

Supplementary Material

The conundrum of species delimitation: a genomic perspective on a mitogenetically super-variable butterfly

Vlad Dincă, Kyung Min Lee, Roger Vila, Marko Mutanen

DOI: 10.1098/rspb.2019.1311

Material and Methods

Dataset used for molecular analyses

The core dataset was based on 93 specimens of *M. didyma* for which both COI sequences and ddRADseq data were obtained (Tables S1-S3). All ddRADseq data, as well as 77 COI sequences, were generated for this study. To this dataset we added two specimens as outgroup (*Melitaea trivia* and *Melitaea deione*) (Leneveu et al. 2009). We followed Pazhenkova & Lukhtanov (2016) to assign the 93 specimens to mtDNA lineages (Fig. S1). For this purpose, as well as to obtain rough estimates of divergence events, we assembled a dataset of 347 COI sequences obtained by combining the 93 COI sequences with data used by two recent studies focused on the *M. didyma* complex (Pazhenkova et al. 2015, Pazhenkova & Lukhtanov 2016). These studies also incorporated sequences originating from Wahlberg & Zimmermann (2000), Vila & Bjorklund (2004), Leneveu et al. (2009), Hausmann et al. (2011), Ashfaq et al. (2013) and Dincă et al. (2011, 2015). Sequence FJ663794 was not used because of the following discrepancy: our preliminary analyses showed that it clusters with taxon *sutschana*, but it was recovered within *latonigena* by Pazhenkova & Lukhtanov (2016).

Mitochondrial DNA sequencing

Thirty-six of the COI sequences generated for this study were obtained at the Butterfly Diversity and Evolution Lab of the Institut de Biologia Evolutiva (CSIC-UPF), Barcelona, Spain. In this case, total genomic DNA was extracted using Chelex 100 resin, 100–200 mesh, sodium form (Biorad), under the following protocol: one leg was removed and introduced into 100 µl of Chelex 10% and 5 µl of Proteinase K (20 mg/ml) were added. The samples were incubated overnight at 55°C and were subsequently incubated at 100°C for 15 minutes. Samples were then centrifuged for 10 s at 3000 rpm. A 658-bp fragment near the 5' end of the mitochondrial gene COI was amplified by polymerase chain reaction (PCR) (primers used and PCR protocol available in Table S4). PCR products were purified and sequenced by Macrogen Inc.

The remaining sequences were generated at the Biodiversity Institute of Ontario, Canada following standard protocols for DNA barcoding (deWaard et al. 2008), and DNA sequencing was performed on an ABI 3730XL capillary sequencer (Applied Biosystems).

Sequences were edited in CodonCode Aligner 3.0 or in GENEIOUS PRO 6.1.8 (Biomatters, <http://www.geneious.com/>) and assembled using the latter.

The 79 COI sequences generated for this study (77 *M. didyma* plus two outgroup specimens) are available in GenBank (see Table S1 for accession numbers), and are also publicly available in the dataset DS-DIDYMA (dx.doi.org/10.5883/DS-DIDYMA) from the Barcode of Life Data Systems (<http://www.boldsystems.org/>).

Analyses of mitochondrial DNA sequences

Phylogenetic relationships for the full dataset (347 COI sequences) were inferred using Bayesian inference (BI) through the CIPRES Science Gateway (Miller et al. 2010). Both BI analyses and the estimation of node ages were run in BEAST 1.8.0 (Drummond & Rambaut 2007). The GTR + I + G substitution model was chosen according to the value of the Akaike information criterion (AIC) obtained in JMODELTEST 2.1.3 (Darriba et al. 2012). Base frequencies were estimated, six gamma rate categories were selected and a randomly generated initial tree was used.

Rough estimates of node ages were obtained by applying two molecular clocks with: 1.5% uncorrected pairwise distance per million years estimated for various invertebrates (Quek et al. 2004), and 2.3% estimated for the entire mitochondrial genome of several arthropods (Brower 1994). A lognormal relaxed clock and a normal prior distribution were used, centred on the mean between the two substitution rates, and the standard deviation was tuned so that the 95% confidence interval of the posterior density coincided with the 1.5% and 2.3% rates, respectively. Parameters were estimated using two independent runs of 40 million generations each, and convergence was checked using the program TRACER 1.6.

