

1 **Electronic Supplementary Material**

3 **Pyrosequencing and sampling localities**

4 The data analyzed in the present paper comes from pyrosequencing of *Lycaeides* populations as
5 described in detail in Gompert *et al.* (2010). The “populations” involved in pyrosequencing each
6 consisted of 15 individuals from three sampling locations, the details of which are given in supplemental
7 table 1. Individuals within populations were pooled to create population samples that were processed
8 using restriction enzymes in a complexity-reduction step (van Orsouw *et al.* 2007; Baird *et al.* 2008) and
9 labeled with unique sequence barcodes prior to being sequenced on a 454 GS XLR70 Titanium platform
10 by SeqWright DNA Technology Services (Houston, TX, USA). Subsequent to pyrosequencing, barcodes
11 were used to associate individual sequence reads with populations. Barcodes were then removed from the
12 sequences and the SeqMan NGen assembler v.2.0.0 (DNASTAR) was used to assemble contigs; further
13 details on the parameterization of assembly are provided in Gompert *et al.* (2010), and the data are
14 available from the NCBI short read archive (accession number SRA010351). Custom Perl scripts were
15 used to select contigs for analyses in IM using the constraint (as discussed in the main text) of at least 3
16 sequences per population per contig for each pairwise comparison between populations.

18 **Implementation of the Isolation with Migration Model**

19 Here we discuss additional details in the implementation of the IM model using the program IMA,
20 and discuss assumptions of the model. Initial runs of IMA were conducted to find effective parameters for
21 the MCMC and to set priors that constrained searches, yet included the full range of reasonable parameter
22 space. We used the geometric increment model for 10 chains, with 0.8 specifying the non-linearity of
23 “heating” across chains, and 0.9 specifying the heating level for the highest numbered chain. Trend lines
24 and effective sample sizes were inspected to verify chain mixing and convergence of parameter estimates.
25 Mutations models implemented in IMA for sequence data (as opposed to microsatellites) are the infinite
26 sites model and HKY (Strasburg & Rieseberg 2010). The infinite sites model was not appropriate for our

27 data given the possibility of recurrent mutation and/or recombination, as suggested by the results of the
28 four gamete test (Hudson & Kaplan 1985) reported in supplemental table 2. The fraction of contigs that
29 failed the four gamete test ranged from 13 to 26%. These failures could be caused by either recurrent
30 mutation or recombination, but in any event are not apparently biased by population (and thus should not
31 introduce a bias to analyses). We note in particular that the variation among sets of contigs used for the
32 same population in different comparisons is similar to the variation across populations (e.g. *L. melissa*
33 east in the comparison with KBB had a failure rate of 14%, while 26% of contigs failed in the comparison
34 with *L. melissa* west). Other details involving molecular variation are shown in supplemental table 2.

35 As discussed in the main text (Materials & Methods), we do not have a genome-wide estimate for
36 mutation rate (μ) in *Lycaeides*. Therefore, we have not converted the parameters estimated by IMA into
37 values of years or individuals. Rather, values are left in the units estimated by the model: for example,
38 migration (m) is a ratio of m / μ (where m is the rate per gene per generation), and the population size
39 estimates (Θ) are $4 \times N_e$ (effective population size) $\times \mu$. Note that μ is per gene (genetic region or locus)
40 and is a geometric mean across loci, thus variation in fragment size could be important. However,
41 fragment sizes are not biased in any way across our populations (the average number of SNPs per contig
42 is quite similar across populations, as shown in supplemental table 2).

43 The IM model as implemented in IMA makes a number of assumptions, including that the genetic
44 regions being studied are not physically linked and are not under selection. Our samples were pooled at
45 the population level, thus we did not test for linkage disequilibrium. If any of our markers are linked, we
46 assume that the large number of markers compensates. Regarding the assumption of neutral variation, we
47 do not, at present, have a means for estimating the extent to which the genomes in question are under
48 selection (though the majority of our data come from non-coding regions (Gompert *et al.* 2010)). Outlier
49 analysis (Beaumont & Nichols 1996; Foll & Gaggiotti 2008) is not appropriate given the average number
50 of sequences in our contigs (see main text for details). However, this is not a limitation for our use of the
51 IM model, as we are interested in genetically-effective rates of gene flow in a comparative sense. If
52 selection acts against admixed individuals, and thereby reduces their contribution to a gene pool, this will

53 be reflected in lower rates of gene flow as estimated by the IM model. Because our data were generated
54 by pyrosequencing of pooled groups of individuals, it is possible that the data includes multiple sequences
55 of the same gene copy for a given individual. This is a violation of the IM model as these replicate
56 sequences should coalesce immediately. However, this problem should be minimal as we have relatively
57 few sequences compared to the number of sampled gene copies in each DNA pool.

58 Finally, the IM model assumes a bifurcating mode of diversification in which the populations
59 being analyzed are each other's closest relatives, and are not exchanging genes with other populations.
60 The history of *Lycaeides* is complex, and still being explored, but sister relationships for KBB and *L.*
61 *melissa*, and for *L. melissa* and *L. idas* are reasonable given the current state of knowledge (Nice *et al.*
62 2005, Gompert *et al.* 2008). Our pairwise comparisons likely violate the assumption of no gene flow with
63 other populations. This assumption is commonly violated in studies that employ the IM model, as few
64 natural pairs of populations will have no connections to other populations. However, simulations have
65 suggested that IM is robust to modest violations of most of the underlying assumptions (Strasburg &
66 Rieseberg 2010).

