

The gene *cortex* controls mimicry and crypsis in butterflies and moths

Nicola J. Nadeau^{1,2}, Carolina Pardo-Diaz³, Annabel Whibley^{4,5}, Megan A. Supple^{2,6}, Suzanne V. Saenko⁴, Richard W. R. Wallbank^{2,7}, Grace C. Wu⁸, Luana Maroja⁹, Laura Ferguson¹⁰, Joseph J. Hanly^{2,7}, Heather Hines¹¹, Camilo Salazar³, Richard M. Merrill^{2,7}, Andrea J. Dowling¹², Richard H. ffrench-Constant¹², Violaine Llaurens⁴, Mathieu Joron^{4,13}, W. Owen McMillan² & Chris D. Jiggins^{7,2}

¹Department of Animal and Plant Sciences, University of Sheffield, Western Bank, Sheffield, S10 2TN UK.

²Smithsonian Tropical Research Institute, Apartado Postal 0843-00153, Panamá, República de Panamá

³Biology Program, Faculty of Natural Sciences and Mathematics, Universidad del Rosario, Cra. 24 No 63C-69, Bogotá D.C., 111221, Colombia.

⁴Institut de Systématique, Evolution et Biodiversité (UMR 7205 CNRS, MNHN, UPMC, EPHE, Sorbonne Université), Museum National d'Histoire Naturelle, CP50, 57 rue Cuvier, 75005 Paris, France.

⁵Cell and Developmental Biology, John Innes Centre, Norwich, Norfolk NR4 7UH, UK.

⁶Research School of Biology, The Australian National University, 134 Linnaeus Way, Acton, ACT, 2601, Australia

⁷Department of Zoology, University of Cambridge, Downing Street, Cambridge, CB2 3EJ, UK

⁸Energy and Resources Group, University of California at Berkeley, California, 94720, USA.

⁹Department of Biology, Williams College, Williamstown, Massachusetts 01267,

USA ¹⁰Department of Zoology, University of Oxford, South Parks Rd, Oxford, OX1 3PS, UK

¹¹Penn State University, 517 Mueller, University Park, Pennsylvania 16802, USA.

¹²School of Biosciences, University of Exeter in Cornwall, Penryn, Cornwall TR10 9FE, UK.

¹³Centre d'Ecologie Fonctionnelle et Evolutive (CEFE, UMR 5175 CNRS, Université de Montpellier, Université Paul-Valéry Montpellier, EPHE), 1919 route de Mende, 34293 Montpellier, France

30 The wing patterns of butterflies and moths (Lepidoptera) are diverse and striking examples of
31 evolutionary diversification by natural selection^{1,2}. Lepidopteran wing colour patterns are a
32 key innovation, consisting of arrays of coloured scales. We still lack a general understanding
33 of how these patterns are controlled and if there is any commonality across the 160,000 moth
34 and 17,000 butterfly species. Here, we identify a gene, *cortex*, through fine-scale mapping
35 using population genomics and gene expression analyses, which regulates pattern switches in
36 multiple species across the mimetic radiation in *Heliconius* butterflies. *cortex* belongs to a
37 fast evolving subfamily of the otherwise highly conserved fizzy family of cell cycle
38 regulators³, suggesting that it most likely regulates pigmentation patterning through
39 regulation of scale cell development. In parallel with findings in the peppered moth (*Biston*
40 *betularia*)⁴, our results suggest that this mechanism is common within Lepidoptera and that
41 *cortex* has become a major target for natural selection acting on colour and pattern variation
42 in this group of insects.

43

44 In *Heliconius*, there is a major effect locus, *Yb*, that controls a diversity of colour pattern
45 elements across the genus. It is the only locus in *Heliconius* that regulates all scale types and
46 colours, including the diversity of white and yellow pattern elements in the two co-mimics *H.*
47 *melpomene* (*Hm*) and *H. erato* (*He*), but also whole wing variation in black, yellow, white,
48 and orange/red elements in *H. numata* (*Hn*)⁵⁻⁷. In addition, genetic variation underlying the
49 *Bigeye* wing pattern mutation in *Bicyclus anynana*, melanism in the peppered moth, *Biston*
50 *betularia*, and melanism and patterning differences in the silkworm, *Bombyx mori*, have all
51 been localised to homologous genomic regions⁸⁻¹⁰ (Fig 1). Therefore, this genomic region
52 appears to contain one or more genes that act as major regulators of wing pigmentation and
53 patterning across the Lepidoptera.

54 Previous mapping of this locus in *He*, *Hm* and *Hn* identified a genomic interval of ~1Mb¹¹⁻¹³
55 (Extended Data Table 1), which also overlaps with the 1.4Mb region containing the
56 *carbonaria* locus in *B. betularia*⁹ and a 100bp non-coding region containing the *Ws* mutation
57 in *B. mori*¹⁰ (Fig 1). We took a population genomics approach to identify single nucleotide
58 polymorphisms (SNPs) most strongly associated with phenotypic variation within the ~1Mb
59 *Heliconius* interval. The diversity of wing patterning in *Heliconius* arises from divergence at
60 wing pattern loci⁷, while convergent patterns generally involve the same loci and sometimes
61 even the same alleles¹⁴⁻¹⁶. We used this pattern of divergence and sharing to identify SNPs
62 associated with colour pattern elements across many individuals from a wide diversity of
63 colour pattern phenotypes (Fig 2).

64 In three separate *Heliconius* species, our analysis consistently implicated the gene *cortex* as
65 being involved in adaptive differences in wing colour pattern. In *He* the strongest associations
66 with the presence of a yellow hindwing bar were centred around the genomic region
67 containing *cortex* (Fig 2A). We identified 108 SNPs that were fixed for one allele in *He*
68 *favorinus*, and fixed for the alternative allele in all individuals lacking the yellow bar, the
69 majority of which were in introns of *cortex* (Extended Data Table 2). 15 SNPs showed a
70 similar fixed pattern for *He demophoon*, which also has a yellow bar. These were non-
71 overlapping with those in *He favorinus*, consistent with the hypothesis that this phenotype
72 evolved independently in the two disjunct populations¹⁷.

73 Previous work has suggested that alleles at the *Yb* locus are shared between *Hm* and the
74 closely related species *H. timareta*, and also the more distantly related species *H. elevatus*,
75 resulting in mimicry between these species¹⁸. Across these species, the strongest associations
76 with the yellow hindwing bar phenotype were again found at *cortex* (Fig 2D, Extended Data
77 Fig 1A and Table 3). Similarly, the strongest associations with the yellow forewing band
78 were found around the 5' UTRs of *cortex* and gene *HM00036*, an orthologue of *D.*

79 *melanogaster washout* gene. A single SNP ~17kb upstream of *cortex* (the closest gene) was
80 perfectly associated with the yellow forewing band across all *Hm*, *H. timareta* and *H.*
81 *elevatus* individuals (Extended Data Fig 1A, Fig 2 and Table 3). We found no fixed coding
82 sequence variants at *cortex* in a larger sample (43-61 individuals) of *Hm aglaope* and *Hm*
83 *amaryllis* (Extended Data Figure 3, Supplementary Information), which differ in *Yb*
84 controlled phenotypes¹⁹, suggesting that functional variants are likely to be regulatory rather
85 than coding. We found extensive transposable element variation around *cortex* but it is
86 unclear if any of these associate with phenotype (Extended Data Figure 3 and Table 4;
87 Supplementary Information).

88 Finally, in *Hn* large inversions at the *P* supergene locus (Fig 1) are associated with different
89 morphs¹³. There is a steep increase in genotype-by-phenotype association at the breakpoint of
90 inversion 1, consistent with the role of these inversions in reducing recombination (Fig 2E).
91 However, the *bicoloratus* morph can recombine with all other morphs across one or the other
92 inversion, permitting finer-scale association mapping of this region. As in *He* and *Hm*, this
93 analysis showed a narrow region of associated SNPs corresponding exactly to the *cortex* gene
94 (Fig 2E), again with the majority of SNPs in introns (Extended Data Table 2). This associated
95 region does not correspond to any other known genomic feature, such as an inversion or
96 inversion breakpoint.

97 To determine whether sequence variants around *cortex* were regulating its expression we
98 investigated gene expression across the *Yb* locus. We used a custom designed microarray
99 including probes from all predicted genes in the *H. melpomene* genome¹⁸, as well as probes
100 tiled across the central portion of the *Yb* locus, focussing on two naturally hybridising *Hm*
101 races (*plesseni* and *malleti*) that differ in *Yb* controlled phenotypes⁷. *cortex* was the only gene
102 across the entire interval to show significant expression differences both between races with
103 different wing patterns and between wing sections with different pattern elements (Fig 3).

104 This finding was reinforced in the tiled probe set, where we observed strong differences in
105 expression of *cortex* exons and introns but few differences outside this region (Extended Data
106 Table 2). *cortex* expression was higher in *Hm malleti* than *Hm plesseni* in all three wing
107 sections used (but not eyes) (Fig 3C; Extended Data Fig 4C). When different wing sections
108 were compared within each race, *cortex* expression in *Hm malleti* was higher in the distal
109 section that contains the *Yb* controlled yellow forewing band, consistent with *cortex*
110 producing this band. In contrast, *Hm plesseni*, which lacks the yellow band, had higher *cortex*
111 expression in the proximal forewing section (Fig 3F; Extended Data Fig 4J). Expression
112 differences were found only in day 1 and day 3 pupal wings rather than day 5 or day 7
113 (Extended Data Fig 4), similar to the pattern observed previously for the transcription factor
114 *optix*²⁰.

115 Differential expression was not confined to the exons of *cortex*; the majority of differentially
116 expressed probes in the tiling array corresponded to *cortex* introns (Fig 3). This does not
117 appear to be due to transposable element variation (Extended Data Table 2), but may be due
118 to elevated background transcription and unidentified splice variants. RT-PCR revealed a
119 diversity of splice variants (Extended Data Fig 5), and sequenced products revealed 8 non-
120 constitutive exons and 6 variable donor/acceptor sites, but this was not exhaustive
121 (Supplementary Information). We cannot rule out the possibility that some of the
122 differentially expressed intronic regions could be distinct non-coding RNAs. However, qRT-
123 PCR in other hybridising races with divergent *Yb* alleles (*aglaope/amaryllis* and
124 *rosina/melpomene*) also identified expression differences at *cortex* and allele-specific splicing
125 differences between both pairs of races (Extended Data Figs 1 and 5, Supplementary
126 Information).

127 Finally, *in situ* hybridisation of *cortex* in final instar larval hindwing discs showed expression
128 in wing regions fated to become black in the adult wing, most strikingly in their

129 correspondence to the black patterns on adult *Hn* wings (Fig 4). In contrast, the array results
130 from pupal wings were suggestive of higher expression in non-melanic regions. This may
131 suggest that *cortex* is upregulated at different time-points in wing regions fated to become
132 different colours.

133 Overall, *cortex* shows significant differential expression and is the only gene in the candidate
134 region to be consistently differentially expressed in multiple race comparisons and between
135 differently patterned wing regions. Coupled with the strong genotype-by-phenotype
136 associations across multiple independent lineages (Extended Data Table 1), this strongly
137 implicates *cortex* as a major regulator of colour and pattern. However, we have not excluded
138 the possibility that other genes in this region also influence pigmentation patterning. A
139 prominent role for *cortex* is also supported by studies in other taxa; our identification of
140 distant 5' untranslated exons of *cortex* (Supplementary Information) suggests that the 100bp
141 interval containing the *Ws* mutation in *B. mori* is likely to be within an intron of *cortex* and
142 not in intergenic space as previously thought¹⁰. In addition, fine-mapping and gene
143 expression also implicate *cortex* as controlling melanism in the peppered moth⁴.

144 It seems likely that *cortex* controls pigmentation patterning through control of scale cell
145 development. The *cortex* gene falls in an insect specific lineage within the fizzy/CDC20
146 family of cell cycle regulators (Extended Data Fig 6A). The phylogenetic tree of the gene
147 family highlighted three major orthologous groups, two of which have highly conserved
148 functions in cell cycle regulation mediated through interaction with the anaphase promoting
149 complex/cyclosome (APC/C)^{3,21}. The third group, *cortex*, is evolving rapidly, with low amino
150 acid identity between *D. melanogaster* and *Hm cortex* (14.1%), contrasting with much higher
151 identities for orthologues between these species in the other two groups (*fzy*, 47.8% and
152 *rap/fzr*, 47.2%, Extended Data Fig 6A). *Drosophila melanogaster cortex* acts through a

153 similar mechanism to *fzy* in order to control meiosis in the female germ line²²⁻²⁴. *Hm cortex*
154 also has some conservation of the fizzy family C-box and IR elements (Supplementary
155 Information) that mediate binding to the APC/C²³, suggesting that it may have retained a cell
156 cycle function, although we found that expressing *Hm cortex* in *D. melanogaster* wings
157 produced no detectable effect (Extended Data Fig 6, Supplementary Information).

158 Previously identified butterfly wing patterning genes have been transcription factors or
159 signalling molecules^{20,25}. Developmental rate has long been thought to play a role in
160 lepidopteran patterning^{26,27}, but *cortex* was not a likely *a priori* candidate, because its
161 *Drosophila* orthologue has a highly specific function in meiosis²³. The recruitment of *cortex*
162 to wing patterning appears to have occurred before the major diversification of the
163 Lepidoptera and this gene has repeatedly been targeted by natural selection^{1,7,9,28} to generate
164 both cryptic⁴ and aposematic patterns.

