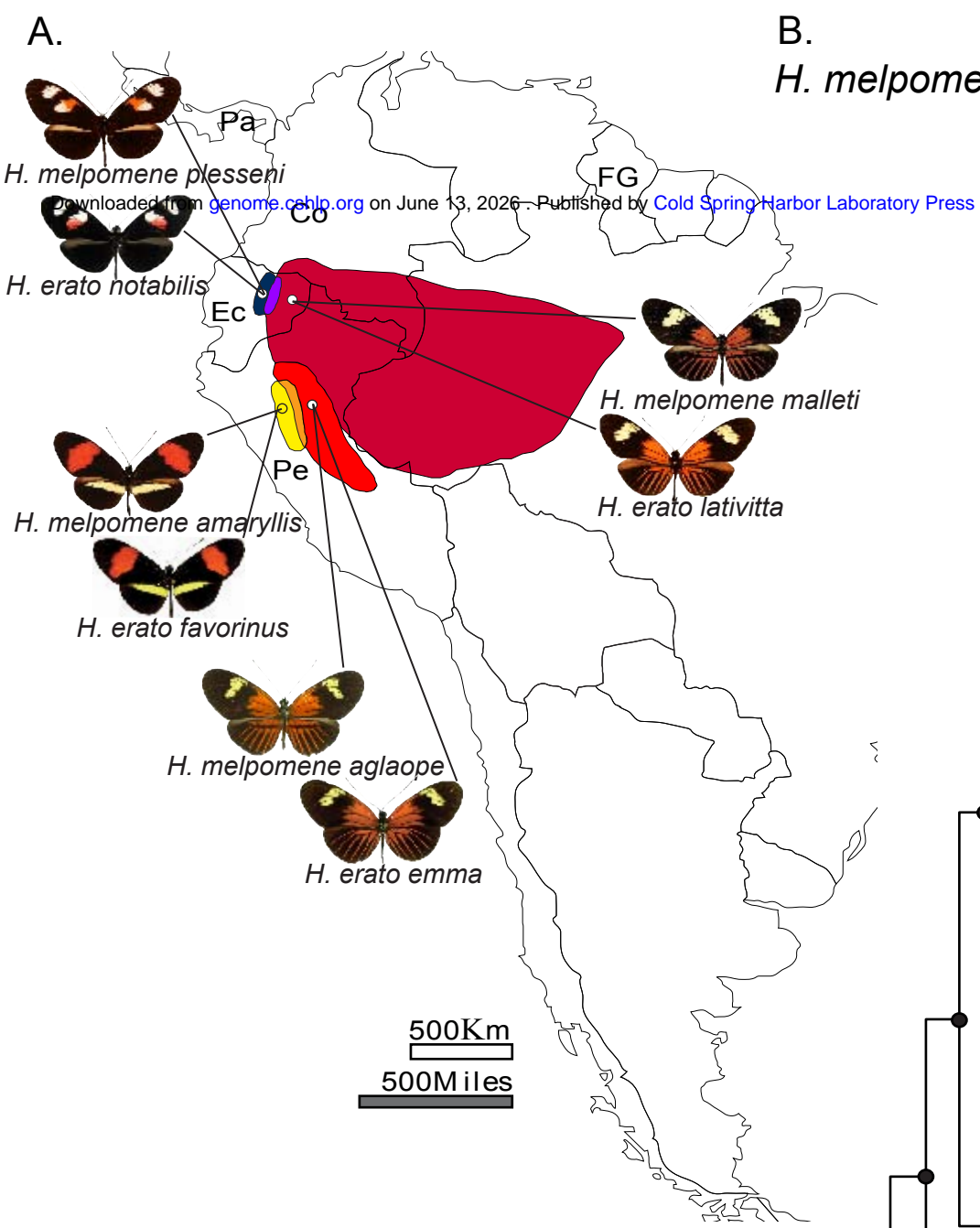
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic Local Alignment Search Tool. *Journal of Molecular Biology* **215**: 403–410. doi:10.1016/S0022-2836(05)80360-2.
- Aulchenko YS, Ripke S, Isaacs A, van Duijn CM. 2007. GenABEL: An R Library for Genome-wide Association Analysis. *Bioinformatics*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