67
68 Baird, N. A., Etter, P. D., Atwood, T. S., Currey, M. C., Shiver, A. L., Lewis, Z. A., Selker, E. U., Cresko, W. A.
69 & Johnson, E. A. 2008 Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers.
70 *Plos One* **3**.

71 Beaumont, M. A. & Nichols, R. A. 1996 Evaluating loci for use in the genetic analysis of population
72 structure. *Proceedings of the Royal Society of London Series B-Biological Sciences* **263**, 1619-
73 1626.

74 Foll, M. & Gaggiotti, O. 2008 A Genome-Scan Method to Identify Selected Loci Appropriate for Both
75 Dominant and Codominant Markers: A Bayesian Perspective. *Genetics* **180**, 977-993.

76 Gompert, Z., Forister, M. L., Fordyce, J. A. & Nice, C. C. 2008 Widespread mito-nuclear discordance with
77 evidence for introgressive hybridization and selective sweeps in *Lycaeides*. *Molecular Ecology* **17**, 5231-
78 5244.

- 79 Gompert, Z., Forister, M. L., Fordyce, J. A., Nice, C. C., Williamson, R. J. & Buerkle, C. A. 2010 Bayesian
80 analysis of molecular variance in pyrosequences quantifies population genetic structure across the
81 genome of *Lycaeides* butterflies. *Molecular Ecology* **19**, 2455-2473.
- 82 Hudson, R. R., & Kaplan N. L. 1985 Statistical properties of the number of recombination events in the history of
83 a sample of DNA sequences. *Genetics* 111:147-164.
- 84 Nice, C. C., Anthony, N., Gelembiuk, G., Raterman, D. & Ffrench-Constant, R. 2005 The history and geography
85 of diversification within the butterfly genus *Lycaeides* in North America. *Molecular Ecology* **14**, 1741-
86 1754.
- 87 Strasburg, J. L. & Rieseberg, L. H. 2010 How Robust Are "Isolation with Migration" Analyses to Violations of
88 the IM Model? A Simulation Study. *Molecular Biology and Evolution* **27**, 297-310.
- 89 Suzek, B. E., Huang, H. Z., McGarvey, P., Mazumder, R. & Wu, C. H. 2007 UniRef: comprehensive and non-
90 redundant UniProt reference clusters. *Bioinformatics* **23**, 1282-1288.
- 91 van Orsouw, N. J., Hogers, R. C. J., Janssen, A., Yalcin, F., Snoeijers, S., Verstege, E., Schneiders, H., van der
92 Poel, H., van Oeveren, J., Verstegen, H. & van Eijk, M. J. T. 2007 Complexity Reduction of Polymorphic
93 Sequences (CRoPS (TM)): A Novel Approach for Large-Scale Polymorphism Discovery in Complex
94 Genomes. *Plos One* **2**.
- 95
96
97
98
99
100
101
102
103
104

105
 106
 107
 108
 109
 110
 111
 112
 113
 114
 115
 116
 117
 118
 119
 120
 121
 122
 123
 124

Supplemental table 1. Locality information for populations studied.

Taxon	Locality	Latitude	Longitude
KBB	Fort McCoy, WI	43°47'59''N	90°49'59''W
	Necedah, WI	44°04'00''N	90°11'20''W
	Saratoga, NY	43°03'24''N	73°48'47''W
<i>L. melissa</i> east	Brandon, SD	43°36'29''N	96°34'39''W
	Victor, ID	43°39'32''N	111°06'41''W
	Indian Bathtubs, WY	41°12'14''N	106°46'17''W
<i>L. melissa</i> west	Garderville, NV	38°48'54''N	119°46'44''W
	Verdi, NV	39°03'01''N	119°55'48''W
	Sierra Valley, CA	39°37'48''N	120°21'40''W
<i>L. idas</i>	Yuba Gap, CA	39°29'34''N	120°35'39''W
	Leek Springs, CA	38°37'59''N	120°14'24''W
	Trap Creek, CA	39°22'43''N	120°40'27''W

Supplemental table 2. Details of data used in analyses. “Comparison” here refers to the pairs of populations involved in IMA analyses, as in table 1. Results from the four gamete test are shown in the “Number failed out of total” column, with number referring to the number of contigs.

Comparison	Total base pairs	Population	Number failed out of total (percent)	SNPs	Proportion variable sites	Average SNPs per contig (SE)	GC content
KBB and <i>L. melissa</i> east	91944	KBB	73 / 317 (23%)	2521	0.0274	7.95 (0.48)	49.7%
		<i>L. melissa</i> east	45 / 317 (14%)	2233	0.0243	7.04 (0.45)	48.3%
<i>L. melissa</i> west and <i>L. idas</i>	60362	<i>L. melissa</i> west	25 / 188 (13%)	1229	0.0204	6.54 (0.52)	47.2%
		<i>L. idas</i>	44 / 188 (23%)	1621	0.0269	8.62 (0.63)	47.5%
<i>L. melissa</i> west and <i>L. melissa</i> east	71031	<i>L. melissa</i> west	36 / 236 (15%)	1512	0.0213	6.41 (0.45)	47.9%
		<i>L. melissa</i> east	62 / 236 (26%)	2100	0.0296	8.90 (0.64)	50.0%