For the core dataset of 93 *M. didyma* COI sequences (and two outgroup samples) (i.e. those specimens for which ddRADseq data were also available) (Tables S1-S3), phylogenetic relationships were inferred using maximum likelihood (ML), to directly compare results with ML analyses based on ddRADseq data. The COI ML tree was inferred in RAxML v.8.2.0 (Stamatakis 2014) with bootstrap support estimated by a 1,000 replicates rapid-bootstrap analysis from the unpartitioned GTR+CAT model. We visualized the resulting phylogeny and assessed bootstrap support using FigTree v.1.4.2 (Rambaut 2015). The tree was rooted with *M. deione* (Leneveu et al. 2009).

ddRADseq library preparation and bioinformatics

In order to proceed with the ddRAD library preparation, genomic DNA (gDNA) was extracted from one or two legs using the DNeasy Blood & Tissue Kit (Qiagen). The quantity of gDNA extracts was checked using PicoGreen kit (Molecular Probes). To reach sufficient gDNA quantity and quality, whole genome amplification was performed using REPLI-g Mini Kit (Qiagen) due to its low concentrations of gDNA in the original extracts. The average original DNA extracts were 8.5 ng/μl and amplified up to 89.6 ng/μl (54.9 ng/μl on average) after whole genome amplification. The ddRADseq library was implemented following

protocols described in Lee *et al.* (2018) with an exception: the size distribution and concentration of the pools were measured with Bioanalyzer (Agilent Technologies). The demultiplexed *Melitaea* fastq data are archived in the NCBI SRA: SRP144304.

Raw paired-end reads were demultiplexed with no mismatches tolerated using their unique barcode and adapter sequences using *ipyrad* v.0.7.23 (Eaton and Overcast 2016). The quality of raw demultiplexed reads was checked with FastQC software (available at <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). The demultiplexed paired-reads were run through PEAR (Zhang et al. 2014) using default setting to merge overlapping reads, and input into the *ipyrad* pipeline. All *ipyrad* defaults were used, with the following exceptions: the minimum depth at which majority rule base calls are made was set to 3, the cluster threshold was set to 0.90, the minimum number of samples that must have data at a given locus for it to be retained was set to 4, 20, 30, 60, and 70, and the assembly method was set to denovo, denovo–reference, and reference for independent testing. The ‘denovo’ method assembles sequences without any reference resources. Homology is inferred during alignment clustering by sequence similarity using the program *vsearch* (<http://github.com/torognes/vsearch>). The ‘denovo–reference’ method was used to exclude sequences that mapped to the *Melitaea cinxia* mitochondrion genome (GenBank accession CM002851). The ‘reference’ assembly method maps sequences to *M. cinxia* whole genome sequences (GenBank, GCA_00071638) with *BWA* using the default *bwa-mem* setting (Li 2013) based on 90% of sequence similarity.

Phylogenetic analysis of ddRADseq data

To study the phylogenetic relationships among taxa and to test the validity of prevailing species hypotheses, we conducted ML analyses. ML trees were inferred in RAxML v.8.2.0 (Stamatakis 2014) for the concatenated RAD data, with bootstrap support estimated by a 1,000 replicates rapid-bootstrap analysis from the unpartitioned GTR+CAT model. We visualized the resulting phylogeny and assessed bootstrap support using FigTree v.1.4.2 (Rambaut 2015).

The unlinked SNP datasets for species tree construction were imported into BEAUti (Bouckaert et al. 2014), where the data were prepared for analyses with the SNAPP v.1.1.16 plugin (Bryant et al. 2012) in BEAST v.2.1.3 (Bouckaert et al. 2014). The priors for forward (u) and reverse (v) mutation rates were set to be estimated and the remaining parameters were left at default values, e.g. species divergence rate $\lambda = 0.00765$, and θ defined by a γ prior with shape parameter $\alpha = 11.750$ and scale parameter $\beta = 109.73$. Runs were carried out for 10 million generations, sampling every 1,000 generations. The output was inspected with Tracer v.1.6 (Rambaut et al. 2013), ensuring stationarity and ESS > 200 for all parameters with a few exceptions for single theta values of internal branches in single replicates. We then visualized the posterior distribution of species trees produced using DensiTree v.2.2.1 (Bouckaert 2010).