- The Population Genetics of Müllerian Mimicry in *Heliconius erato*. *PLoS Genetics* **6**: e1000796. doi:10.1371/journal.pgen.1000796.
- Crawford JE, Nielsen R. 2013. Detecting Adaptive Trait Loci in Non-model Systems: Divergence or Admixture Mapping? *Molecular Ecology* in press. doi:10.1111/mec.12562.
- Cuthill JH, Charleston M. 2012. Phylogenetic Codivergence Supports Coevolution of Mimetic *Heliconius* Butterflies. *PLoS One* **7**: e36464. doi:10.1371/journal.pone.0036464.
- DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, del Angel G, Rivas MA, Hanna M, et al. 2011. A Framework for Variation Discovery and Genotyping Using Next-generation DNA Sequencing Data. *Nat Genet* **43**: 491–498. doi:10.1038/ng.806.
- Feder JL, Gejji R, Yeaman S, Nosil P. 2012. Establishment of New Mutations Under Divergence and Genome Hitchhiking. *Phil. Trans. Roy. Soc. B* **367**: 461–474. doi:10.1098/rstb.2011.0256.
- Ferguson L, Lee SF, Chamberlain N, Nadeau NJ, Joron M, Baxter S, Wilkinson P, Papanicolaou A, Kumar S, Kee T-J, et al. 2010. Characterization of a Hotspot for Mimicry: Assembly of a Butterfly Wing Transcriptome to Genomic Sequence at the *HmYb/Sb* Locus. *Mol. Ecol.* **19**: 240–254. doi:10.1111/j.1365-294X.2009.04475.x.
- Flanagan NS, Tobler, Davison AA, Pybus OG, Kapan DD, Planas S, Linares M, Heckel D, McMillan WO. 2004. Historical Demography of Müllerian Mimicry in the Neotropical *Heliconius* Butterflies. *PNAS* **101**: 9704–9709. doi:10.1073/pnas.0306243101.
- Foll M, Gaggiotti O. 2008. A Genome-Scan Method to Identify Selected Loci Appropriate for Both Dominant and Codominant Markers: A Bayesian Perspective. *Genetics* **180**: 977–993. doi:10.1534/genetics.108.092221.
- Giraldo N., Salazar C, Jiggins C, Bermingham E, and Linares M. 2008. Two Sisters in the Same Dress: *Heliconius* Cryptic Species. *BMC Evolutionary Biology* **8**: 324. doi:10.1186/1471-2148-8-324.
- Gompert Z, Buerkle CA. 2011. A Hierarchical Bayesian Model for Next-Generation Population Genomics. *Genetics* **187**: 903–917. doi: 10.1534/genetics.110.124693.
- Gompert, Z, Lucas LK, Nice CC, Fordyce JA, Forister ML, Buerkle AC. 2012. Genomic regions with a history of divergent selection affect fitness of hybrids between two butterfly species. *Evolution* **66**: 2167–2181. doi:10.1111/j.1558-5646.2012.01587.x.
- Guindon S, Gascuel O. 2003. A Simple, Fast, and Accurate Algorithm to Estimate Large Phylogenies by Maximum Likelihood. *Systematic Biology* **52**: 696–704. doi:10.1080/10635150390235520.
- Harrison RG. 1993. *Hybrid Zones and the Evolutionary Process*. Oxford University Press.
- Hines HM, Counterman BA, Papa R, Albuquerque de Moura P, Cardoso MZ, Linares M, Mallet J, Reed RD, Jiggins CD, Kronforst MR, et al. 2011. Wing Patterning Gene Redefines the Mimetic History of *Heliconius* Butterflies. *PNAS* **108**: 19666–19671. doi:10.1073/pnas.1110096108.
- Jiggins, CD. 2008. Ecological Speciation in Mimetic Butterflies. *BioScience* **58**: 541–548.
- Joron M, Papa R, Beltrán M, Chamberlain N, Mavárez J, Baxter S, Abanto M, Bermingham E, Humphray SJ, Rogers J, et al. 2006. A Conserved Supergene Locus Controls Colour Pattern Diversity in *Heliconius* Butterflies. *PLoS Biology* **4**: e303. doi:10.1371/journal.pbio.0040303.
- Kapan DD. 2001. Three-butterfly System Provides a Field Test of Mullerian Mimicry. *Nature* **409**: 338–340. doi:10.1038/35053066.
- Kawakami T, Butlin RK. 2001. Hybrid Zones. *eLS*. John Wiley & Sons, Ltd. <http://onlinelibrary.wiley.com/doi/10.1002/9780470015902.a0001752.pub2/abstract>.
- Keller I, Wagner CE, Greuter L, Mwaiko S, Selz OM, Sivasundar A, Wittwer S, Seehausen O. 2013. Population Genomic Signatures of Divergent Adaptation, Gene Flow and Hybrid Speciation in the Rapid Radiation of Lake Victoria Cichlid Fishes. *Molecular Ecology* **22**: 2848–2863. doi:10.1111/mec.12083.
- Kocher TD, Sage RD. 1986. Further Genetic Analyses of a Hybrid Zone Between Leopard Frogs (*Rana pipiens* Complex) in Central Texas. *Evolution* **40**: 21–33. doi:10.2307/2408600.
- Kronforst MR, Hansen MEB, Crawford NG, Gallant JR, Zhang W, Kulathinal RJ, Kapan DD, Mullen SP. 2013. Hybridization Reveals the Evolving Genomic Architecture of Speciation. *Cell Reports* **5**: 666–677. doi:10.1016/j.celrep.2013.09.042.