165 **References**

- 166 1. Cook, L. M., Grant, B. S., Saccheri, I. J. & Mallet, J. Selective bird predation on the
167 peppered moth: the last experiment of Michael Majerus. *Biol. Lett.* **8**, 609–612 (2012).
- 168 2. Jiggins, C. D. Ecological Speciation in Mimetic Butterflies. *BioScience* **58**, 541–548
169 (2008).
- 170 3. Dawson, I. A., Roth, S. & Artavanis-Tsakonas, S. The *Drosophila* Cell Cycle Gene fizzy
171 Is Required for Normal Degradation of Cyclins A and B during Mitosis and Has
172 Homology to the CDC20 Gene of *Saccharomyces cerevisiae*. *J. Cell Biol.* **129**, 725–737
173 (1995).
- 174 4. Van't Hof, A. E. *et al.* The industrial melanism mutation in British peppered moths is a
175 transposable element. *Nature* **This issue**,

- 176 5. Joron, M. *et al.* A Conserved Supergene Locus Controls Colour Pattern Diversity in
177 Heliconius Butterflies. *PLoS Biol.* **4**, (2006).
- 178 6. Sheppard, P. M., Turner, J. R. G., Brown, K. S., Benson, W. W. & Singer, M. C.
179 Genetics and the Evolution of Mullerian Mimicry in Heliconius Butterflies. *Philos.*
180 *Trans. R. Soc. Lond. B. Biol. Sci.* **308**, 433–610 (1985).
- 181 7. Nadeau, N. J. *et al.* Population genomics of parallel hybrid zones in the mimetic
182 butterflies, *H. melpomene* and *H. erato*. *Genome Res.* **24**, 1316–1333 (2014).
- 183 8. Beldade, P., Saenko, S. V., Pul, N. & Long, A. D. A Gene-Based Linkage Map for
184 *Bicyclus anynana* Butterflies Allows for a Comprehensive Analysis of Synteny with the
185 Lepidopteran Reference Genome. *PLoS Genet* **5**, e1000366 (2009).
- 186 9. van't Hof, A. E., Edmonds, N., Dalíková, M., Marec, F. & Saccheri, I. J. Industrial
187 Melanism in British Peppered Moths Has a Singular and Recent Mutational Origin.
188 *Science* **332**, 958–960 (2011).
- 189 10. Ito, K. *et al.* Mapping and recombination analysis of two moth colour mutations, Black
190 moth and Wild wing spot, in the silkworm *Bombyx mori*. *Heredity* (2015).
191 doi:10.1038/hdy.2015.69
- 192 11. Counterman, B. A. *et al.* Genomic Hotspots for Adaptation: The Population Genetics of
193 Müllerian Mimicry in *Heliconius erato*. *PLoS Genet.* **6**, e1000796 (2010).
- 194 12. Ferguson, L. *et al.* Characterization of a hotspot for mimicry: assembly of a butterfly
195 wing transcriptome to genomic sequence at the HmYb/Sb locus. *Mol. Ecol.* **19**, 240–254
196 (2010).
- 197 13. Joron, M. *et al.* Chromosomal rearrangements maintain a polymorphic supergene
198 controlling butterfly mimicry. *Nature* **477**, 203–206 (2011).
- 199 14. Hines, H. M. *et al.* Wing patterning gene redefines the mimetic history of *Heliconius*
200 butterflies. *Proc. Natl. Acad. Sci.* **108**, 19666–19671 (2011).

- 201 15. Pardo-Diaz, C. *et al.* Adaptive Introgression across Species Boundaries in Heliconius
202 Butterflies. *PLoS Genet* **8**, e1002752 (2012).
- 203 16. Wallbank, R. W. R. *et al.* Evolutionary Novelty in a Butterfly Wing Pattern through
204 Enhancer Shuffling. *PLoS Biol* **14**, e1002353 (2016).
- 205 17. Maroja, L. S., Alschuler, R., McMillan, W. O. & Jiggins, C. D. Partial Complementarity
206 of the Mimetic Yellow Bar Phenotype in Heliconius Butterflies. *PLoS ONE* **7**, e48627
207 (2012).
- 208 18. The Heliconius Genome Consortium. Butterfly genome reveals promiscuous exchange of
209 mimicry adaptations among species. *Nature* **487**, 94–98 (2012).
- 210 19. Mallet, J. The Genetics of Warning Colour in Peruvian Hybrid Zones of *Heliconius erato*
211 and *H. melpomene*. *Proc. R. Soc. Lond. B Biol. Sci.* **236**, 163–185 (1989).
- 212 20. Reed, R. D. *et al.* optix Drives the Repeated Convergent Evolution of Butterfly Wing
213 Pattern Mimicry. *Science* **333**, 1137–1141 (2011).
- 214 21. Barford, D. Structural insights into anaphase-promoting complex function and
215 mechanism. *Philos. Trans. R. Soc. B Biol. Sci.* **366**, 3605–3624 (2011).
- 216 22. Chu, T., Henrion, G., Haegeli, V. & Strickland, S. Cortex, a *Drosophila* gene required to
217 complete oocyte meiosis, is a member of the Cdc20/fizzy protein family. *genesis* **29**,
218 141–152 (2001).
- 219 23. Pesin, J. A. & Orr-Weaver, T. L. Developmental Role and Regulation of cortex, a
220 Meiosis-Specific Anaphase-Promoting Complex/Cyclosome Activator. *PLoS Genet* **3**,
221 e202 (2007).
- 222 24. Swan, A. & Schüpbach, T. The Cdc20/Cdh1-related protein, Cort, cooperates with
223 Cdc20/Fzy in cyclin destruction and anaphase progression in meiosis I and II in
224 *Drosophila*. *Dev. Camb. Engl.* **134**, 891–899 (2007).

- 225 25. Martin, A. *et al.* Diversification of complex butterfly wing patterns by repeated
226 regulatory evolution of a Wnt ligand. *Proc. Natl. Acad. Sci.* **109**, 12632–12637 (2012).
- 227 26. Koch, P. B., Lorenz, U., Brakefield, P. M. & ffrench-Constant, R. H. Butterfly wing
228 pattern mutants: developmental heterochrony and co-ordinately regulated phenotypes.
229 *Dev. Genes Evol.* **210**, 536–544 (2000).
- 230 27. Gilbert, L. E., Forrest, H. S., Schultz, T. D. & Harvey, D. J. Correlations of ultrastructure
231 and pigmentation suggest how genes control development of wing scales of *Heliconius*
232 butterflies. *J. Res. Lepidoptera* **26**, 141–160 (1988).
- 233 28. Mallet, J. & Barton, N. H. Strong Natural Selection in a Warning-Color Hybrid Zone.
234 *Evolution* **43**, 421–431 (1989).
- 235 29. Wahlberg, N., Wheat, C. W. & Peña, C. Timing and Patterns in the Taxonomic
236 Diversification of Lepidoptera (Butterflies and Moths). *PLoS ONE* **8**, e80875 (2013).
- 237 30. Surridge, A. *et al.* Characterisation and expression of microRNAs in developing wings of
238 the neotropical butterfly *Heliconius melpomene*. *BMC Genomics* **12**, 62 (2011).

239

240 **Supplementary Information** is linked to the online version of the paper at

241 www.nature.com/nature.

242 **Acknowledgements** We thank Christopher Sasaki, Clemson University, for assembly of the
243 *He* BACs. Moises Abanto and Adriana Tapia assisted with raising butterflies. Thanks to
244 Mathieu Chouteau, Jake Morris and Kanchon Dasmahapatra for providing larvae for *in situ*
245 hybridisations. Anna Morrison, Robert Tetley, Sarah Carl and Hanna Wegener assisted with
246 lab work at the University of Cambridge. Simon Baxter made the *Hm* fosmid libraries. We
247 thank the governments of Colombia, Ecuador, Panama and Peru for permission to collect
248 butterflies. This work was funded by a Leverhulme Trust award and BBSRC grant
249 (H01439X/1) to CDJ, NSF grants (DEB 1257689, IOS 1052541) to WOM, an ERC starting

250 grant to MJ and a French National Agency for Research (ANR) grant to VL (ANR-13-JSV7-
251 0003-01). NJN is funded by a NERC fellowship (NE/K008498/1).

252

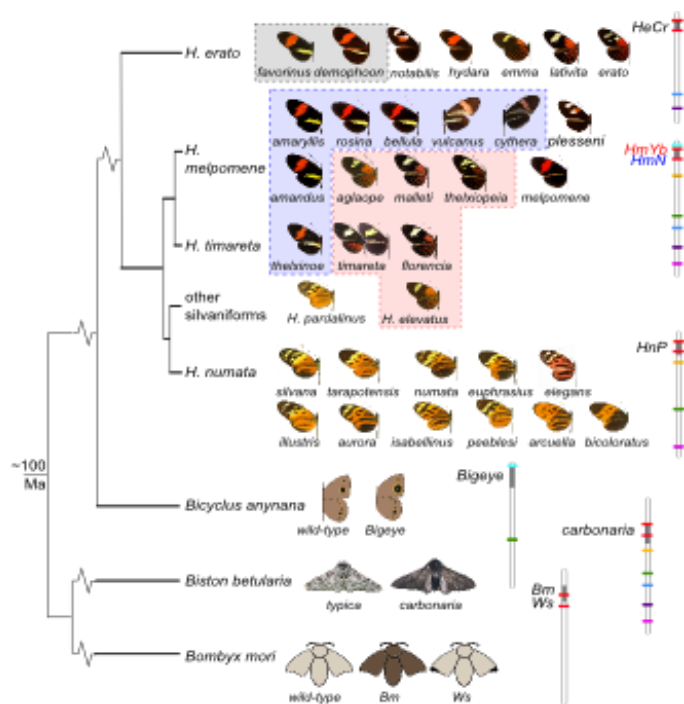
253 **Author Contributions** NJN performed the association analyses, 5' RACE, RT-PCR, qRT-
254 PCR and prepared the manuscript. NJN and CDJ co-ordinated the research. CP-D performed
255 and analysed the microarray and RNAseq experiments. AW performed the *Hn* association
256 analysis. MS assembled and annotated the *HeCr* BAC reference and the *He* alignments. SVS
257 performed *in situ* hybridizations. RWRW performed the transgenic experiments and analysis
258 of *de novo* assembled sequences and fosmids together with JJH. GW and LF initially
259 identified splicing variants of *cortex*. LM performed crosses between *Hm* races. HH screened
260 the *HeCr* BAC library. CS and RM provided samples. AD contributed to the *Hm* BAC
261 sequencing and annotation. R-fC, MJ, VL, WOM and CDJ are PIs who obtained funding and
262 led the project elements. All authors commented on the manuscript.

263

264 **Author Information** Short read sequence data generated for this study are available from
265 ENA (<http://www.ebi.ac.uk/ena>) under study accession PRJEB8011 and PRJEB12740 (see
266 Supplementary Table 1 for previously published data accessions). The updated Cr contig is
267 deposited in Genbank with accession KC469893. The assembled *Hm* fosmid sequences are
268 deposited in Genbank with accessions KU514430-KU514438. The microarray data are
269 deposited in GEO with accessions GSM1563402- GSM1563497. Reprints and permissions
270 information is available at www.nature.com/reprints. Correspondence and requests for
271 materials should be addressed to n.nadeau@sheffield.ac.uk or c.jiggins@zoo.cam.ac.uk