Population structure and admixture

We inferred population clustering with admixture from SNP frequency data to better visualize genomic variation between individuals with STRUCTURE v.2.3.1 (Pritchard et al. 2000). Ten replicates were run with each value of *K* between 1 and 8. Each run had a burn-in

of 20K generations followed by 200K generations of sampling. Replicates were permuted using CLUMPP (Jakobsson and Rosenberg 2007) and the optimal K value inferred using StructureHarvest (Earl and VonHoldt 2012) according to the ad hoc ΔK statistics (Evanno et al. 2005), which is the second-order rate of change of the likelihood function. STRUCTURE results were visualized using DISTRUCT (Rosenberg 2004).

FineRADstructure was used to investigate the genetic structure at population level within the *M. didyma* complex (Malinsky et al. 2018). The package includes RADpainter, a program designed to infer the co-ancestry matrix and estimate the number of populations within the dataset. The input file used was an allele.loci matrix (20% of missing data) generated by *ipyrad* program. The allele data was converted using a python script available at <https://github.com/edgardomortiz/fineRADstructure-tools> (last accessed April 20, 2018). Then, the individuals were assigned to populations and the phylogenetic tree was built using the fineSTRUCTURE MCMC clustering algorithm.

TreeMix was used to identify patterns of divergence and admixtures, testing for migration events ranging from one to five including possible admixture with two closely related outgroup samples included in the analyses (Pickrell and Pritchard 2012). This method constructs a bifurcating tree of populations using 100 bootstrap replicates and it identifies potential episodes of gene flow from the residual covariance matrix. This analysis was applied to a subset of 27 specimens that were also used for D-statistics (see below). This subset consisted of five representative samples from each of the five lineages, as well as two outgroups (*Melitaea trivia* and *M. deione*), that were chosen in order to equalize the sample size based on the highest number of recovered loci in the final data matrix (Table S2).

We used four-taxon D-statistics (Durand et al. 2011) to distinguish introgression from incomplete lineage sorting based on the same subset of 27 specimens used for *TreeMix*. The full dataset could not be used due to computational limitations. The test is based on the assumption of a true four-taxon asymmetric phylogeny (((P1, P2) P3,) O). All sites considered in the alignment of sequences from these taxa must be either mono or biallelic, with the outgroup defining the ancestral state 'A' relative to the derived state 'B'. If two alleles are present in a site, the possible combinations are ABBA and BABA. The D-statistics compares the occurrence of these two discordant site patterns, representing sites where an allele is derived in P3 relative to outgroup (O), and is derived in one but not both of the sister lineages P1 and P2. These discordant sites can arise through the sorting of ancestral polymorphisms. In absence of introgression, the frequencies for these two outcomes are expected to be equal. This finding would support incomplete lineage sorting (ILS) being responsible for barcode sharing, while deviation from it would support introgression (Durand et al. 2011). For the test, 1,000 bootstrap replicates were performed to measure the standard deviation of the D-statistics. Significance was evaluated by converting the Z-score (which represents the number of standard deviations from zero from D-statistics) into two tailed P-values, and using $\alpha=0.01$ as a conservative cutoff for significance after correcting for multiple comparisons using Holm-Bonferroni correction. All D-statistics were calculated in pyRAD

v.3.0.64 (Eaton 2014). In order to run interactive data analysis, the Python Jupyter notebooks (<https://jupyter.org>) were used.

The python script that we applied for D-statistics has been uploaded and shared via Dryad (DOI: <http://doi.org/10.5061/dryad.b883mf8>).

Pairwise F_{ST} values were calculated using Arlequin v.3.5 (Excoffier and Lischer 2010). Statistical significance of the F_{ST} values was tested by permutation analysis with 1,000 permutations. The proportion of missing data was calculated using Mesquite (Maddison and Maddison 2017).