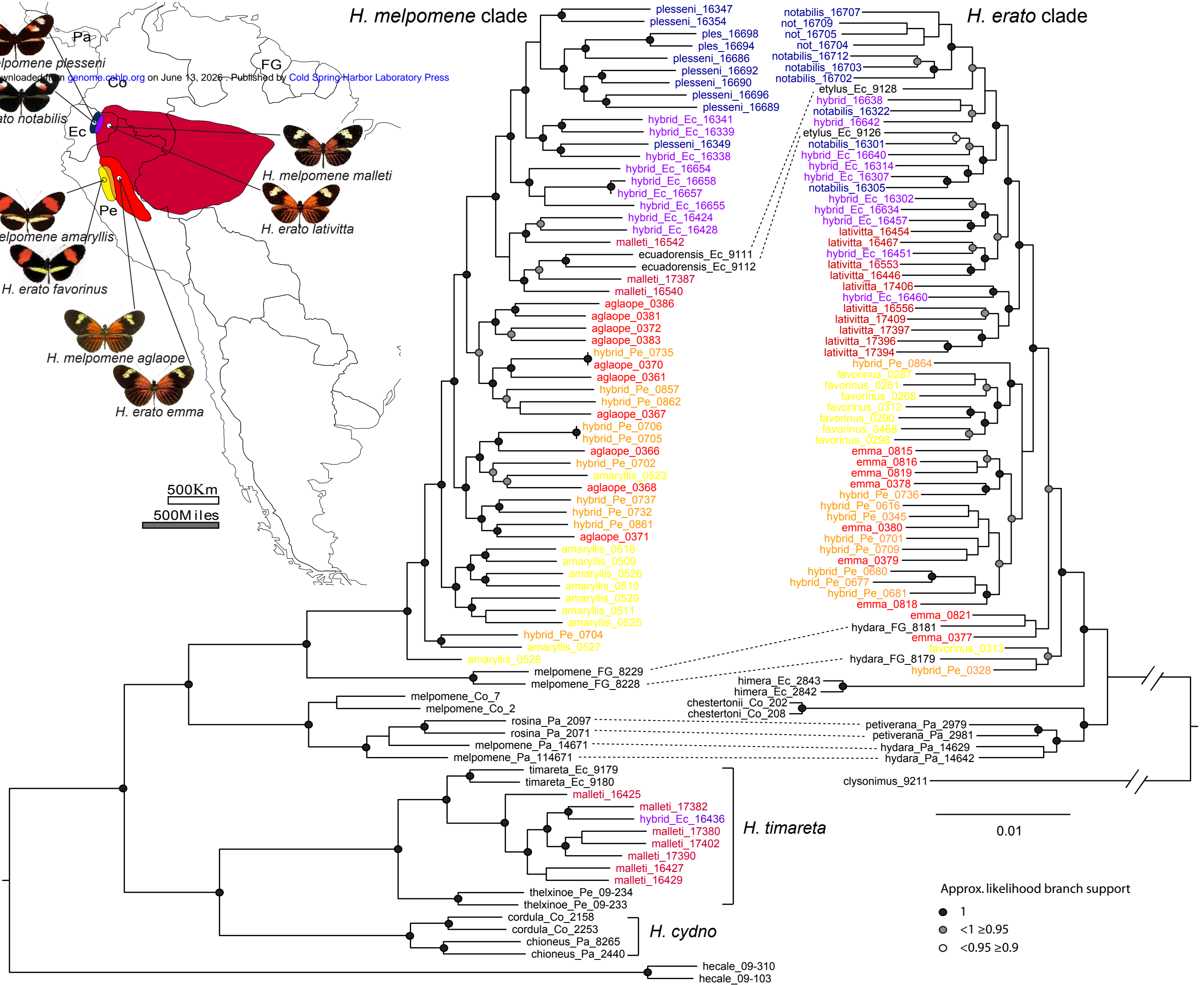
- Lamas G. 1997. Comentarios Taxonomicos y Nomenclaturales Sobre Heliconiini Neotropicales, Con Designacion de Lectotipos y Descripcion de Cuatro Subespecies Nuevas (Lepidoptera: Nymphalidae: Heliconiinae). *Rev. Per. Ent.* **40**: 111–125.
- Langham GM. 2004. Specialized Avian Predators Repeatedly Attack Novel Color Morphs of Heliconius Butterflies. *Evolution* **58**: 2783–2787. doi:10.1111/j.0014-3820.2004.tb01629.x.
- Lexer C, Buerkle CA, Joseph JA, Heinze B, Fay MF. 2006. Admixture in European *Populus* Hybrid Zones Makes Feasible the Mapping of Loci That Contribute to Reproductive Isolation and Trait Differences. *Heredity* **98**: 74–84. doi:10.1038/sj.hdy.6800898.
- Lunter G, Goodson M. 2011. Stampy: A Statistical Algorithm for Sensitive and Fast Mapping of Illumina Sequence Reads. *Genome Research* **21**: 936–939. doi:10.1101/gr.111120.110.
- Mallet J. 1989. The Genetics of Warning Colour in Peruvian Hybrid Zones of *Heliconius erato* and *H. melpomene*. *Proc. Roy. Soc. B* **236**: 163–185.
- Mallet J. 1999. Causes and Consequences of a Lack of Coevolution in Müllerian Mimicry. *Evolutionary Ecology* **13**: 777–806. doi:10.1023/A:1011060330515.
- Mallet J. 2010. Shift Happens! Shifting Balance and the Evolution of Diversity in Warning Colour and Mimicry. *Ecological Entomology* **35**: 90–104. doi:10.1111/j.1365-2311.2009.01137.x.
- Mallet J, Barton N. 1989a. Strong Natural Selection in a Warning-Color Hybrid Zone. *Evolution* **43**: 421–431. doi:10.2307/2409217.
- Mallet J, Barton N. 1989b. Inference from Clines Stabilized by Frequency-dependent Selection. *Genetics* **122**: 967–976.
- Mallet, J, Barton N, Lamas G, Santisteban J, Muedas M, and Eeley H. 1990. Estimates of Selection and Gene Flow from Measures of Cline Width and Linkage Disequilibrium in *Heliconius* Hybrid Zones. *Genetics* **124**: 921–936.
- Martin, A, Papa R, Nadeau NJ, Hill RI, Counterman BA, Halder G, Jiggins CD, Kronforst MR, Long AD, McMillan WO, et al. 2012. Diversification of Complex Butterfly Wing Patterns by Repeated Regulatory Evolution of a Wnt Ligand. *PNAS* **109**: 12632–12637. doi:10.1073/pnas.1204800109.
- Martin SH, Dasmahapatra KK, Nadeau NJ, Salazar C, Walters JR, Simpson F, Blaxter M, Manica A, James Mallet, Jiggins CD. 2013. Genome-wide Evidence for Speciation with Gene Flow in *Heliconius* Butterflies. *Genome Research* **23**: 1817–1828. doi: 10.1101/gr.159426.113.
- Mérot C., Mavárez J, Evin A, Dasmahapatra KK, Mallet J, Lamas G, and Joron M. 2013. Genetic Differentiation Without Mimicry Shift in a Pair of Hybridizing *Heliconius* Species (Lepidoptera: Nymphalidae). *Biological Journal of the Linnean Society* doi:10.1111/bij.12091.
- Merrill RM, Gompert Z, Dembeck LM, Kronforst MR, McMillan WO, Jiggins CD. 2011. Mate Preference Across the Speciation Continuum in a Clade of Mimetic Butterflies. *Evolution* **65**: 1489–1500. doi:10.1111/j.1558-5646.2010.01216.x.
- Merrill RM, Schooten BV, Scott JA, and Jiggins CD. 2011. Pervasive Genetic Associations Between Traits Causing Reproductive Isolation in *Heliconius* Butterflies. *Proc. Roy Soc B* **278**: 511–518. doi:10.1098/rspb.2010.1493.
- Müller F. 1879. Ituna and Thyridia; a Remarkable Case of Mimicry in Butterflies. *Trans. Entomol. Soc. Lond.* 1879: xx–xxix.
- Nadeau NJ, Jiggins CD. 2010. A Golden Age for Evolutionary Genetics? Genomic Studies of Adaptation in Natural Populations. *Trends in Genetics* **26**: 484–492. doi:16/j.tig.2010.08.004.
- Nadeau NJ., Martin SH, Kozak KM, Salazar C, Dasmahapatra KK, Davey JW, Baxter SW, Blaxter ML, Mallet J, Jiggins CD. 2013. Genome-wide Patterns of Divergence and Gene Flow Across a Butterfly Radiation. *Molecular Ecology* **22**: 814–826. doi:10.1111/j.1365-294X.2012.05730.x.
- Nadeau NJ, Whibley A, Jones RT, Davey JW, Dasmahapatra KK, Baxter SW, Quail MA, Joron M, French-Constant RH, Blaxter, M, et al. 2012. Genomic Islands of Divergence in Hybridizing *Heliconius* Butterflies Identified by Large-scale Targeted Sequencing. *Phil. Trans. Roy. Soc. B* **367**: 343–353. doi:10.1098/rstb.2011.0198.
- Orr HA. 2005. The Genetic Theory of Adaptation: a Brief History. *Nat Rev Genet* **6**: 119–127.