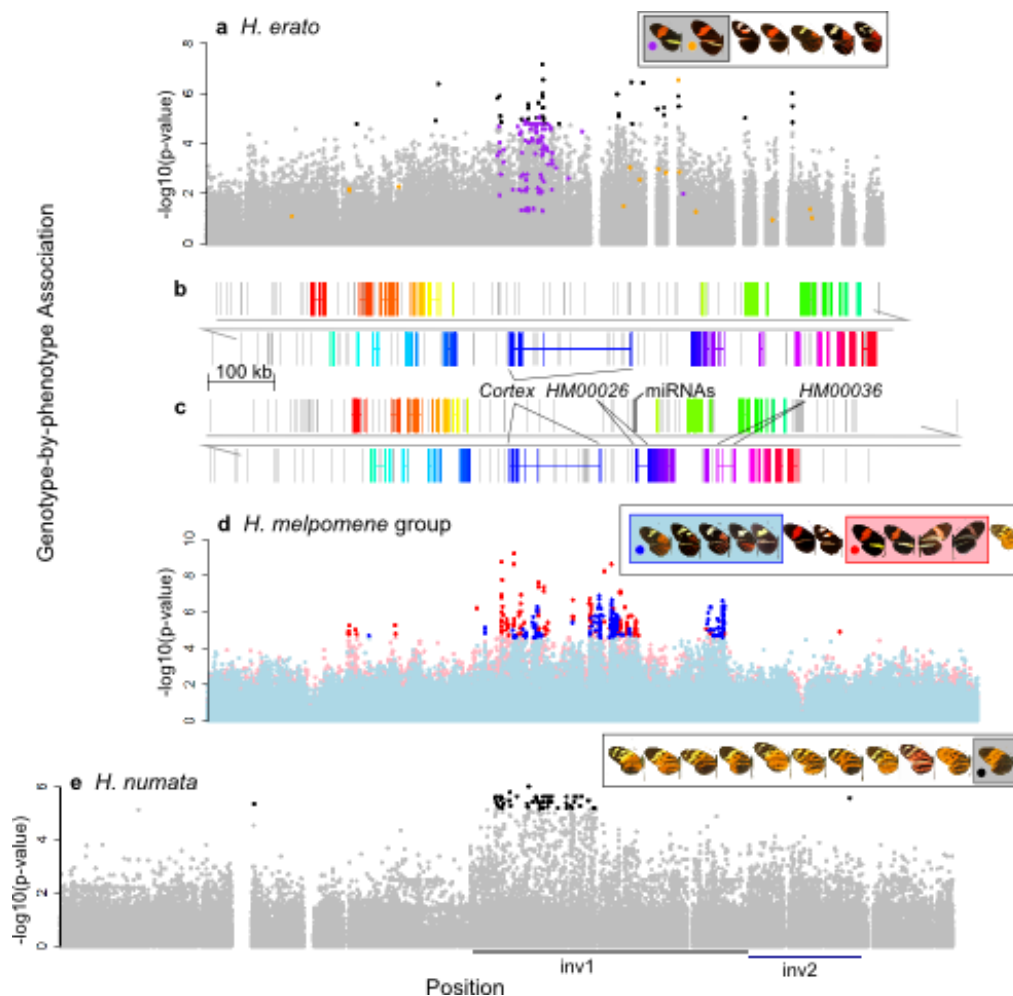
272

273



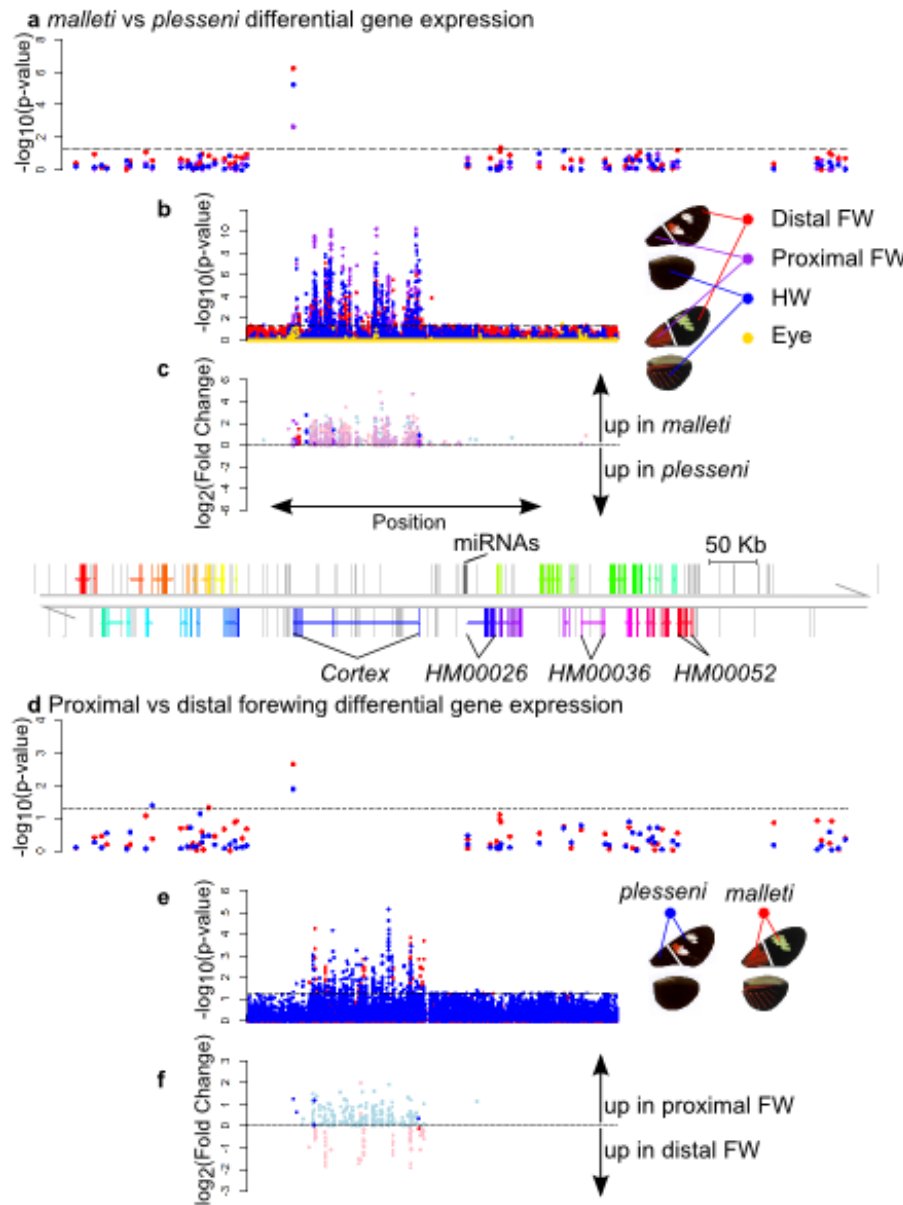
274

275 Figure 1. A homologous genomic region controls a diversity of phenotypes across the
 276 Lepidoptera. Left: phylogenetic relationships²⁹. Right: chromosome maps with colour pattern
 277 intervals in grey, coloured bars represent markers used to assign homology^{5,8-10}, the first and
 278 last genes from Fig 2 shown in red. In *He* the *HeCr* locus controls the yellow hind-wing bar
 279 phenotype (grey boxed races). In *Hm* it controls both the yellow hind-wing bar (*HmYb*, pink
 280 box) and the yellow forewing band (*HmN*, blue box). In *Hn* it modulates black, yellow and
 281 orange elements on both wings (*HnP*), producing phenotypes that mimic butterflies in the
 282 genus *Melinaea*. Morphs/races of *Heliconius* species included in this study are shown with
 283 names.



284

285 Figure 2. Association analyses across the genomic region known to contain major colour
 286 pattern loci in *Heliconius*. A) Association in *He* with the yellow hind-wing bar (n=45).
 287 Coloured SNPs are fixed for a unique state in *He demophoon* (orange) or *He favorinus*
 288 (purple). B) Genes in *He* with direct homologs in *Hm*. Genes are in different colours with
 289 exons (coding and UTRs) connected by a line. Grey bars are transposable elements. C) *Hm*
 290 genes and transposable elements: colours correspond to homologous *He* genes; MicroRNAs³⁰
 291 in black. D) Association in the *Hm/timareta/silvaniform* group with the yellow hind-wing bar
 292 (red) and yellow forewing band (blue) (n=49). E) Association in *Hn* with the *bicoloratus*
 293 morph (n=26); inversion positions¹³ shown below. In all cases black/dark coloured points are
 294 above the strongest associations found outside the colour pattern scaffolds (*He* p=1.63e-05;
 295 *Hm* p=2.03e-05 and p=2.58e-05; *Hn* p=6.81e-06).

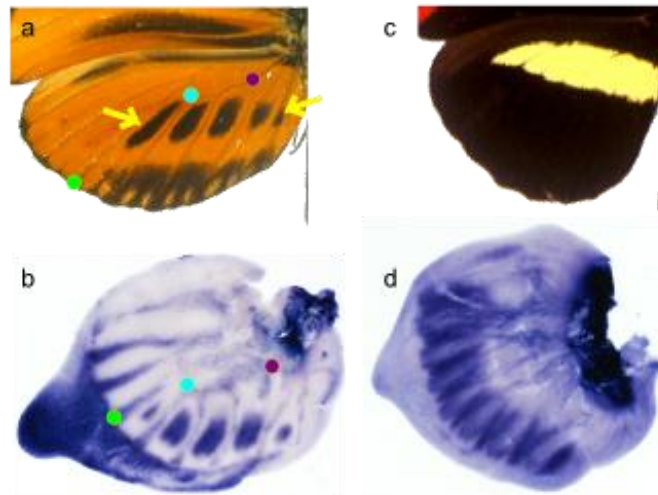


296

297 Figure 3. Differential gene expression across the genomic region known to contain major
 298 colour pattern loci in *Heliconius melpomene*. Expression differences in day 3 pupae, for all
 299 genes in the *Yb* interval (A,D) and tiling probes spanning the central portion of the interval
 300 (B,C,E,F). Expression is compared between races for each wing region (A,B,C) and between
 301 proximal and distal forewing sections for each race (D,E,F). C and F: magnitude and
 302 direction of expression difference (\log_2 fold-change) for tiling probes showing significant
 303 differences ($p \leq 0.05$); probes in known *cortex* exons shown in dark colours. Gene *HM00052*

304 was differentially expressed between other races in RNA sequence data (Supplementary
 305 Information) but is not differentially expressed here.

306



307

308 Figure 4. *In situ* hybridisations of *cortex* in hind-wings of final instar larvae. B) *Hn*
 309 *tarapotensis*; adult wing shown in A, coloured points indicate landmarks, yellow arrows
 310 highlight adult pattern elements corresponding to the *cortex* staining. D) *Hm rosina*; adult
 311 wing shown in C, staining patterns in other *Hm* races (*meriana* and *aglaope*) appeared
 312 similar. The probe used was complementary to the *cortex* isoform with the longest open
 313 reading frame (also the most common, Supplementary Information).

314

315 **Methods**

316 *He Cr* reference

317 *Cr* is the homologue of *Yb* in *He* (Fig 1). An existing reference for this region was available
 318 in 3 pieces (467,734bp, 114,741bp and 161,149bp, GenBank: KC469893.1)³¹. We screened
 319 the same BAC library used previously^{11,31} using described procedures¹¹ with probes designed

320 to the ends of the existing BAC sequences and the *HmYb* BAC reference sequence. Two
321 BACs (04B01 and 10B14) were identified as spanning one of the gaps and sequenced using
322 Illumina 2x250 bp paired-end reads collected on the Illumina MiSeq. The raw reads were
323 screened to remove vector and *E. coli* bases. The first 50k read pairs were taken for each
324 BAC and assembled individually with the Phrap³² software and manually edited with
325 consed³³. Contigs with discordant read pairs were manually broken and properly merged
326 using concordant read data. Gaps between contig ends were filled using an in-house
327 finishing technique where the terminal 200bp of the contig ends were extracted and queried
328 against the unused read data for spanning pairs, which were added using the
329 addSolexaReads.perl script in the consed package. Finally, a single reference contig was
330 generated by identifying and merging overlapping regions of the two consensus BAC
331 sequences.

332 In order to fill the remaining gap (between positions 800,387 and 848,446) we used the
333 overhanging ends to search the scaffolds from a preliminary *He* genome assembly of five
334 Illumina paired end libraries with different insert sizes (250, 500, 800, 4300 and 6500bp)
335 from two related *He demophon* individuals. We identified two scaffolds (scf1869 and
336 scf1510) that overlapped and spanned the gap (using 12,257bp of the first scaffold and
337 35,803bp of the second).

338 The final contig was 1,009,595bp in length of which 2,281bp were unknown (N's). The *HeCr*
339 assembly was verified by aligning to the *HmYb* genome scaffold (HE667780) with mummer
340 and blast. The *HeCr* contig was annotated as described previously³², with some minor
341 modifications. Briefly this involved first generating a reference based transcriptome assembly
342 with existing *H. erato* RNA-seq wing tissue (GenBank accession SRA060220). We used
343 Trimmomatic³⁴ (v0.22), and FLASH³⁵ (v1.2.2) to prepare the raw sequencing reads, checking
344 the quality with FastQC³⁶ (v0.10.0). We then used the Bowtie/TopHat/Cufflinks³⁷⁻³⁹ pipeline

345 to generate transcripts for the unmasked reference sequence. We generated gene predictions
346 with the MAKER pipeline⁴⁰ (v2.31). Homology and synteny in gene content with the *Hm Yb*
347 reference were identified by aligning the *Hm* coding sequences to the *He* reference with
348 BLAST. Homologous genes were present in the same order and orientation in *He* and *Hm*
349 (Fig 2B,C). Annotations were manually adjusted if genes had clearly been merged or split in
350 comparison to *H. melpomene* (which has been extensively manually curated¹²). In addition
351 *He cortex* was manually curated from the RNA-seq data and using *Exonerate*⁴¹ alignments of
352 the *H. melpomene* protein and mRNA transcripts, including the 5' UTRs.

353 ***Genotype-by-phenotype association analyses***

354 Information on the individuals used and ENA accessions for sequence data are given in
355 Supplementary Table 1. We used shotgun Illumina sequence reads from 45 *He* individuals
356 from 7 races that were generated as part of a previous study³¹ (Supplementary Information).
357 Reads were aligned to an *He* reference containing the *Cr* contig and other sequenced *He*
358 BACs^{11,31} with BWA⁴², which has previously been found to work better than Stampy⁴³
359 (which was used for the alignments in the other species) with an incomplete reference
360 sequence³¹. The parameters used were as follows: Maximum edit distance (n), 8; maximum
361 number of gap opens (o), 2; maximum number of gap extensions (e), 3; seed (l), 35;
362 maximum edit distance in seed (k), 2. We then used Picard tools to remove PCR and optical
363 duplicate sequence reads and GATK⁴⁴ to re-align indels and call SNPs using all individuals
364 as a single population. Expected heterozygosity was set to 0.2 in GATK. 132,397 SNPs were
365 present across *Cr*. A further 52,698 SNPs not linked to colour pattern loci were used to
366 establish background association levels.

367 For the *Hm / Hn* clade we used previously published sequence data from 19 individuals from
368 enrichment sequencing targeting of the *Yb* region, the unlinked *HmB/D* region that controls

369 the presence/absence of red colour pattern elements, and ~1.8Mb of non-colour pattern
370 genomic regions⁴⁵, as well as 9 whole genome shotgun sequenced individuals^{18,46}. We added
371 targeted sequencing and shotgun whole genome sequencing of an additional 47 individuals
372 (Supplementary Information). Alignments were performed using Stampy⁴³ with default
373 parameters except for substitution rate which was set to 0.01. We again removed duplicates
374 and used GATK to re-align indels and call SNPs with expected heterozygosity set to 0.1.