Coalescent-based species delimitation with Bayes factors

We performed Bayes factor species delimitation using the BFD* method (Leaché et al. 2014) based on a subset of specimens (see Table S2) plus 2 outgroup specimens, due to computational limitations. This method allows for the comparison of alternative species delimitation models in an explicit multispecies coalescent framework using genome-wide SNP data, implemented using SNAPP (Bryant et al. 2012). We tested nine (when assuming 5 taxa) and ten competing species models (when assuming 8 taxa) for *M. didyma*: 37 specimens (including 2 outgroup specimens) when assuming 5 taxa; 41 specimens (including 2 outgroup specimens) when assuming 8 taxa. The specimens were selected based on the highest number of recovered loci in the final data matrix. The full dataset could not be used due to computational limitations. For all species models, we conducted path sampling for a total of 24 steps (200,000 MCMC steps, 10,000 burn-in steps each) to calculate marginal likelihood estimates (MLE) for each competing model. Bayes factor (BF) support was compared between models to identify the best-supported species model. We assessed the strength of support of alternative species delimitation models following the scale of Kass & Raftery (1995).

When assuming 5 taxa, that hypothesis was best supported, but when assuming 8 taxa, that hypothesis was best supported (Table S6). The results should be interpreted with caution as there is recent evidence suggesting that hypotheses considering more species/lineages are better supported than those assuming less splitting (e.g. Sukumaran & Knowles 2017; O'Connell & Smith 2018). Only few studies have examined the effect of missing data when using BFD*, which does not accommodate missing data between assigned species. The inclusion of more loci, even at the expense of very high amounts of missing data, led to higher BF and better resolved species trees than datasets with less missing data but fewer loci. This is because less stringent filtering retains lineage-specific loci, which may help coalescent methods to better delimit lineages. Therefore, caution should be taken in interpreting BFD* results on the basis of different levels of missing data.

Wolbachia infection analyses

All 95 specimens for which COI and ddRADseq data were available were surveyed for the presence of the bacterium *Wolbachia* (Table S1).

The presence of *Wolbachia* was tested using PCR and sequencing primers specific to *Wolbachia* genes *wsp* and *ftsZ*, which are extensively used to detect *Wolbachia* infection in a wide array of insects (Baldo et al. 2006). Primers used and PCR protocols are available in Table S4.

Samples with amplicons of the expected size (as visualized on agarose gels) were scored as positive for *Wolbachia* and PCR products of all infected specimens were sequenced. Sequences were then compared to existing records using the *Wolbachia* MLST Database (pubmlst.org/wolbachia/) in order to identify the sequence type for each gene locus. Sequences obtained during the screening are available in GenBank (see Table S1 for accession numbers) and in the dataset DS-DIDYMA (dx.doi.org/10.5883/DS-DIDYMA) from the Barcode of Life Data Systems (<http://www.boldsystems.org/>).

References

- Ashfaq M, Akhtar S, Khan AM, Adamowicz SJ, Hebert PDN. 2013. DNA barcode analysis of butterfly species from Pakistan points towards regional endemism. *Molecular Ecology Resources* **13**, 832–843. (doi: 10.1111/1755-0998.12131)
- Baldo L, Dunning Hotopp JC, Jolley KA, Bordenstein SR, Biber SA, Choudhury RR, Hayashi C, Maiden MC, Tettelin H, Werren JH. 2006. Multilocus Sequence Typing System for the Endosymbiont *Wolbachia pipientis*. *Applied and Environmental Microbiology* **72**(11), 7098–7110. (doi:10.1128/AEM.00731-0)
- Bouckaert R. 2010. DensiTree: making sense of sets of phylogenetic trees. *Bioinformatics*, **26**, 1372–1373.
- Bouckaert R, Heled J, Kühnert D *et al.* 2014. BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Computational Biology*, **10**, e1003537.
- Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, Suchard MA, Rambaut A, Drummond AJ. 2014. BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Computational Biology*, **10**, e1003537. (<https://doi.org/10.1371/journal.pcbi.1003537>)
- Brower AVZ. 1994. Rapid morphological radiation and convergence among races of the butterfly *Heliconius erato* inferred from patterns of mitochondrial DNA evolution. *PNAS* **91**, 6491–6495. (<https://doi.org/10.1073/pnas.91.14.6491>)
- Bryant D, Bouckaert R, Felsenstein J, Rosenberg NA, Roychoudhury A. 2012. Inferring species trees directly from biallelic genetic markers: Bypassing gene trees in a full coalescent analysis. *Molecular Biology and Evolution*, **29**, 1917–1932.
- Darriba D, Taboada GL, Doallo R, Posada D. 2012. jModelTest 2: more models, new heuristics and parallel computing. *Nature Methods* **9**, 772. (doi:10.1038/nmeth.2109)
- deWaard JR, Ivanova NV, Hajibabaei M, Hebert PDN. 2008. Assembling DNA Barcodes: Analytical Protocols. In *Methods in Molecular Biology: Environmental Genetics* (ed. Cristofre M.), pp. 275–293. Totowa, USA: Humana Press Inc.
- Dincă V, Zakharov EV, Hebert PDN, Vila R. 2011. Complete DNA barcode reference library for a country's butterfly fauna reveals high performance for temperate Europe. *Proceedings of the Royal Society B* **278**, 347–355. (DOI: 10.1098/rspb.2010.1089)
- Dincă V, Montagud S, Talavera G, Hernández-Roldán J, Munguira ML, García-Barros E, Hebert PDN, Vila R. 2015. DNA barcode reference library for Iberian butterflies enables a continental-scale preview of potential cryptic diversity. *Scientific Reports* **5**, 12395. (doi:10.1038/srep12395)
- Drummond AJ, Rambaut A. 2007. BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evolutionary Biology* **7**, 214. (<https://doi.org/10.1186/1471-2148-7-214>)