- Papa R, Kapan DD, Counterman BA, Maldonado K, Lindstrom DP, Reed RD, Nijhout HF, Hrbek T, McMillan WO. 2013. Multi-Allelic Major Effect Genes Interact with Minor Effect QTLs to Control Adaptive Color Pattern Variation in *Heliconius erato*. *PLoS ONE* **8**: e57033. doi:10.1371/journal.pone.0057033.
- Papa R, Martin A, and Reed RD. 2008. Genomic Hotspots of Adaptation in Butterfly Wing Pattern Evolution. *Current Opinion in Genetics & Development* **18**: 559–564. doi:10.1016/j.jgde.2008.11.007.
- Papa R, Morrison CM, Walters JR, Counterman BA, Chen R, Halder G, Ferguson L, Chamberlain N, French-Constant R, Kapan DD et al. 2008. Highly Conserved Gene Order and Numerous Novel Repetitive Elements in Genomic Regions Linked to Wing Pattern Variation in *Heliconius* Butterflies. *BMC Genomics* **9**: 345. doi:10.1186/1471-2164-9-345.
- Pardo-Diaz C, Salazar C, Baxter SW, Merot C, Figueiredo-Ready W, Joron M, McMillan WO, Jiggins CD. 2012. Adaptive Introgression Across Species Boundaries in *Heliconius* Butterflies. *PLoS Genet* **8**: e1002752. doi:10.1371/journal.pgen.1002752.
- Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. 2006. Principal Components Analysis Corrects for Stratification in Genome-wide Association Studies. *Nature Genetics* **38**: 904–909. doi:10.1038/ng1847.
- Pritchard JK, Stephens M, Donnelly P. 2000. Inference of Population Structure Using Multilocus Genotype Data. *Genetics* **155**: 945–959.
- Quek S-P, Counterman BA, Albuquerque de Moura P, Cardoso MZ, Marshall C, McMillan WO, Kronforst MR. 2010. Dissecting Mimetic Radiations in *Heliconius* Reveals Divergent Histories of Convergent Butterflies. *PNAS* **107**: 7365–7370. doi:10.1073/pnas.0911572107.
- R Development Core Team. 2011. *R: A Language and Environment for Statistical Computing* (version 2.14). Vienna, Austria: R Foundation for Statistical Computing. <http://www.R-project.org/>.
- Reed RD, Papa R, Martin A, Hines HM, Counterman BA, Pardo-Diaz C, Jiggins CD, Chamberlain NI, Kronforst MR, Chen R, et al. 2011. Optix Drives the Repeated Convergent Evolution of Butterfly Wing Pattern Mimicry. *Science* **333**: 1137–1141. doi:10.1126/science.1208227.
- Rieseberg L, Buerkle CA. 2002. Genetic Mapping in Hybrid Zones. *The American Naturalist* **159**: S36–S50. doi:10.1086/338371.
- Rosser N, Phillimore AB, Huertas B, Willmott KR, Mallet J. 2012. Testing Historical Explanations for Gradients in Species Richness in Heliconiine Butterflies of Tropical America. *Biological Journal of the Linnean Society* **105**: 479–497. doi:10.1111/j.1095-8312.2011.01814.x.
- Salazar PA. 2012. Hybridization and the Genetics of Wing Colour-pattern Diversity in *Heliconius* Butterflies. PhD, Cambridge, UK: University of Cambridge.
- Sheppard PM, Turner JRG, Brown KS, Benson WW, Singer MC. 1985. Genetics and the Evolution of Mullerian Mimicry in *Heliconius* Butterflies. *Phil. Trans. Roy. Soc. B* **308**: 433–610. doi:10.2307/2398716.
- Stewart C-B, Schilling JW, Wilson AC. 1987. Adaptive Evolution in the Stomach Lysozymes of Foregut Fermenters. *Nature* **330**: 401–404. doi:10.1038/330401a0.
- Stinchcombe JR, Hoekstra HE. 2007. Combining Population Genomics and Quantitative Genetics: Finding the Genes Underlying Ecologically Important Traits. *Heredity* **100**: 158–170.
- Supple MA, Hines HM, Dasmahapatra KK, Lewis JJ, Nielsen DM, Lavoie C, Ray DA, Salazar C, McMillan WO, Counterman BA. 2013. Genomic Architecture of Adaptive Color Pattern Divergence and Convergence in *Heliconius* Butterflies. *Genome Research* **23**: 1248–1257. doi:10.1101/gr.150615.112.
- The *Heliconius* Genome Consortium. 2012. Butterfly Genome Reveals Promiscuous Exchange of Mimicry Adaptations Among Species. *Nature* **487**: 94–98. doi:10.1038/nature11041.
- Tobler A, Kapan D, Flanagan NS, Gonzalez C, Peterson E, Jiggins CD, Johnstone JS, Heckel DG, McMillan WO. 2004. First-generation Linkage Map of the Warningly Colored Butterfly *Heliconius erato*. *Heredity* **94**: 408–417. doi:10.1038/sj.hdy.6800619.