375 The analysis of the *Hm/timareta/silvaniform* included 49 individuals, which were aligned to
376 v1.1 of the *Hm* reference genome with the scaffolds containing *Yb* and *HmB/D* swapped with
377 reference BAC sequences¹⁸, which contained fewer gaps of unknown sequence than the
378 genome scaffolds. 232,631 SNPs were present in the *Yb* region and a further 370,079 SNPs
379 were used to establish background association levels.

380 The *Hn* analysis included 26 individuals aligned to unaltered v1.1 of the *Hm* reference
381 genome, because the genome scaffold containing *Yb* is longer than the BAC reference
382 making it easier to compare the inverted and non-inverted regions present in this species. We
383 tested for associations at 262,137 SNPs on the *Yb* scaffold with the *Hn bicoloratus* morph,
384 which had a sample of 5 individuals.

385 We measured associations between genotype and phenotype using a score test (qtscore) in the
386 GenABEL package in R⁴⁷. This was corrected for background population structure using a
387 test specific inflation factor, λ , calculated from the SNPs unlinked to the major colour pattern
388 controlling loci (described above), as the colour pattern loci are known to have different
389 population structure to the rest of the genome^{14,15,18}. We used a custom perl script to convert
390 GATK vcf files to Illumina SNP format for input to genABEL⁴⁷. genABEL does not accept
391 multiallelic sites, so the script also converted the genotype of any individuals for which a
392 third (or fourth) allele was present to a missing genotype (with these defined as the lowest

393 frequency alleles). Custom R scripts were used to identify sites showing perfect associations
394 with calls for >75% of individuals.

395 *Microarray Gene Expression Analyses*

396 We designed a Roche NimbleGen microarray (12x135K format) with probes for all annotated
397 *Hm* genes¹⁸ and tiling the central portion of the *Yb* BAC sequence contig that was previously
398 identified as showing the strongest differentiation between *Hm* races⁴⁵. In addition to the
399 *HmYb* tiling array probes there were 6,560 probes tiling *HmAc* (a third unlinked colour
400 pattern locus) and 10,716 probes tiling *HmB/D*, again distanced on average at 10bp intervals.
401 The whole-genome gene expression array contained 107,898 probes in total.

402 This was interrogated with Cy3 labelled double stranded cDNA generated from total RNA
403 (with a SuperScript double-stranded cDNA synthesis kit, Invitrogen, and a one-colour DNA
404 labelling kit, Niblegen) from four pupal developmental stages of *Hm plesseni* and *malleti*.
405 Pupae were from captive stocks maintained in insectary facilities in Gamboa, Panama. Tissue
406 was stored in RNA later at -80°C prior to RNA extraction. RNA was extracted using TRIZOL
407 (Invitrogen) followed by purification with RNeasy (Qiagen) and DNase treated with DNA-
408 free (Ambion). Quantification was performed using a Qubit 2.0 fluorometer (Invitrogen) and
409 purity and integrity assessed using a Bioanalyzer 2100 (Agilent). Samples were randomised
410 and each hybridised to a separate array. The *HmYb* probe array contained 9,979 probes
411 distanced on average at 10bp. The whole-genome expression array contained on average 9
412 probes per annotated gene in the genome (v1.1¹⁸) as well as any transcripts not annotated but
413 predicted from RNA-seq evidence.

414 Background corrected expression values for each probe were extracted using NimbleScan
415 software (version 2.3). Analyses were performed with the LIMMA package implemented in
416 R/Bioconductor⁴⁸. The tiling array and whole-genome data sets were analysed separately.

417 Expression values were extracted and quantile-normalised, log₂-transformed, quality
418 controlled and analysed for differences in expression between individuals and wing regions.
419 P-values were adjusted for multiple hypotheses testing using the False Discovery Rate (FDR)
420 method⁴⁹.

421 *In situ hybridisations*

422 *Hn* and *Hm* larvae were reared in a greenhouse at 25-30°C and sampled at the last instar. In
423 situ hybridizations were performed according to previously described methods²⁵ with a *cortex*
424 riboprobe synthesized from a 831-bp cDNA amplicon from *Hn*. Wing discs were incubated in
425 a standard hybridization buffer containing the probe for 20-24 h at 60°C. For secondary
426 detection of the probe, wing discs were incubated in a 1:3000 dilution of anti-digoxigenin
427 alkaline phosphatase Fab fragments and stained with BM Purple for 3-6 h at room
428 temperature. Stained wing discs were photographed with a Leica DFC420 digital camera
429 mounted on a Leica Z6 APO stereomicroscope.

430 *De novo assembly of short read data in Hm and related taxa*

431 In order to better characterise indel variation from the short-read sequence data used for the
432 genotype-by-phenotype association analysis, we performed *de novo* assemblies of a subset of
433 *Hm* individuals and related taxa with a diversity of phenotypes (Extended Data Figure 2).
434 Assemblies were performed using the *de novo* assembly function of CLCGenomics
435 Workbench v.6.0 under default parameters. The assembled contigs were then BLASTed
436 against the *Yb* region of the *Hm melpomene* genome¹⁸, using Geneious v.8.0. The contigs
437 identified by BLAST were then concatenated to generate an allele sequence for each
438 individual. Occasionally two unphased alleles were generated when two contigs were
439 matched to a given region. If more than two contigs of equal length matched then this was

440 considered an unresolvable repeat region and replaced with Ns. The assembled alleles were
441 then aligned using the MAFFT alignment plugin in Geneious v.8.0.

442 ***Long-range PCR targeted sequencing of cortex in Hm aglaope and Hm amaryllis***

443 We generated two long-range PCR products covering 88.8% of the 1,344bp coding region of
444 *cortex* (excluding 67bp at the 5' end and 83bp at the 3' end, further details in Supplementary
445 Information). A product spanning coding exons 5 to 9 (the final exon) was obtained from 29
446 *Hm amaryllis* individuals and 29 *Hm aglaope* individuals; a product spanning coding exons 2
447 to 5 was obtained from 32 *Hm amaryllis* individuals and 14 *Hm aglaope*. In addition, a
448 product spanning exons 4 to 6 was obtained from 6 *Hm amaryllis* and 5 *Hm aglaope* that
449 failed to amplify one or both of the larger products. Long-range PCR was performed using
450 Extensor long-range PCR mastermix (Thermo Scientific) following manufacturers guidelines
451 with a 60°C annealing temperature in a 10-20µl volume. The product spanning coding exons
452 5 to 9 was obtained with primers HM25_long_F1 and HM25_long_R4 (see Supplementary
453 Table 2 for primer sequences); the product spanning coding exons 2 to 5 was obtained with
454 primers HM25_long_F4 and HM25_long_R2; the product spanning exons 4 to 6 was
455 obtained with primers 25_ex5-ex7_r1 and 25_ex5-ex7_f1. Products were pooled for each
456 individual, including 5 additional products from the *Yb* locus and 7 products in the region of
457 the *HmB/D* locus. They were then cleaned using QIAquick PCR purification kit (QIAGEN)
458 before being quantified with a Qubit Fluorometer (Life Technologies) and pooled in
459 equimolar amounts for each individual, taking into account variation in the length and
460 number of PCR products included for each individual (because of some PCR failures, ie.
461 proportionally less DNA was included if some PCR products were absent for a given
462 individual).

463 Products were pooled within individuals (including additional products for other genes not
464 analysed here) and then quantified and pooled in equimolar amounts for each individual
465 within each race. The pooled products for each race (*Hm aglaope* and *amaryllis*) were then
466 prepared as two separate libraries with molecular identifiers and sequenced on a single lane
467 of an Illumina GAIIX. Analysis was performed using Galaxy and the history is available at
468 <https://usegalaxy.org/u/njnadeau/h/long-pcr-final>. Reads were quality filtered with a
469 minimum quality of 20 required over 90% of the read, which resulted in 5% of reads being
470 discarded. Reads were then quality trimmed to remove bases with quality less than 20 from
471 the ends. They were then aligned to the target regions using the fosmid sequences from
472 known races⁴⁵ with sequence from the *Yb* BAC walk¹² used to fill any gaps. Alignments were
473 performed with BWA v0.5.6⁴² and converted to pileup format using Samtools v0.1.12 before
474 being filtered based on quality (≥ 20) and coverage (≥ 10). BWA alignment parameters were
475 as follows: fraction of missing alignments given 2% uniform base error rate (aln -n) 0.01;
476 maximum number of gap opens (aln -o) 2; maximum number of gap extensions (aln -e) 12;
477 disallow long deletion within 12 bp towards the 3'-end (aln -d); number of first subsequences
478 to take as seed (aln -l) 100. We then calculated coverage and minor allele frequencies for
479 each race and the difference between these using custom scripts in R⁵⁰.

480 ***Sequencing and analysis of Hm fosmid clones***

481 Fosmid libraries had previously been made from single individuals of 3 *Hm* races (*rosina*,
482 *amaryllis* and *aglaope*) and several clones overlapping the *Yb* interval had been sequenced⁴⁵.
483 We extended the sequencing of this region, particularly the region overlapping *cortex* by
484 sequencing an additional 4 clones from *Hm rosina* (1051_83D21, accession KU514430;
485 1051_97A3, accession KU514431; 1051_65N6, accession KU514432; 1051_93D23,
486 accession KU514433) 2 clones from *Hm amaryllis* (1051_13K4, accession KU514434;
487 1049_8P23, accession KU514435) and 3 clones from *Hm aglaope* (1048_80B22, accession

488 KU514437; 1049_19P15, accession KU514436; 1048_96A7, accession KU514438). These
489 were sequenced on a MiSeq 2000, and assembled using the *de novo* assembly function of
490 CLCGenomics Workbench v.6.0. The individual clones (including existing clones 1051-
491 143B3, accession FP578990; 1049-27G11, accession FP700055; 1048-62H20, accession
492 FP565804) were then aligned to the BAC and genome scaffold¹⁸ references using the
493 MAFFT alignment plugin of Geneious v.8.0. Regions of general sequence similarity were
494 identified and visualised using MAUVE⁵¹. We merged overlapping clones from the same
495 individual if they showed no sequence differences, indicating that they came from the same
496 allele. We identified transposable elements (TEs) using nBLAST with an insect TE list
497 downloaded from Repbase Update⁵² including known *Heliconius* specific TEs⁵³.

498 **5' RACE, RT-PCR and qRT-PCR**

499 All tissues used for gene expression analyses were dissected from individuals from captive
500 stocks derived from wild caught individuals of various races of *Hm* (*aglaope*, *amaryllis*,
501 *melpomene*, *rosina*, *plesseni*, *malleti*) and F2 individuals from a *Hm rosina* (female) x *Hm*
502 *melpomene* (male) cross. Experimental individuals were reared at 28°C-31°C. Developing
503 wings were dissected and stored in RNAlater (Ambion Life Technologies). RNA was
504 extracted using a QIAGEN RNeasy Mini kit following the manufacturer's guidelines and
505 treated with TURBO DNA-free DNase kit (Ambion Life Technologies) to remove remaining
506 genomic DNA. RNA quantification was performed with a Nanodrop spectrophotometer, and
507 the RNA integrity was assessed using the Bioanalyzer 2100 system (Agilent).

508 Total RNA was thoroughly checked for DNA contamination by performing PCR for EF1 α
509 (using primers ef1-a_RT_for and ef1-a_RT_rev, Table S2) with 0.5 μ l of RNA extract (50ng-
510 1 μ g of RNA) in a 20 μ l reaction using a polymerase enzyme that is not functional with RNA

511 template (BioScript, Bioline Reagents Ltd.). If a product amplified within 45 cycles then the
512 RNA sample was re-treated with DNase.

513 Single stranded cDNA was synthesised using BioScript MMLV Reverse Transcriptase
514 (Bioline Reagents Ltd.) with random hexamer (N6) primers and 1µg of template RNA from
515 each sample in a 20 µl reaction volume following the manufacturer's protocol. The resulting
516 cDNA samples were then diluted 1:1 with nuclease free water and stored at -80°C.

517 5' RACE was performed using RNA from hind-wing discs from one *Hm aglaope* and one
518 *Hm amaryllis* final instar larvae with a SMARTer RACE kit from Clontech (California,
519 USA). The gene specific primer used for the first round of amplification was anchored in
520 exon 4 (fzl_raceex5_R1, Supplementary Table 2). Secondary PCR of these products was then
521 performed using a primer in exon 2 (HM25_long_F2, Supplementary Table 2) and the nested
522 universal primer A. Other isoforms were detected by RT-PCR using primers within exons 2
523 and 9 (gene25_for_full1 and gene25_rev_ex3). We identified isoforms from 5' RACE and
524 RT-PCR products by cutting individual bands from agarose gels and if necessary by cloning
525 products before Sanger sequencing. Cloning of products was performed using TOPO TA
526 (Invitrogen) or pGEM-T (Promega) cloning kits. Sanger sequencing was performed using
527 BigDye terminator v3.1 (Applied Biosystems) run on an ABI13730 capillary sequencer.
528 Primers fzl_ex1a_F1 and fzl_ex4_R1 were used to confirm expression of the furthest 5'
529 UTR. For isoforms that appeared to show some degree of race specificity we designed
530 isoform specific PCR primers spanning specific exon junctions (Extended Data Fig 2, 4,
531 Supplementary Table 2) and used these to either qualitatively (RT-PCR) or quantitatively
532 (qRT-PCR) assess differences in expression between races.