- Durand EY, Patterson N, Reich D, Slatkin M. 2011. Testing for ancient admixture between closely related populations. *Molecular Biology and Evolution* **28(8)**, 2239–2252. (<https://doi.org/10.1093/molbev/msr048>)
- Earl DA, VonHoldt BM. 2012. STRUCTURE HARVESTER: a website and program for visualizing STRUCTURE output and implementing the Evanno method. *Conserv. Genet. Resour.* **4**, 359–361. (doi:10.1007/s12686-011-9548-7)
- Eaton DAR, 2014. PyRAD: assembly of de novo RADseq loci for phylogenetic analyses. *Bioinformatics* **30(13)**, 1844–1849. (doi: 10.1093/bioinformatics/btu121)
- Eaton DAR, Overcast I. 2016. ipyrad: interactive assembly and analysis of RADseq data sets. Available from: <http://ipyrad.readthedocs.io/>.
- Evanno G, Regnaut S, Goudet J. 2005. Detecting the number of clusters of individuals using the software structure: a simulation study. *Mol. Ecol.* **14**, 2611–2620. (doi:10.1111/j.1365-294X.2005.02553.x)
- Excoffier L, Lischer H. 2010. Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows. *Molecular Ecology Resources*, **10**, 564–567. (doi: 10.1111/j.1755-0998.2010.02847.x.)
- Hausmann A., Haszprunar G., Segerer AH, Speidel W, Behounek G, Hebert PDN. 2011. Now DNA-barcoded: the butterflies and larger moths of Germany. *Spixiana* **34(1)**, 47–58.
- Jakobsson M, Rosenberg NA. 2007. CLUMPP: A cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure. *Bioinformatics* **23**, 1801–1806. (doi:10.1093/bioinformatics/btm233)
- Kass R, Raftery A. 1995. Bayes factors. *Journal of the American Statistical Association* **90**, 773–795.
- Lee KM, Kivelä SM, Ivanov V, Hausmann A, Kaila L, Wahlberg N, Mutanen M. 2018. Information Dropout Patterns in Restriction Site Associated DNA Phylogenomics and a Comparison with Multilocus Sanger Data in a Species-Rich Moth Genus. *Systematic Biology* **67(6)**, 925–939. (<https://doi.org/10.1093/sysbio/syy029>)
- Leneveu J, Chichvarkhin A, Wahlberg N. 2009. Varying rates of diversification in the genus *Melitaea* (Lepidoptera: Nymphalidae) during the past 20 million years. *Biological Journal of the Linnean Society* **97**, 346–361. (doi: 10.1111/j.1095-8312.2009.01208.x)
- Leaché AD, Fujita MK, Minin VN, Bouckaert RR. 2014. Species delimitation using genome-wide SNP Data. *Systematic Biology*, **63**, 534–542. (<https://doi.org/10.1093/sysbio/syu018>)
- Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*, 1303.3997.
- Maddison WP, Maddison DR. 2017. Mesquite: a modular system for evolutionary analysis. Version 3.2.
- Malinsky M, Trucchi E, Lawson DJ, Falush D. 2018. RADpainter and fineRADstructure: Population Inference from RADseq Data. *Molecular Biology and Evolution* **35**, 1284–1290. (doi: 10.1093/molbev/msy023)
- Miller MA, Pfeiffer W, Schwartz T. 2010. Creating the CIPRES Science Gateway for inference of large phylogenetic trees. In *Proceedings of the Gateway Computing*