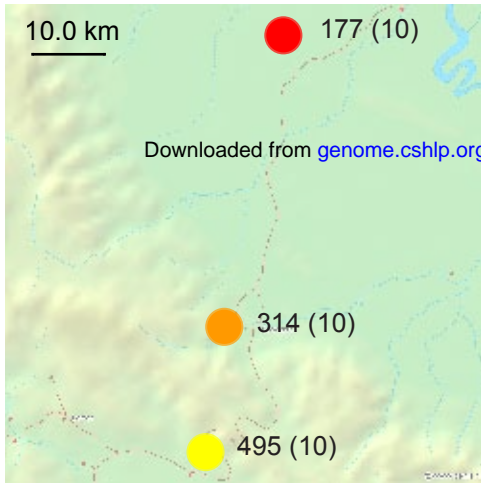
- Turner JR, Johnson MS, Eanes WF. 1979. Contrasted Modes of Evolution in the Same Genome: Allozymes and Adaptive Change in *Heliconius*. *Proceedings of the National Academy of Sciences* **76**: 1924–1928.
- Via S. 2012. Divergence Hitchhiking and the Spread of Genomic Isolation During Ecological Speciation-with-gene-flow. *Phil. Trans. Roy. Soc. B* **367**: 451–460. doi:10.1098/rstb.2011.0260.
- Visscher PM, Brown MA, McCarthy MI, Yang J. 2012. Five Years of GWAS Discovery. *The American Journal of Human Genetics* **90**: 7–24. doi:10.1016/j.ajhg.2011.11.029.
- Wood TE, Burke TM, Rieseberg LH. 2005. Parallel Genotypic Adaptation: When Evolution Repeats Itself. *Genetica* **123**: 157–170. doi:10.1007/s10709-003-2738-9.