533 We performed qRT-PCR using SensiMix SYBR green (Bioline Reagents Ltd.) with 0.2-
534 0.25µM of each primer and 1µl of the diluted product from the cDNA reactions. Reactions

535 were performed in an Opticon 2 DNA engine (MJ Research), with the following cycling
536 parameters: 95°C for 10min, 35-50 x: (95°C for 15sec, 55-60°C for 30sec, 72° for 30sec),
537 72°C for 5min. Melting curves were generated between 55°C and 90°C with readings taken
538 every 0.2°C for each of the products to check that a single product was generated. At least
539 one product from each set of primers was also run on a 1% agarose gel to check that a single
540 product of the expected size was produced and the identity of the product confirmed by direct
541 sequencing (See Supplementary Table 2 for details of primers for each gene). We used two
542 housekeeping genes (*EF1α* and *Ribosomal Protein S3A*) for normalisation and all results
543 were taken as averages of triplicate PCR reactions for each sample.

544 C_t values were defined as the point at which fluorescence crossed a threshold (R_{Ct}) adjusted
545 manually to be the point at which fluorescence rose above the background level.

546 Amplification efficiencies (E) were calculated using a dilution series of clean PCR product.

547 Starting fluorescence, which is proportional to the starting template quantity, was calculated

548 as $R_0 = R_{Ct} (1+E)^{-Ct}$. Normalized values were then obtained by dividing R_0 values for the

549 target loci by R_0 values for *EF1α* and *RPS3A*. Results from both of these controls were

550 always very similar, therefore the results presented are normalized to the mean of *EF1α* and

551 *RPS3A*. All results were taken as averages of triplicate PCR reactions. If one of the triplicate

552 values was more than one cycle away from the mean then this replicate was excluded.

553 Similarly any individuals that were more than two standard deviations away from the mean of

554 all individuals for the target or normalization genes were excluded (these are not included in

555 the numbers of individuals reported). Statistical significance was assessed by Wilcoxon rank

556 sum tests performed in R⁵⁰.

557 ***RNAseq analysis of *Hm amaryllis/aglaope****

558 RNA-seq data for hind-wings from three developmental stages had previously been obtained
559 for two individuals of each race at each stage (12 individuals in total) and used in the
560 annotation of the *Hm* genome¹⁸ (deposited in ENA under study accessions ERP000993 and
561 PRJEB7951). Four samples were multiplexed on each sequencing lane with the fifth instar
562 larval and day 2 pupal samples sequenced on a GAIIX sequencer and the day 3 pupal wings
563 sequenced on a HiSeq 2000 sequencer.

564 Two methods were used for alignment of reads to the reference genome and inferring read
565 counts, Stampy⁴³ and RSEM (RNAseq by Expectation Maximisation)⁵⁴. In addition we used
566 two different R/Bioconductor packages for estimation of differential gene expression,
567 DESeq⁵⁵ and BaySeq⁵⁶. Read bases with quality scores < 20 were trimmed with FASTX-
568 Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/index.html). Stampy was run with default
569 parameters except for mean insert size, which was set to 500, SD 100 and substitution rate,
570 which was set to 0.01. Alignments were filtered to exclude reads with mapping quality <30
571 and sorted using Samtools⁵⁷. We used the HT seq-count script in with HTseq⁵⁸ to infer counts
572 per gene from the BAM files.

573 RSEM⁵⁴ was run with default parameters to infer a transcriptome and then map RNAseq
574 reads against this using Bowtie³⁷ as an aligner. This was run with default parameters except
575 maximum number of mismatches, which was set to 3.

576 *Annotation and alignment of fizzy family proteins*

577 In the arthropod genomes, some fizzy family proteins were found to be poorly annotated
578 based on alignments to other family members. In these cases annotations were improved
579 using well annotated proteins from other species as references in the program Exonerate⁴¹
580 and the outputs were manually curated. Specifically, the annotation of *B. mori* *fzr* was
581 extended based on alignment of *D. plexippus* *fzr*; the annotation of *B. mori* *fzy* was altered

582 based on alignment of *Drosophila melanogaster* and *D. plexippus fzy*; *H. melpomene fzy* was
583 identified as part of the annotated gene HMEL017486 on scaffold HE671623 (Hmel v1.1)
584 based on alignment of *D. plexippus fzy*; the *Apis mellifera fzs* annotation was altered based
585 on alignment of *D. melanogaster fzs*; the annotation of *Acyrtosiphon pisum fzs* was altered
586 based on alignment of *D. melanogaster fzs*. No one-to-one orthologues of *D. melanogaster*
587 *fzs2* were found in any of the other arthropod genera, suggesting that this gene is *Drosophila*
588 specific. Multiple sequence alignment of all the fizzy family proteins was then performed
589 using the Expresso server⁵⁹ within T-coffee⁶⁰, and this alignment was used to generate a
590 neighbour joining tree in Geneious v8.1.7.

591 **Expression of *H. melpomene cortex* in *D. melanogaster* wings**

592 *D. melanogaster Cortex* is known to generate an irregular microchaete phenotype when
593 ectopically expressed in the posterior compartment of the adult fly wing²⁴. We performed the
594 same assay using *H. melpomene cortex* in order to test if this functionality was conserved.
595 Following the methods of Swan and Schüpbach²⁴ a UAS-GAL4 construct was created using
596 the coding region for the long isoform of *Hm cortex*, plus a *Drosophila cortex* version to act
597 as positive control. The HA-tagged *H. melpomene UAS-cortex* expression construct was
598 generated using cDNA reverse transcribed (Revert-Aid, Thermo-Scientific) from RNA
599 extracted (Qiagen RNeasy) from pre-ommochrome pupal wing material. An HA-tagged
600 *D.melanogaster UAS-cortex* version was also constructed, following the methods of Swan
601 and Schüpbach, (2007). Expression was driven by hsp70 promoter. Constructs were injected
602 into ϕ C31-attP40 flies (#25709, Bloomington stock centre, Indiana; Cambridge University
603 Genetics Department, UK, fly injection service) by site directed insertion into CII via an attB
604 site in the construct. Homozygous transgenic flies were crossed with w,y';en-GAL4;UAS-
605 GFP (gift of M. Landgraf lab, Cambridge University Zoology Department) to drive

606 expression in the engrailed posterior domain of the wing, and adult offspring wings
607 photographed (Extended Data Fig 6B-D). Expression of the construct was confirmed by IHC
608 (standard *Drosophila* protocol) of final instar larval wing discs using mouse anti-HA and goat
609 anti-mouse alexa-fluor 568 secondary antibodies (Abcam), imaged by Leica SP5 confocal.
610 Successful expression of *Hm_Cortex* was confirmed by IHC against an HA tag inserted at the
611 N terminal of either protein (Extended Data Fig 6E).

612

613 **References**

- 614 31. Supple, M. A. *et al.* Genomic architecture of adaptive color pattern divergence and
615 convergence in *Heliconius* butterflies. *Genome Res.* **23**, 1248–1257 (2013).
- 616 32. de la Bastide, M. & McCombie, W. R. Assembling genomic DNA sequences with
617 PHRAP. *Curr. Protoc. Bioinforma. Ed. Board Andreas Baxevanis Al Chapter 11,*
618 *Unit11.4* (2007).
- 619 33. Gordon, D., Abajian, C. & Green, P. Consed: a graphical tool for sequence finishing.
620 *Genome Res.* **8**, 195–202 (1998).
- 621 34. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
622 sequence data. *Bioinformatics* btu170 (2014). doi:10.1093/bioinformatics/btu170
- 623 35. Magoč, T. & Salzberg, S. L. FLASH: fast length adjustment of short reads to improve
624 genome assemblies. *Bioinformatics* **27**, 2957–2963 (2011).
- 625 36. Andrews, S. *FastQC*. (2011).
- 626 37. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient
627 alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).
- 628 38. Trapnell, C., Pachter, L. & Salzberg, S. L. TopHat: discovering splice junctions with
629 RNA-Seq. *Bioinformatics* **25**, 1105–1111 (2009).

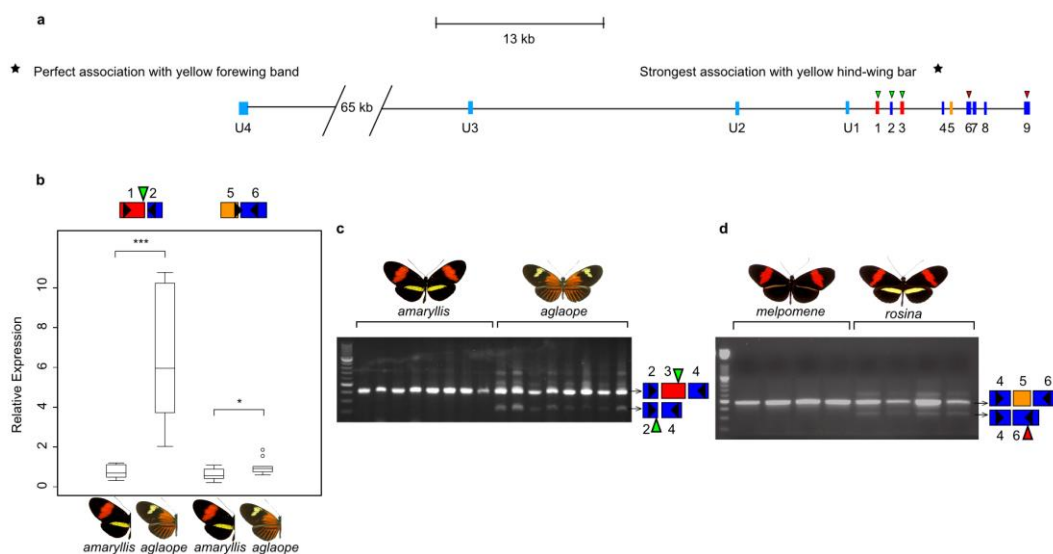
- 630 39. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals
631 unannotated transcripts and isoform switching during cell differentiation. *Nat.*
632 *Biotechnol.* **28**, 511–515 (2010).
- 633 40. Holt, C. & Yandell, M. MAKER2: an annotation pipeline and genome-database
634 management tool for second-generation genome projects. *BMC Bioinformatics* **12**, 491
635 (2011).
- 636 41. Slater, G. S. & Birney, E. Automated generation of heuristics for biological sequence
637 comparison. *BMC Bioinformatics* **6**, 31 (2005).
- 638 42. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
639 transform. *Bioinforma. Oxf. Engl.* **25**, 1754–1760 (2009).
- 640 43. Lunter, G. & Goodson, M. Stampy: A statistical algorithm for sensitive and fast mapping
641 of Illumina sequence reads. *Genome Res.* **21**, 936–939 (2011).
- 642 44. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-
643 generation DNA sequencing data. *Nat Genet* **43**, 491–498 (2011).
- 644 45. Nadeau, N. J. *et al.* Genomic islands of divergence in hybridizing *Heliconius* butterflies
645 identified by large-scale targeted sequencing. *Philos. Trans. R. Soc. B Biol. Sci.* **367**,
646 343–353 (2012).
- 647 46. Martin, S. H. *et al.* Genome-wide evidence for speciation with gene flow in *Heliconius*
648 butterflies. *Genome Res.* **23**, 1817–1828 (2013).
- 649 47. Aulchenko, Y. S., Ripke, S., Isaacs, A. & van Duijn, C. M. GenABEL: an R library for
650 genome-wide association analysis. *Bioinforma. Oxf. Engl.* **23**, 1294–1296 (2007).
- 651 48. Smyth, G. K. in *Bioinformatics and Computational Biology Solutions Using R and*
652 *Bioconductor* (eds. Gentleman, R., Carey, V. J., Huber, W., Irizarry, R. A. & Dudoit, S.)
653 397–420 (Springer New York, 2005).

- 654 49. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and
655 Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300
656 (1995).
- 657 50. R Development Core Team. *R: A language and environment for statistical computing.*
658 (R Foundation for Statistical Computing, 2011).
- 659 51. Darling, A. C. E., Mau, B., Blattner, F. R. & Perna, N. T. Mauve: Multiple Alignment of
660 Conserved Genomic Sequence With Rearrangements. *Genome Res.* **14**, 1394–1403
661 (2004).
- 662 52. Jurka, J. *et al.* Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet.*
663 *Genome Res.* **110**, 462–467 (2005).
- 664 53. Lavoie, C. A., Platt, R. N., Novick, P. A., Counterman, B. A. & Ray, D. A. Transposable
665 element evolution in *Heliconius* suggests genome diversity within Lepidoptera. *Mob.*
666 *DNA* **4**, 21 (2013).
- 667 54. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data
668 with or without a reference genome. *BMC Bioinformatics* **12**, 323 (2011).
- 669 55. Anders, S. & Huber, W. Differential expression analysis for sequence count data.
670 *Genome Biol.* **11**, 1–12 (2010).
- 671 56. Hardcastle, T. J. & Kelly, K. A. baySeq: Empirical Bayesian methods for identifying
672 differential expression in sequence count data. *BMC Bioinformatics* **11**, 422 (2010).
- 673 57. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinforma. Oxf. Engl.*
674 **25**, 2078–2079 (2009).
- 675 58. Anders, S., Pyl, P. T. & Huber, W. HTSeq - A Python framework to work with high-
676 throughput sequencing data. *bioRxiv* (2014). doi:10.1101/002824

- 677 59. Armougom, F. *et al.* Espresso: automatic incorporation of structural information in
 678 multiple sequence alignments using 3D-Coffee. *Nucleic Acids Res.* **34**, W604–608
 679 (2006).
- 680 60. Di Tommaso, P. *et al.* T-Coffee: a web server for the multiple sequence alignment of
 681 protein and RNA sequences using structural information and homology extension.
 682 *Nucleic Acids Res.* **39**, W13–17 (2011).