- Environments Workshop (GCE), New Orleans, LA, Nov. 14, 2010*. Institute of Electrical and Electronics Engineers, pp 1–8. (doi: 10.1109/GCE.2010.5676129)
- O’Connell KA, Smith EN. 2018. The effect of missing data on coalescent species delimitation and a taxonomic revision of whipsnakes (Colubridae: *Masticophis*). *Molecular Phylogenetics and Evolution* **127**, 356–366. (<https://doi.org/10.1016/j.ympev.2018.03.018>)
- Pazhenkova EA, Zakharov EV, Lukhtanov VA. 2015. DNA barcoding reveals twelve lineages with properties of phylogenetic and biological species within *Melitaea didyma* sensu lato (Lepidoptera, Nymphalidae). *ZooKeys* **538**, 35–46. (<https://doi.org/10.3897/zookeys.538.6605>)
- Pazhenkova EA, Lukhtanov VA. 2016. Chromosomal and mitochondrial diversity in *Melitaea didyma* complex (Lepidoptera, Nymphalidae): eleven deeply diverged DNA barcode groups in one non-monophyletic species? *Comparative Cytogenetics* **10(4)**, 697–717. (doi: 10.3897/CompCytogen.v10i4.11069)
- Pickrell JK, Pritchard JK. 2012. Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS Genetics*, **8**, e1002967. (<https://doi.org/10.1371/journal.pgen.1002967>)
- Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using multilocus genotype data. *Genetics* **155**, 945–959. (doi:10.1111/j.1471-8286.2007.01758.x)
- Quek SP, Davies SJ, Itino T, Pierce NE. 2004. Codiversification in an ant-plant mutualism: stem texture and the evolution of host use in *Crematogaster* (Formicidae: Myrmicinae) inhabitants of *Macaranga* (Euphorbiaceae). *Evolution* **58**, 554–570. (doi: 10.1111/j.0014-3820.2004.tb01678.x)
- Rambaut A. 2015. FigTree, v1.4.2: Tree Figure Drawing Tool. Molecular evolution, phylogenetics and epidemiology. Available from: <http://tree.bio.ed.ac.uk/software/figtree/>.
- Rambaut A, Suchard M., Drummond A. 2013. Tracer v. 1.6. Available from: <http://tree.bio.ed.ac.uk/software/tracer/>.
- Rosenberg N. 2004. DISTRUCT: a program for the graphical display of population structure. *Mol. Ecol. Notes* **4**, 137–138. (doi:10.1046/j.1471-8286.2003.00566.x)
- Stamatakis A. 2014. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313. (doi: 10.1093/bioinformatics/btu033)
- Sukumaran J, Knowles LL. 2017. Multispecies coalescent delimits structure, not species. *PNAS* **114(7)**, 1607–1612. (<https://doi.org/10.1073/pnas.1607921114>)
- Vila R, Bjorklund M. 2004. The utility of the neglected mitochondrial control region for evolutionary studies in Lepidoptera (Insecta). *Journal of Molecular Evolution* **58(3)**, 280–290. (doi: 10.1007/s00239-003-2550-2)
- Wahlberg N, Zimmermann M. 2000. Pattern of phylogenetic relationships among members of the tribe Melitaeini (Lepidoptera: Nymphalidae) inferred from mitochondrial DNA sequences. *Cladistics* **16**, 347–363. (doi: 10.1006/clad.2000.0136)
- Zhang J, Kobert K, Flouri T, Stamatakis A. 2014. PEAR: a fast and accurate Illumina Paired-End reAd mergeR. *Bioinformatics* **30**, 614–620. (doi: 10.1093/bioinformatics/btt593)