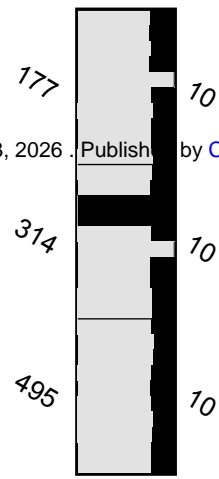
B.
H. melpomene clade



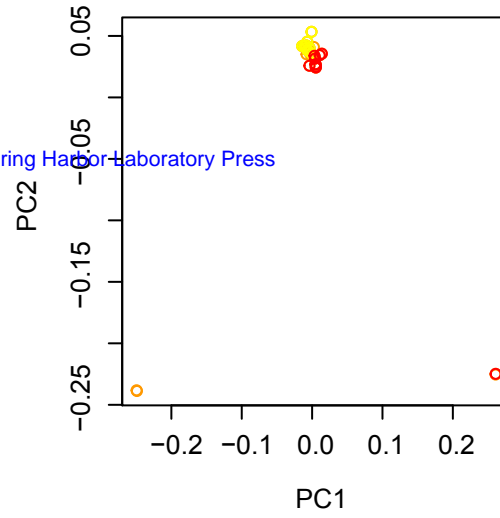
A.
***H. melpomene* Peru**



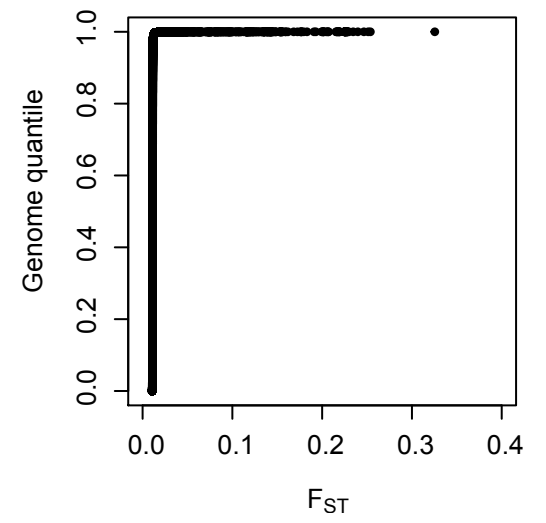
B. Structure
K=2



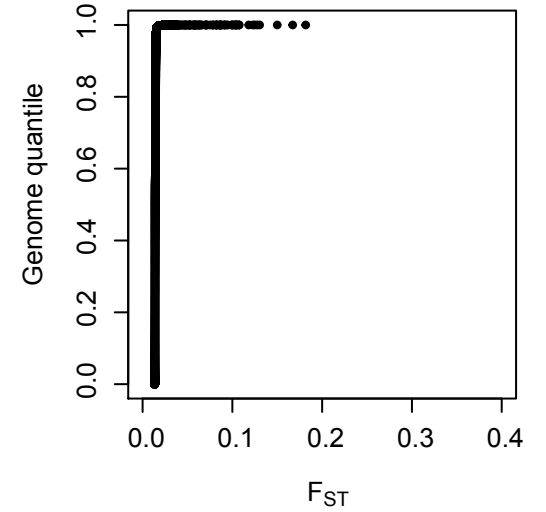
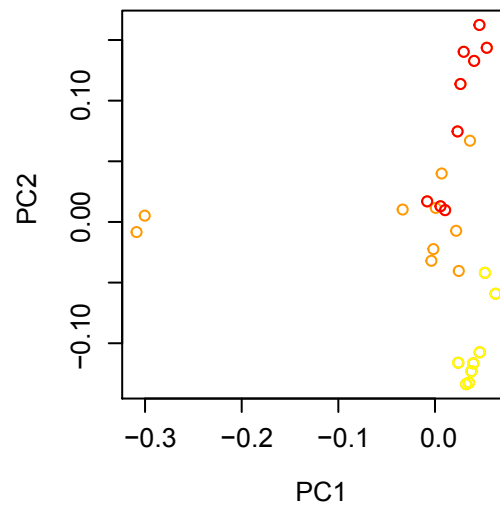
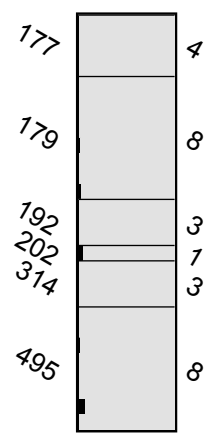
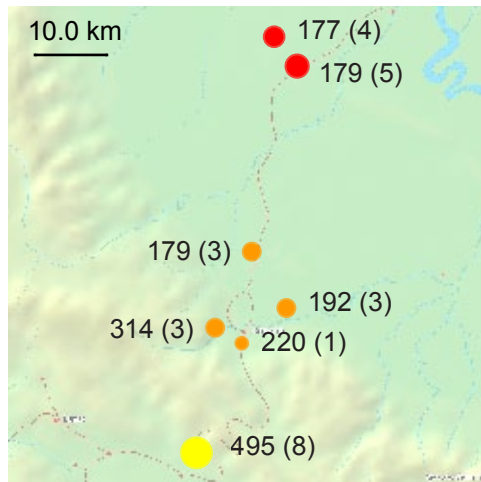
C. PCA



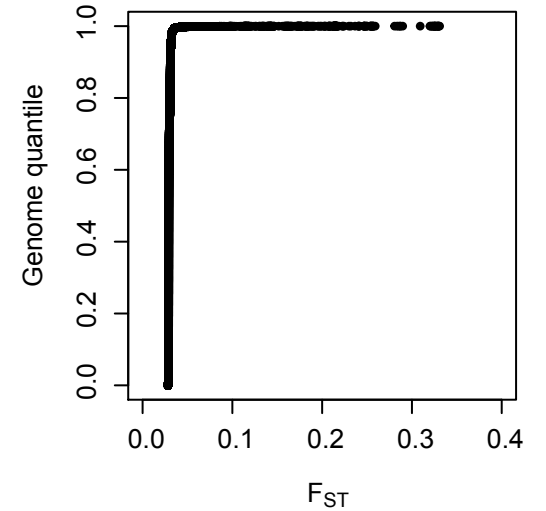
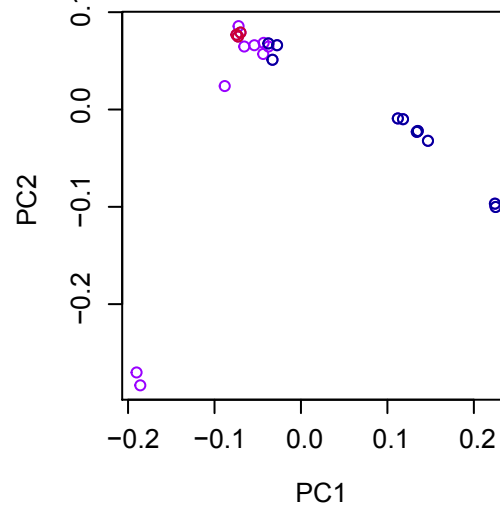
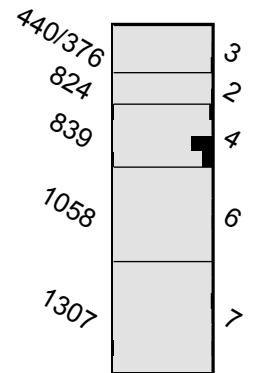
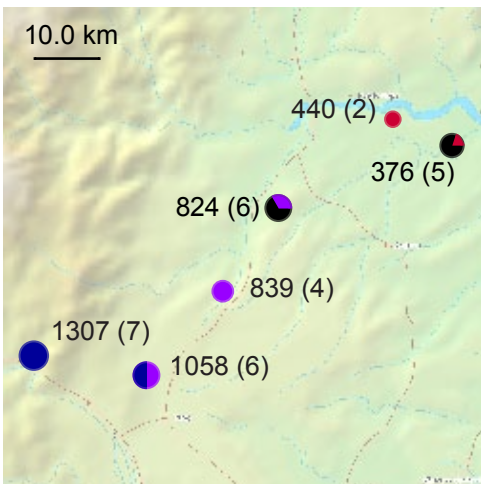
D. Cumulative Distribution of
between-subspecies F_{ST}