683

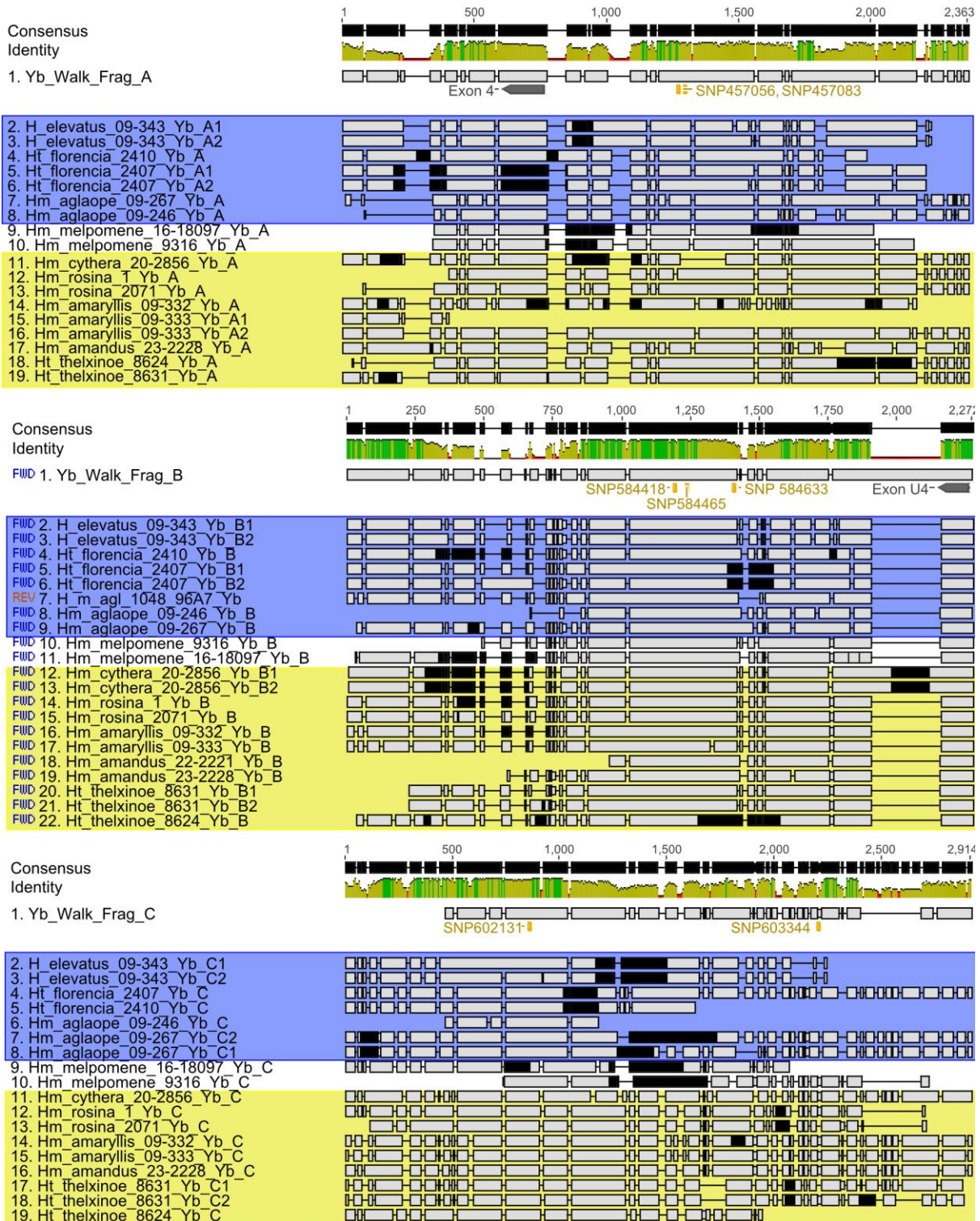
684

685 **Extended Data**

686

- 687 Extended Data Figure 1. A) Exons and splice variants of *cortex* in *Hm*. Orientation is
 688 reversed with respect to figures 2 and 4, with transcription going from left to right. SNPs
 689 showing the strongest associations with phenotype are shown with stars. B) Differential
 690 expression of two regions of *cortex* between *Hm amaryllis* and *Hm aglaope* whole hindwings
 691 (N=11 and N=10 respectively). Boxplots are standard (median; 75th and 25th percentiles;
 692 maximum and minimum excluding outliers – shown as discrete points) C) Expression of a

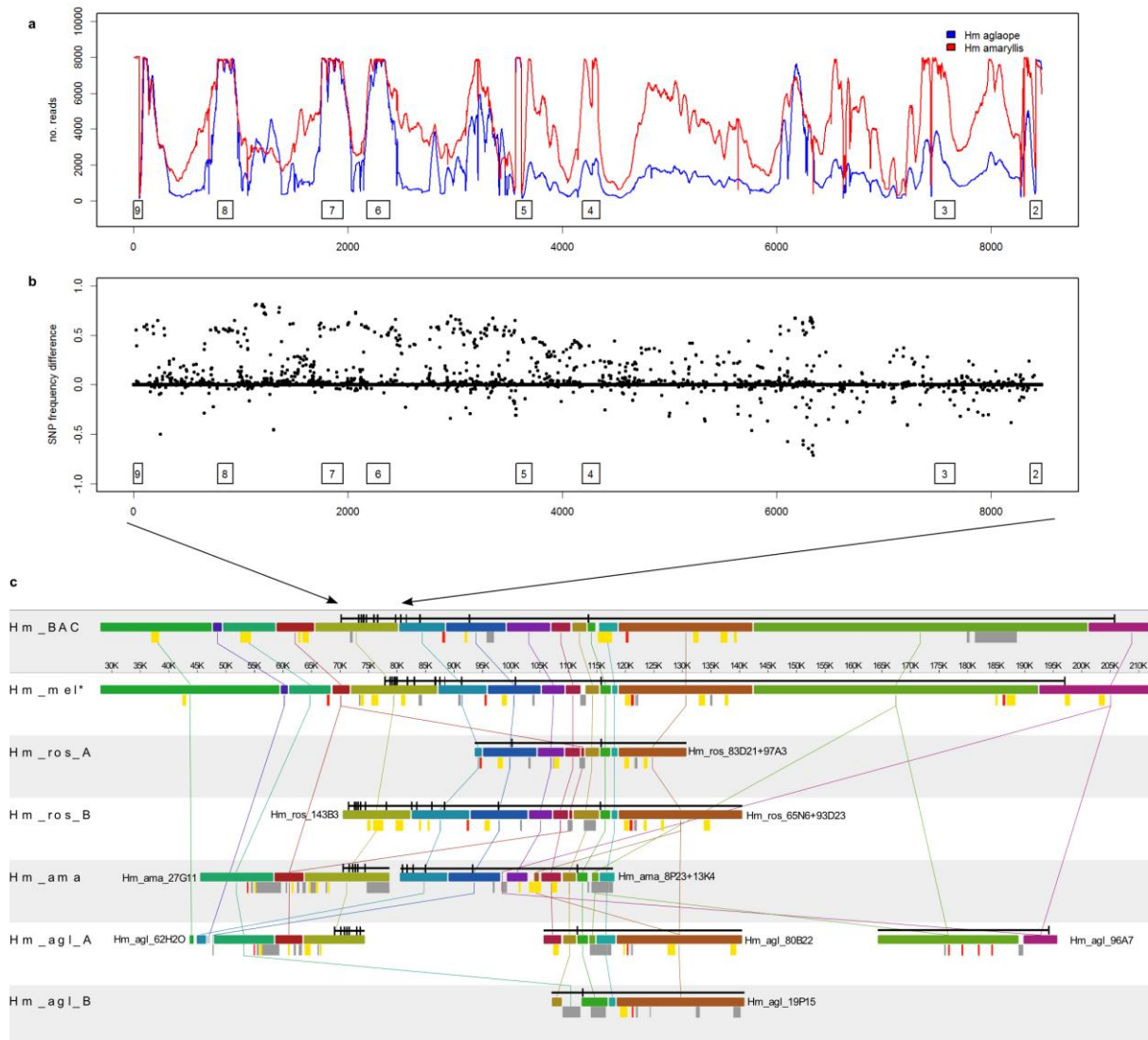
693 *cortex* isoform lacking exon 3 is found in *Hm aglaope* but not *Hm amaryllis* hindwings. D)
694 Expression of an isoform lacking exon 5 is found in *Hm rosina* but not *Hm melpomene*
695 hindwings. Green triangles indicate predicted start codons and red triangles predicted stop
696 codons, with usage dependent on which exons are present in the isoform. Schematics of the
697 targeted exons are shown for each (q)RT-PCR product, black triangles indicate the position
698 of the primers used in the assay.



699

700 Extended Data Figure 2. Alignments of *de novo* assembled fragments containing the top
 701 associated SNPs from *Hm* and related taxa short-read data. Identified indels do not show
 702 stronger associations with phenotype that those seen at SNPs (as shown in Extended Data
 703 Table 2), although some near-perfect associations are seen in fragment C. Black regions =

704 missing data; yellow box = individuals with a hindwing yellow bar; blue box = individuals
 705 with a yellow forewing band.



706

707 Extended Data Figure 3. Sequencing of long-range PCR products and fosmids spanning

708 *cortex*. A) Sequence read coverage from long-range PCR products across the *cortex* coding

709 region from 2 *Hm* races. B) Minor allele frequency difference from these reads between *Hm*

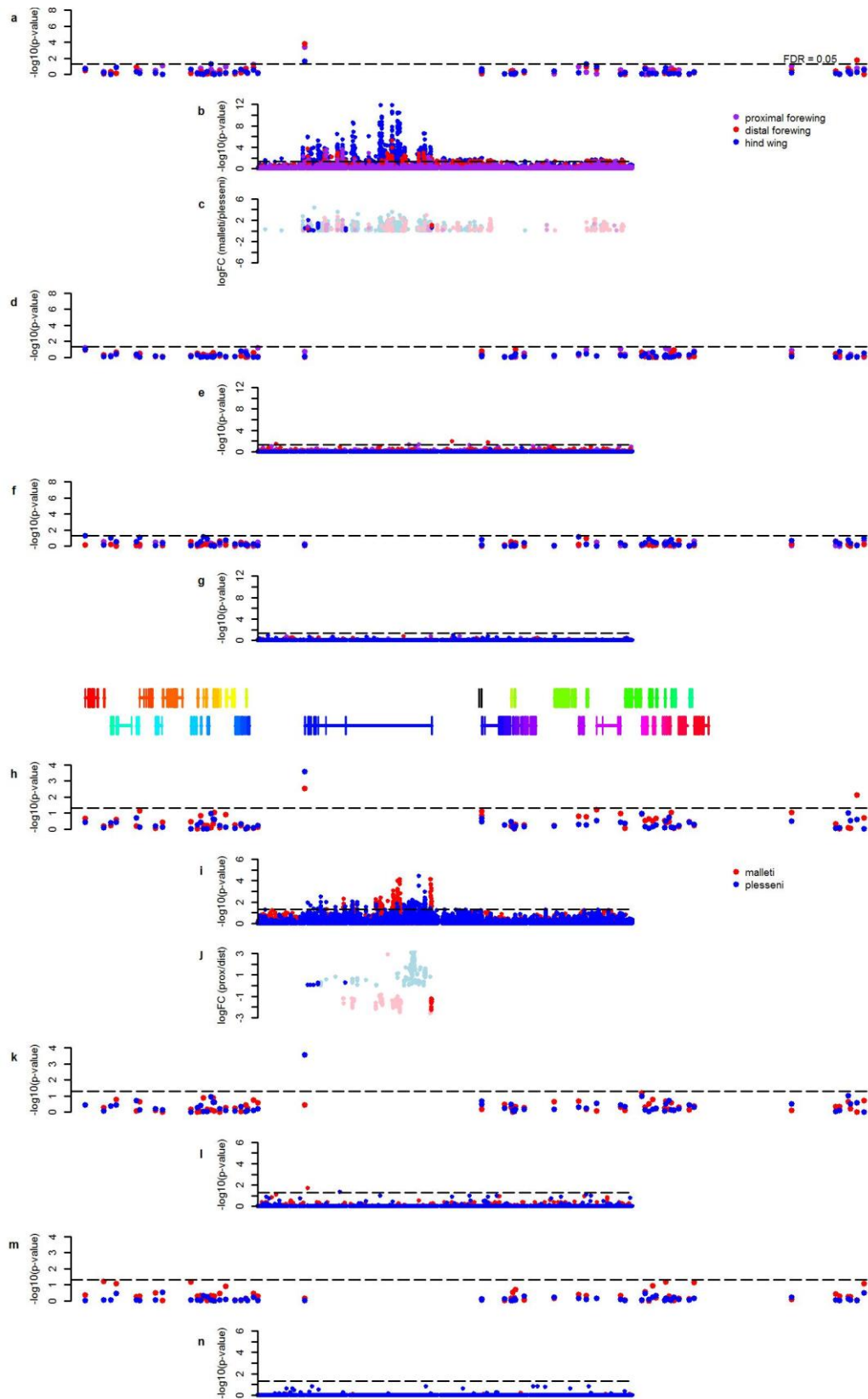
710 *aglaope* and *Hm amaryllis*. Exons of *cortex* are indicated by boxes, numbered as in Extended

711 Data Figure 2. C) Alignments of sequenced fosmids overlapping *cortex* from 3 *Hm*

712 individuals of different races. No major rearrangements are observed, nor any major

713 differences in transposable element (TE) content between closely related races with different

714 colour patterns (*melpomene/rosina* or *amaryllis/aglaope*). *Hm amaryllis* and *rosina* have the
715 same phenotype, but do not share any TEs that are not present in the other races. Hm_BAC =
716 BAC reference sequence, Hm_mel = *melpomene* from new unpublished assembly of *Hm*
717 genome⁵¹, Hm_ros = *rosina* (2 different alleles were sequenced from this individual),
718 Hm_ama = *amaryllis* (2 non-overlapping clones were sequenced in this individual), Hm_agla
719 = *aglaope* (4 clones were sequenced in this individual 2 of which represent alternative
720 alleles). Alignments were performed with Mauve: coloured bars represent homologous
721 genomic regions. *cortex* is annotated in black above each clone. Variable TEs are shown as
722 coloured bars below each clone: red = Metulj-like non-LTR, yellow = Helitron-like DNA,
723 grey = other.