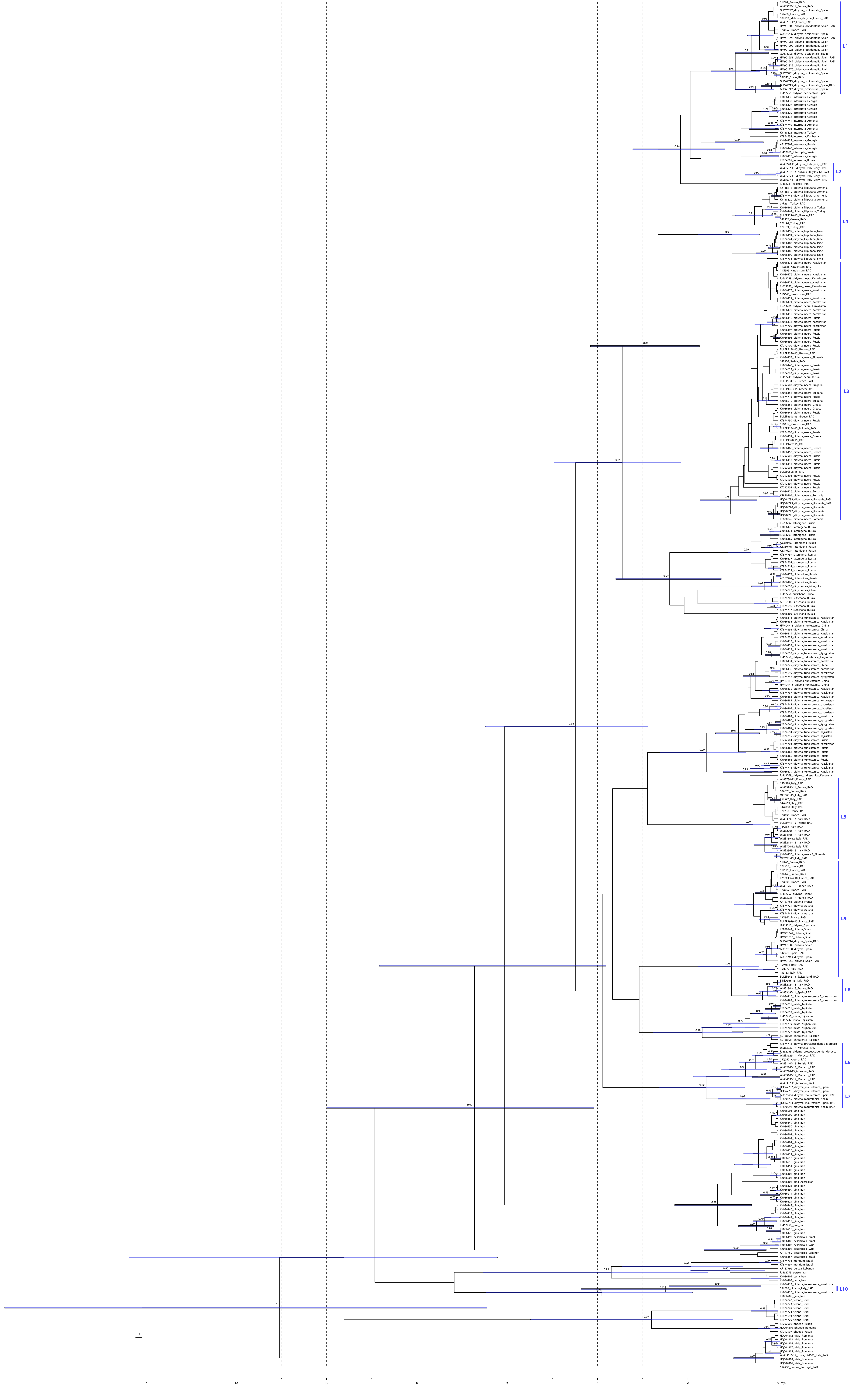


Fig. S1. Bayesian ultrametric tree for *Melitaea didyma* and related taxa based on cytochrome c oxidase subunit 1 sequences (numbers above recovered nodes, Bayesian posterior probabilities >0.7; node bars, 95% highest posterior density for age estimations; blue vertical bars, main COI lineages of *M. didyma* included in the comparison with ddRADseq data; specimens followed by "RAD" have been used for ddRADseq analyses).