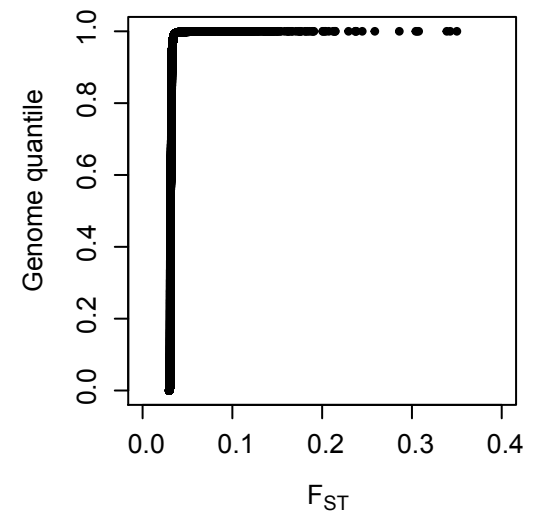
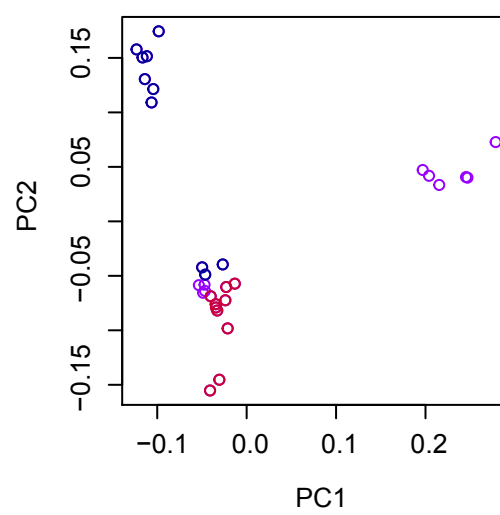
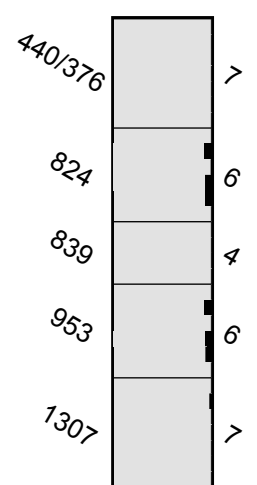
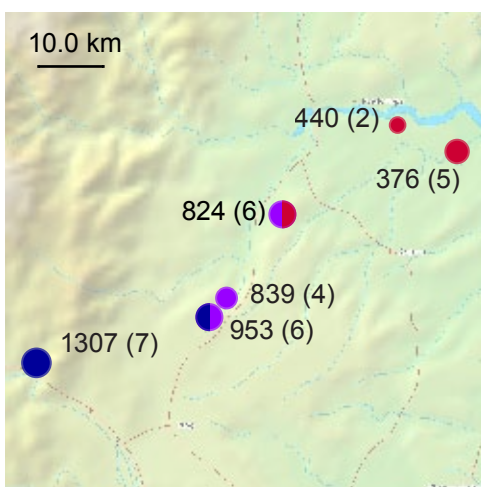
***H. erato* Peru**