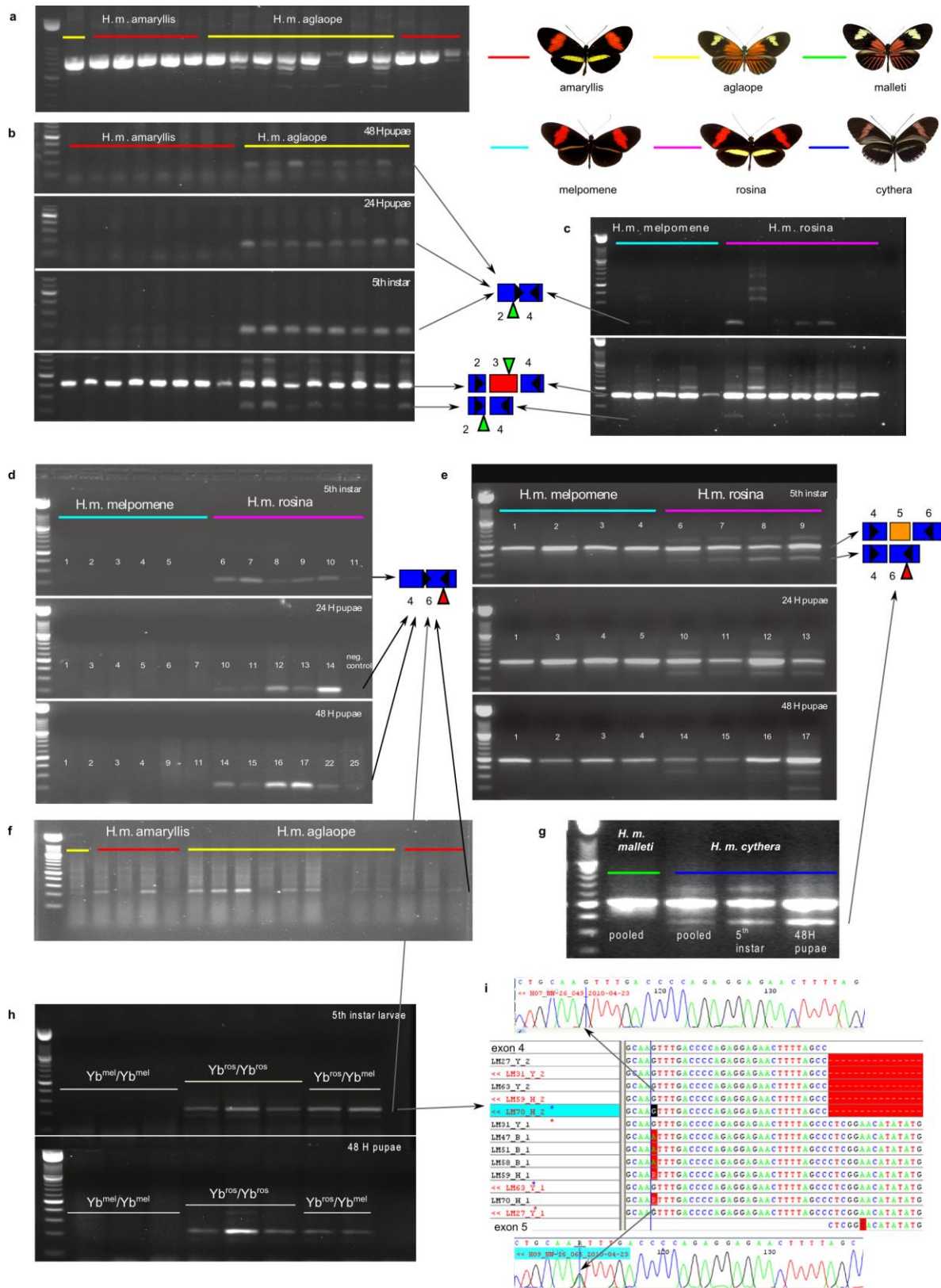


724

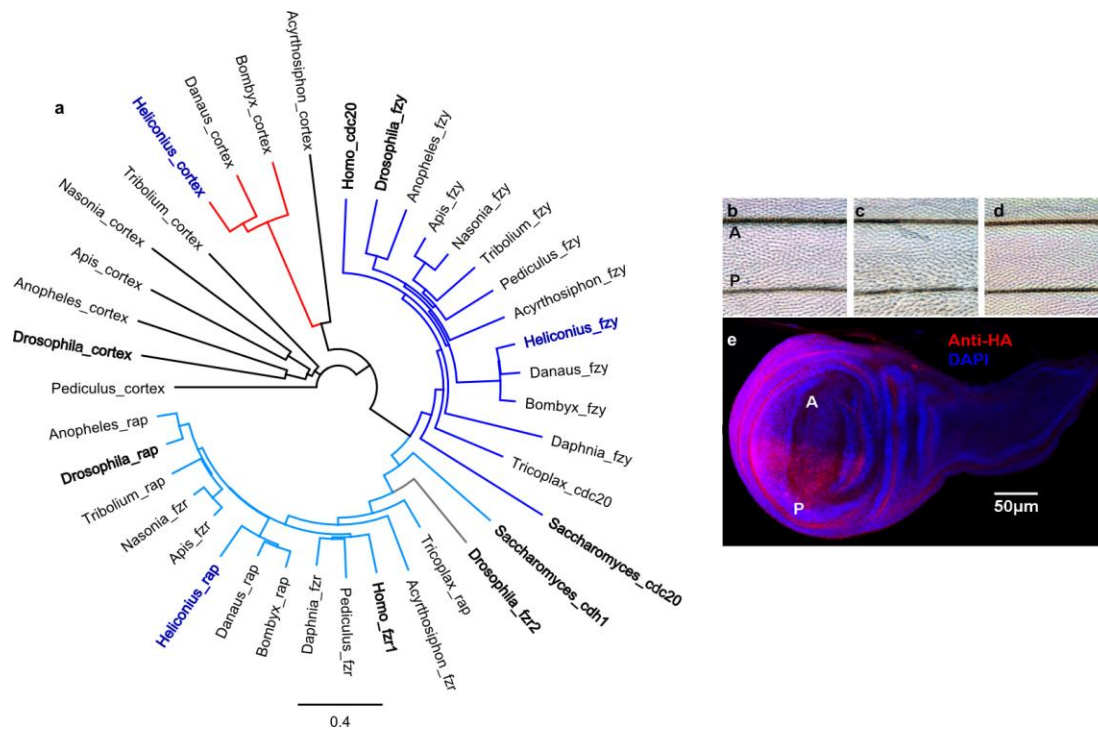
725 Extended Data Figure 4. Expression array results for additional stages, related to Figure 4. A-

726 G: comparisons between races (*H. m. plesse* and *H. m. malleti*) for 3 wing regions. H-N:

727 comparisons between proximal and distal forewing regions for each race. Significance values
728 ($-\log_{10}(\text{p-value})$) are shown separately for genes in the *HmYb* region from the gene array
729 (A,D,F,H,K,M) and for the *HmYb* tiling array (B,E,G,I,L,N) for day 1 (A,B,H,I), day 5
730 (D,E,K,L) and day 7 (F,G,M,N) after pupation. The level of expression difference (log fold
731 change) for tiling probes showing significant differences ($p \leq 0.05$) is shown for day 1 (C and
732 J) with probes in known *cortex* exons shown in dark colours and probes elsewhere shown as
733 pale colours.



737 Differences in splicing of exon 3 between *H. m. aglaope* and *H. m. amaryllis*. Products
738 amplified with a primer spanning the exon 2/4 junction at 3 developmental stages. The lower
739 panel shows verification of this assay by amplification between exons 2 and 4 for the same
740 final instar larval samples (replicated in Extended Data Figure 2C) C) Lack of consistent
741 differences between *H. m. melpomene* and *H. m. rosina* in splicing of exon 3. Top panel
742 shows products amplified with a primer spanning the exon 2/4 junction, lower panel is the
743 same samples amplified between exons 2 and 4. D) Differences in splicing of exon 5 between
744 *H. m. melpomene* and *H. m. rosina*. Products amplified with a primer spanning the exon 4/6
745 junction at 3 developmental stages. E) Subset of samples from D amplified with primers
746 between exons 4 and 6 for verification (middle, 24hr pupae samples are replicated in
747 Extended Data Figure 2D). F) Lack of consistent differences between *H. m. aglaope* and *H.*
748 *m. amaryllis* in splicing of exon 5. Products amplified with a primer spanning the exon 4/6
749 junction. G) *H. m. cythera* also expresses the isoform lacking exon 5, while a pool of 6 *H. m.*
750 *malleti* individuals do not. H) Expression of the isoform lacking exon 5 from an F2 *H. m.*
751 *melpomene* x *H. m. rosina* cross. Individuals homozygous or heterozygous for the *H. m.*
752 *rosina HmYb* allele express the isoform while those homozygous for the *H. m. melpomene*
753 *HmYb* allele do not. I) Allele specific expression of isoforms with and without exon 5.
754 Heterozygous individuals (indicated with blue and red stars) express only the *H. m. rosina*
755 allele in the isoform lacking exon 5 (G at highlighted position), while they express both
756 alleles in the isoform containing exon 5 (G/A at this position).



757

758 Extended Data Figure 6. Phylogeny of fizzy family proteins and effects of expressing *cortex*759 in the *Drosophila* wing. A) Neighbour joining phylogeny of Fizzy family proteins including760 functionally characterised proteins (in bold) from *Saccharomyces cerevisiae*, *Homo sapiens*761 and *Drosophila melanogaster* as well as copies from the basal metazoan *Trichoplax*762 *adhaerens* and a range of annotated arthropod genomes (*Daphnia pulex*, *Acyrthosiphon*763 *pisum*, *Pediculus humanus*, *Apis mellifica*, *Nasonia vitripennis*, *Anopheles gambiae*,764 *Tribolium castaneum*) including the lepidoptera *H. melpomene* (in blue), *Danaus plexippus*765 and *Bombyx mori*. Branch colours: dark blue, CDC20/fzy; light blue, CDH1/fzr/rap; red,766 lepidopteran cortex. B-E) Ectopic expression of *cortex* in *Drosophila melanogaster*.767 *Drosophila cortex* produces an irregular microchaete phenotype when expressed in the768 posterior compartment of the fly wing (C) whereas *Heliconius cortex* does not (D), when769 compared to no expression (B). A, anterior; P, posterior. Successful *Heliconius cortex*770 expression was confirmed by anti-HA IHC in the last instar *Drosophila* larva wing imaginal

771 disc (D, red), with DAPI staining in blue.