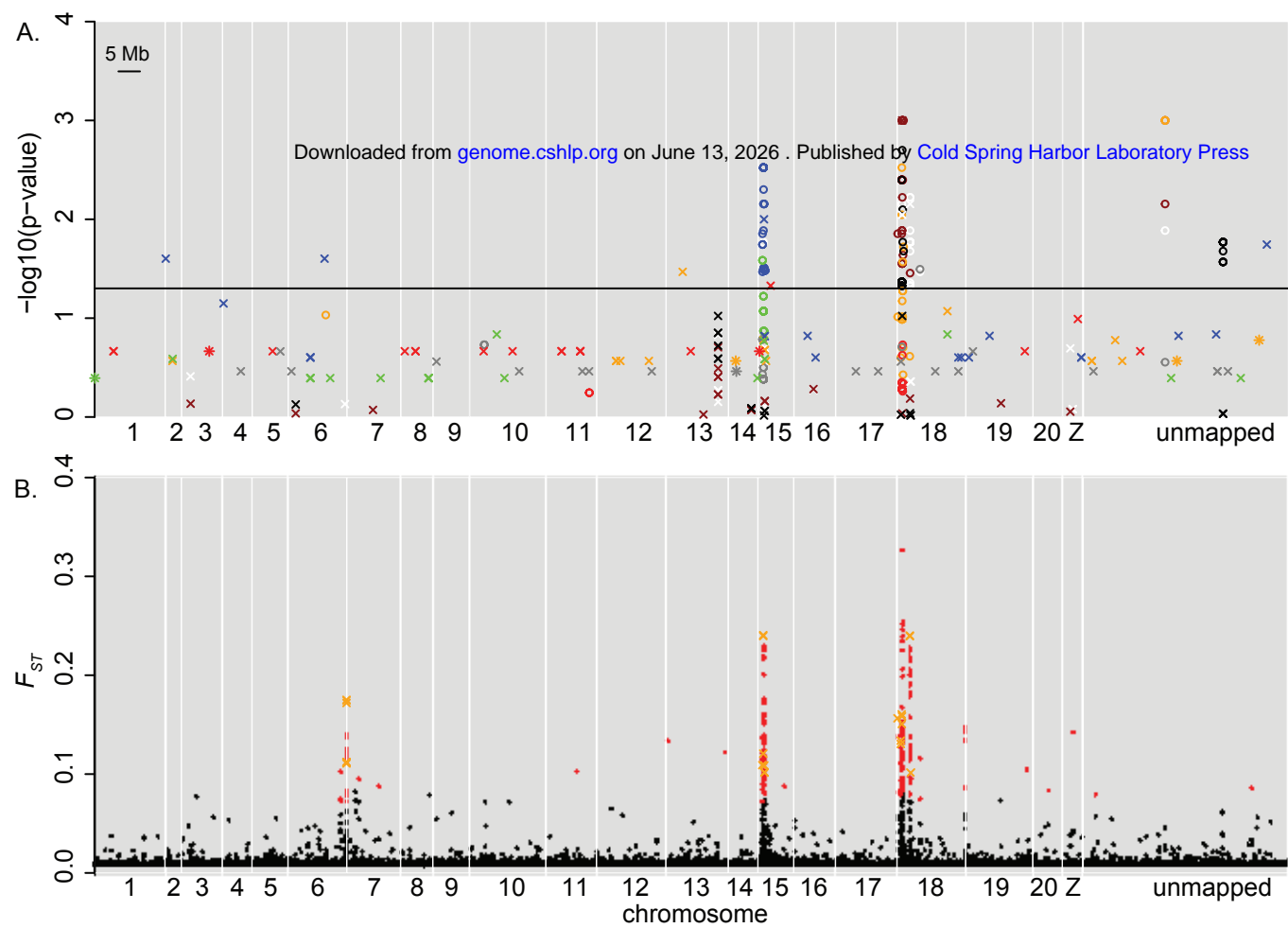


***H. melpomene* Ecuador**

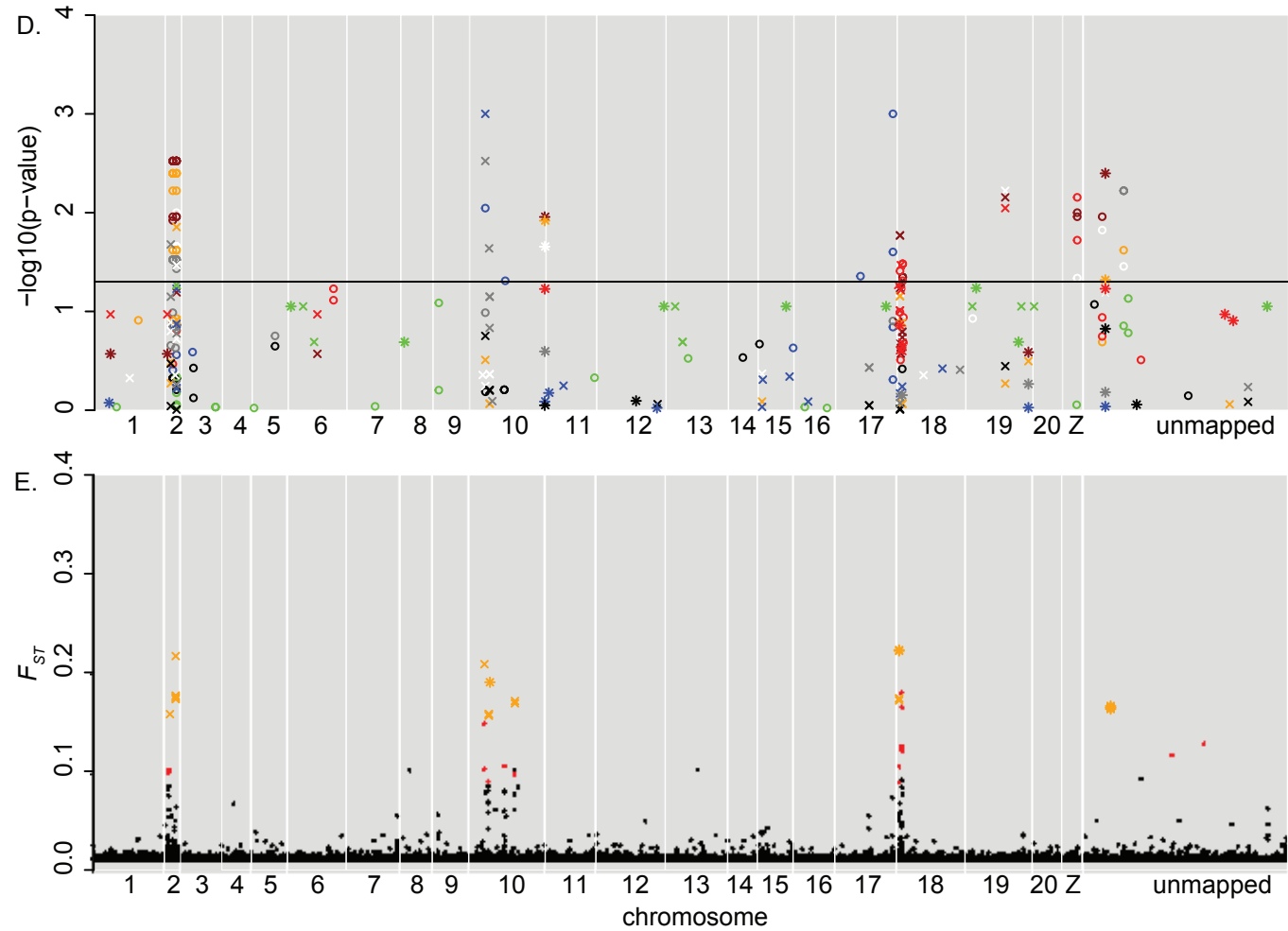
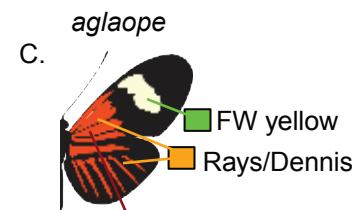


***H. erato* Ecuador**

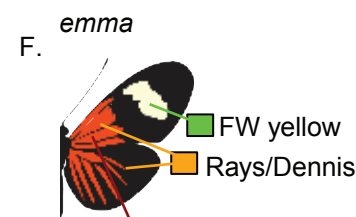




H. melpomene Peru



H. erato Peru



favorinus