772 Extended Data Table 1. Genes in the *Yb* region and evidence for wing patterning control in

773 *Heliconius*

<i>Hm</i> gene ID	<i>He</i> gene ID	Putative gene name	<i>Heliconius melpomene</i>										<i>H. erato</i>			<i>Hn</i>		
			<i>Yb</i> ^l	<i>Sb</i> ^l	<i>A</i> ^{Yb}	<i>A</i> ^N	<i>E</i> ^l	<i>E</i> ^{9w}	<i>E</i> ^{9r}	<i>E</i> ^{hw}	<i>E</i> ^{lr}	<i>Cr</i> ^l	<i>A</i> ^{pat}	<i>A</i> ^{fav}	<i>P</i> ^l	<i>A</i> ^{bic}		
HM00002	HERA000036	Acylpeptide hydrolase			2									x				
HM00003	HERA000037	HM00003												x				
HM00004	HERA000038	Trehalase-1B	x											x				
HM00006	HERA000038.1	Trehalase-1A	x											x				
HM00007	HERA000039	B9 protein	x											x				
HM00008	HERA000040	HM00008	x		2									x				
HM00010	HERA000041	WD40 repeat domain 85	x											x				
HM00012	HERA000042	CG2519	x					x						x				
HM00013	HERA000045	Unkempt	x											x				
HM00014	HERA000046	Histone H3	x											x				
HM00015	HERA000047	HM00015	x											x				
HM00016	HERA000048	HM00016	x											x				
HM00017	HERA000049	RecQ Helicase	x											x				
HM00018	HERA000051	HM00018	x											x				
HM00019	HERA000052	BmSuc2	x					x						x				
HM00020	HERA000053	CG5796	x											x				
HM00021	HERA000054	HM00021	x											x				
HM00022	HERA000055	Enoyl-CoA hydratase	x											x				
HM00023	HERA000056	ATP binding protein	x											x				
HM00024	HERA000057	HM00024	x											x				
HM00025	HERA000059	cortex	x	x	56	74	x	x	x	603	1796	x	2	99	x	51		
HM00026	HERA000077	Poly(A)-specific ribonuclease (parr)		x	10					1	34	x					x	
HM00027	HERA000079	CG31320		x								x					x	
HM00028	HERA000080	ARP-like		x								x					x	
HM00029	HERA000081	CG4692		x								x					x	
HM00030	HERA000082	Proteasome 26S non ATPase subunit 4		x								x					x	
HM00031	HERA000083	HM00031		x					x			x					x	
HM00032	HERA000084	Zinc phosphodiesterase		x								1	x				x	
HM00033	HERA000085	Serine/threonine-protein kinase (LMTK1)		x								8	x				x	
HM00034	HERA000086	WD repeat domain 13 (Wdr13)			1	4						5	x				x	
HM00035	HERA000087	Domeless			1	2							x				x	
HM00036	HERA000061	WAS protein family homologue 1			5	36						37	x				x	
HM00038	HERA000062	Lethal (2) k05819 CG3054											x	2			x	
HM00039	HERA000064	Mitogen-activated protein kinase (MAPKK)											x				x	
HM00040	HERA000064.1	DNA excision repair protein ERCC-6											x				x	
HM00041	HERA000065	Penguin											x				x	
HM00042	HERA000066	Thymidylate kinase											x				x	
HM00043	HERA000067	Caspase-activated DNase											x				x	
HM00044	HERA000068	Regulator of ribosome biosynthesis											x				x	
HM00045	HERA000069	CG12659											x				x	
HM00046	HERA000070	CG33505											x				x	
HM00047	HERA000071	Sr protein											x				x	
HM00048	HERA000073	HM00048											x				x	
HM00049	HERA000073.1	HM00049											x				x	
HM00050	HERA000074	Shuttle craft											x				x	
HM00051	HERA000075	HM00051											x				x	
774	HM00052	HERA000076	HM00052					x					x				x	

775 *Yb*^l, within the previously mapped *Yb* interval¹². *Sb*^l, within the previously mapped *Sb*

776 interval¹². *Sb* controls a white/yellow hindwing margin and is not investigated in this study.

777 The *N* locus has not been fine-mapped previously. *A*^{Yb}, number of above background SNPs

778 associated with the hindwing yellow bar in this study. A^N , number of above background
779 SNPs associated with the forewing yellow band in this study. E^1 , detected as differentially
780 expressed between *Hm aglaope* and *amaryllis* from RNAseq data in this study
781 (Supplementary Information). E^{gw} , detected as differentially expressed between forewing
782 regions in the gene array in this study. E^{gr} , detected as differentially expressed between *Hm*
783 *plesseni* and *malleti* in in the gene array in this study. E^{tw} , numbers of probes showing
784 differential expression between forewing regions in the tiling array in this study. E^{tr} ,
785 numbers of probes showing differential expression between *Hm plesseni* and *malleti* in in the
786 tiling array in this study. Cr^I , within the previously mapped *HeCr* interval¹¹. A^{pet} , number of
787 SNPs fixed for the alternative allele in *He demophoon*. A^{fav} , number of SNPs fixed for the
788 alternative allele in *He favorinus*. P^I , within the previously mapped P interval¹³. A^{bic} , number
789 of above background SNPs associated with the *Hn bicoloratus* phenotype in this study.
790

791 Extended Data Table 2. Locations of fixed/above background SNPs and differentially
 792 expressed (DE) tiling array probes

		Positions of SNPs in the <i>He</i> and <i>Hn</i> association analyses								
		<i>cortex</i> coding exons	<i>cortex</i> UTR exons	<i>cortex</i> introns (nonTE)	<i>cortex</i> flanking intergenic (nonTE)	TEs	Other genes (exons or introns)	Other intergenic	Total	
<i>erato favorinus</i> fixed		2	0	96	8	2	0	0	108	
<i>erato demophoon</i> fixed		0	0	1	5	1	2	6	15	
<i>numata bicoloratus</i> above background		1	3	47	16	0	2	0	69	
		Positions of DE tiling array probes								
		Known <i>cortex</i> coding exons	<i>cortex</i> UTR exons	<i>cortex</i> introns (nonTE)	miRNAs	TEs	Other gene exons	Other introns/intergenic	Total	
Day3	malleti vs plesseni	Forewing proximal	8	7	323	0	13	1	7	359
		Forewing distal	12	2	327	0	8	0	8	357
		Hindwing	5	14	378	0	9	1	6	413
	Proximal vs distal	malleti	0	1	68	0	0	0	12	81
		plesseni	2	4	222	0	10	0	4	242
	Day1	malleti vs plesseni	Forewing proximal	1	0	22	0	3	0	7
Forewing distal			2	3	116	1	9	5	112	248
Hindwing			9	10	500	1	20	2	80	622
Proximal vs distal		malleti	0	12	95	0	1	0	0	108
		plesseni	3	3	81	0	99	0	0	186

793

794

795 Extended Data Table 3. SNPs showing the strongest phenotypic associations in the *H.*
 796 *melpomene/timareta/silvaniform* comparison.

Species	Race	Sample Code	SNP pos		SNP pos		SNP pos		SNP pos		SNP pos		SNP pos	
			HW 457083† bar (p=6.07E-10)	439063* (p=1.72E-09)	602131‡ (p=2.42E-09)	457056† (p=2.42E-09)	FW band	584465§ (p=1.37E-07)	584418§ (p=1.41E-07)	584633§ (p=2.10E-07)	603344‡ (p=2.19E-07)			
<i>H. melpomene</i>	<i>aglaope</i>	09-246	0	A/A	A/G	A/A	C/C	1	T/T	A/A	NA	T/T		
<i>H. melpomene</i>	<i>aglaope</i>	09-267	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. melpomene</i>	<i>aglaope</i>	09-268	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. melpomene</i>	<i>aglaope</i>	09-357	0	A/A	G/G	G/A	C/C	1	T/T	NA	C/C	T/T		
<i>H. melpomene</i>	<i>aglaope</i>	aglaope.1	0	A/A	G/G	NA	C/C	1	C/T	T/A	T/C	T/T		
<i>H. melpomene</i>	<i>amandus</i>	2221	1	A/A	NA	G/G	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>amandus</i>	2228	1	A/A	NA	G/G	C/C	0	C/T	T/A	T/C	A/A		
<i>H. melpomene</i>	<i>amaryllis</i>	09-332	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>amaryllis</i>	09-333	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>amaryllis</i>	09-075	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>amaryllis</i>	09-079	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>amaryllis</i>	amaryllis.1	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>bellula</i>	228	1	T/T	NA	G/G	T/T	0	C/C	T/T	T/T	NA		
<i>H. melpomene</i>	<i>bellula</i>	231	1	T/T	NA	G/A	T/T	0	C/T	T/A	T/C	NA		
<i>H. melpomene</i>	<i>cythera</i>	2856	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>cythera</i>	2857	1	NA	NA	NA	NA	0	NA	NA	NA	NA		
<i>H. melpomene</i>	<i>malleti</i>	17162	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. melpomene</i>	<i>melpomene</i>	melpomene18038	0	A/A	G/G	G/G	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>melpomene</i>	melpomene18097	0	NA	G/G	NA	C/C	0	C/C	T/T	T/T	NA		
<i>H. melpomene</i>	<i>melpomene</i>	melpomenem0.06	0	A/A	G/G	G/G	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>melpomene</i>	melpomenegen_ref	0	A/A	G/G	NA	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>melpomene</i>	melpomene13435	0	A/A	G/G	A/A	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>melpomene</i>	melpomene9315	0	A/A	G/G	A/A	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>melpomene</i>	melpomene9316	0	A/A	G/G	A/A	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>melpomene</i>	melpomene9317	0	A/A	G/G	A/A	C/C	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>plesseni</i>	9156	0	A/A	G/G	A/A	C/C	0	C/C	T/T	T/T	NA		
<i>H. melpomene</i>	<i>plesseni</i>	16293	0	A/A	G/G	A/A	C/C	0	C/C	T/T	T/T	NA		
<i>H. melpomene</i>	<i>rosina</i>	rosina.1	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>rosina</i>	2071	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>rosina</i>	531	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>rosina</i>	533	1	T/T	NA	G/G	T/T	0	C/C	T/T	T/T	NA		
<i>H. melpomene</i>	<i>rosina</i>	546	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. melpomene</i>	<i>thelxiopoeia</i>	13566	0	A/A	G/G	A/A	C/C	1	C/T	T/A	T/C	T/T		
<i>H. melpomene</i>	<i>vulcanus</i>	14632	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	NA		
<i>H. melpomene</i>	<i>vulcanus</i>	519	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. timareta</i>	<i>florencia</i>	2403	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>florencia</i>	2406	0	A/A	A/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>florencia</i>	2407	0	A/A	A/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>florencia</i>	2410	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>timareta</i>	8533	0	A/A	G/G	A/A	C/C	1	C/T	T/A	T/C	T/T		
<i>H. timareta</i>	<i>timareta</i>	9184	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>timareta</i>	8520	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>timareta</i>	8523	0	A/A	G/G	A/A	C/C	1	T/T	A/A	C/C	T/T		
<i>H. timareta</i>	<i>thelxinoe</i>	09-312	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. timareta</i>	<i>thelxinoe</i>	8624	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. timareta</i>	<i>thelxinoe</i>	8628	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. timareta</i>	<i>thelxinoe</i>	8631	1	T/T	A/A	G/G	T/T	0	C/C	T/T	T/T	A/A		
<i>H. elevatus</i>		09-343	0	A/T	G/G	A/A	T/T	1	C/T	NA	C/C	T/T		
<i>H. pardalinus</i>	<i>sergestus</i>	09-326	0	A/A	A/A	A/A	NA	0	C/C	T/T	T/T	NA		

797

798 *downstream of *cortex*, †between exons 3 and 4 of *cortex*, ‡upstream of *cortex*, §between799 exons U4 and U3 of *cortex*. None of these SNPs are within known TEs. Colours show

800 phenotypic associations: yellow = yellow hindwing bar; pink = no yellow hindwing bar;

801 green = yellow forewing band; blue = no yellow forewing band; grey = allele does not match
 802 expected pattern.

803

804 Extended Data Table 4. Transposable Elements (TEs) found within the *Yb* region.

Unique Occurrences					No.	TE name	Superfamily	Type
BAC	mel	ros	ama	agl				
1					1	BEL-1	BEL	LTR retrotransposon
					1	CR1-2	Jockey	Non-LTR retrotransposon
	1				1	Daphne-1	Jockey	Non-LTR retrotransposon
1					1	Daphne-6	Jockey	Non-LTR retrotransposon
1					1	DNA-like-8		DNA transposon
					1	Helitron-like-14	Helitron_A	DNA transposon
	1	2			4	Helitron-like-12	Helitron_A	DNA transposon
1	2				5	Helitron-like-12b	Helitron_A	DNA transposon
	1	1	1	1	7	Helitron-like-4a	Helitron_A	DNA transposon
						Helitron-like-4b	Helitron_A	DNA transposon
						Helitron-N2	Helitron_A	DNA transposon
					3	Helitron-like-7	Helitron_A	DNA transposon
5	3	3	1	2	16	Helitron-like-6a	Helitron_B	DNA transposon
						Helitron-like-6b	Helitron_B	DNA transposon
						Helitron-like-11	Helitron_B	DNA transposon
2	2	1		1	11	Helitron-like-15	Helitron_B	DNA transposon
6	5	3	1		18	Helitron-like-5	Helitron_B	DNA transposon
		1			2	Hmel_Unknown_50		
	1		1		2	Hmel_Unknown_174a/b		
	1				1	Hmel_Unknown_187b		
			1	1	2	Hmel_Unknown_230		
					1	Hmel_Unknown_234a		
					1	Hmel_Unknown_236a		
	1				1	Jockey-4	Jockey	Non-LTR retrotransposon
	1				1	LTR-3_gypsy	Gypsy	LTR retrotransposon
				1	1	Mariner-4	Mariner/Tc1	DNA transposon
1				3	29	Metulj-0	Metulj	SINE Non-LTR retrotransposon
						Metulj-1	Metulj	SINE Non-LTR retrotransposon
						Metulj-2	Metulj	SINE Non-LTR retrotransposon
						Metulj-3	Metulj	SINE Non-LTR retrotransposon
						Metulj-4	Metulj	SINE Non-LTR retrotransposon
						Metulj-5	Metulj	SINE Non-LTR retrotransposon
						Metulj-6	Metulj	SINE Non-LTR retrotransposon
						Metulj-7	Metulj	SINE Non-LTR retrotransposon
						nTc3-4	Mariner/Tc1	DNA transposon
						SINE-1	SINE	Non-LTR retrotransposon
1	1				2	nMar-3	Mariner/Tc1	DNA transposon
1					1	nMar-16	Mariner/Tc1	DNA transposon
			1		1	nMar-12/20	Mariner/Tc1	DNA transposon
				1	1	nPIF-3	PIF/Harbinger	DNA transposon
1					1	nTc3-2	Mariner/Tc1	DNA transposon
1					2	nTc3-3	Mariner/Tc1	DNA transposon
	1				2	R4-1	R2	Non-LTR retrotransposon
			1	1	6	Rep-1	REP	Non-LTR retrotransposon
2		1		1	4	RTE-3	RTE	Non-LTR retrotransposon
				1	2	RTE-11	RTE	Non-LTR retrotransposon
	1				3	Zenon-1	Jockey	Non-LTR retrotransposon
			1		1	Zenon-3	Jockey	Non-LTR retrotransposon

805