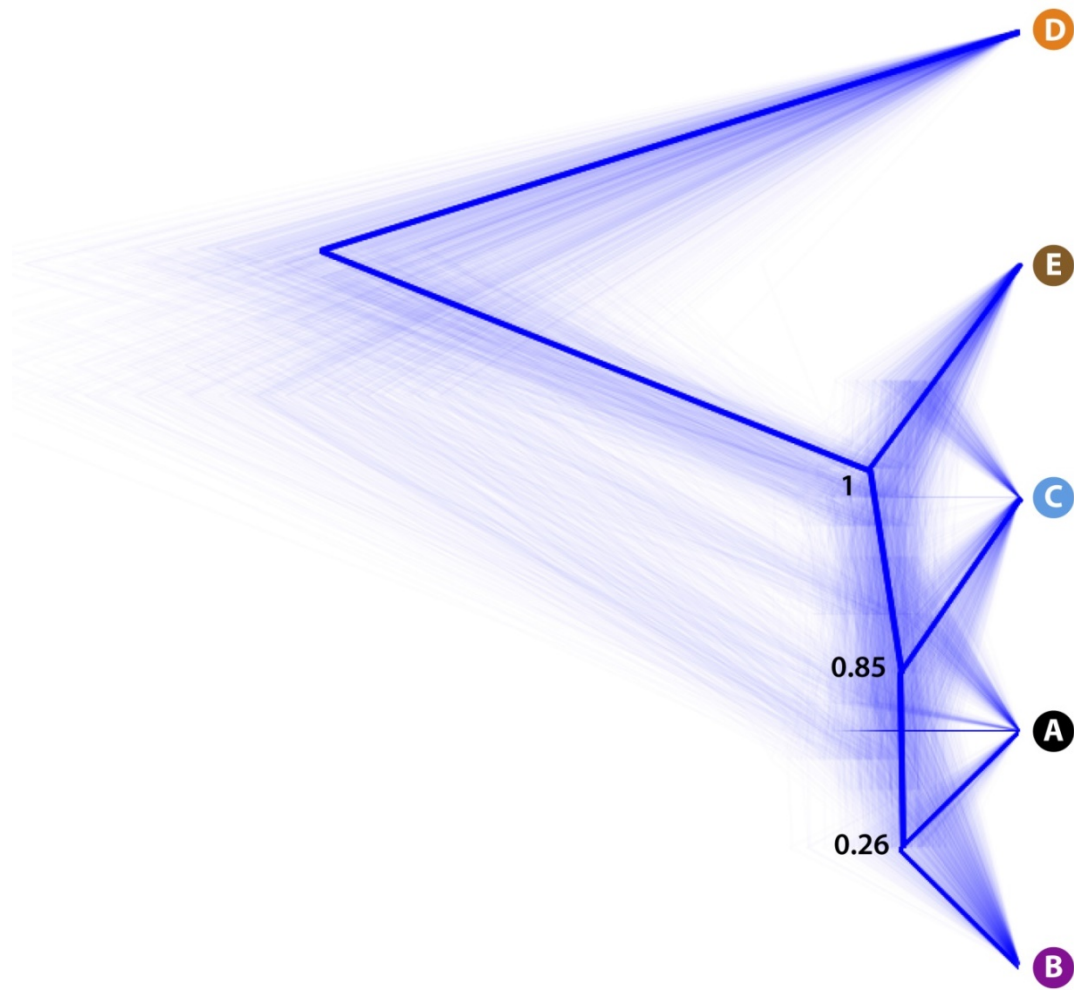


Fig. S2. DensiTree of the posterior distribution of SNAPP trees from 295 unlinked SNPs mined from ddRAD loci. The five species/clades were identified by the species delimitation analysis. The consensus is represented in thick blue lines within the species tree and the posterior probability is indicated for each node.

(a) *de novo* assembly

(b) Reference assembly

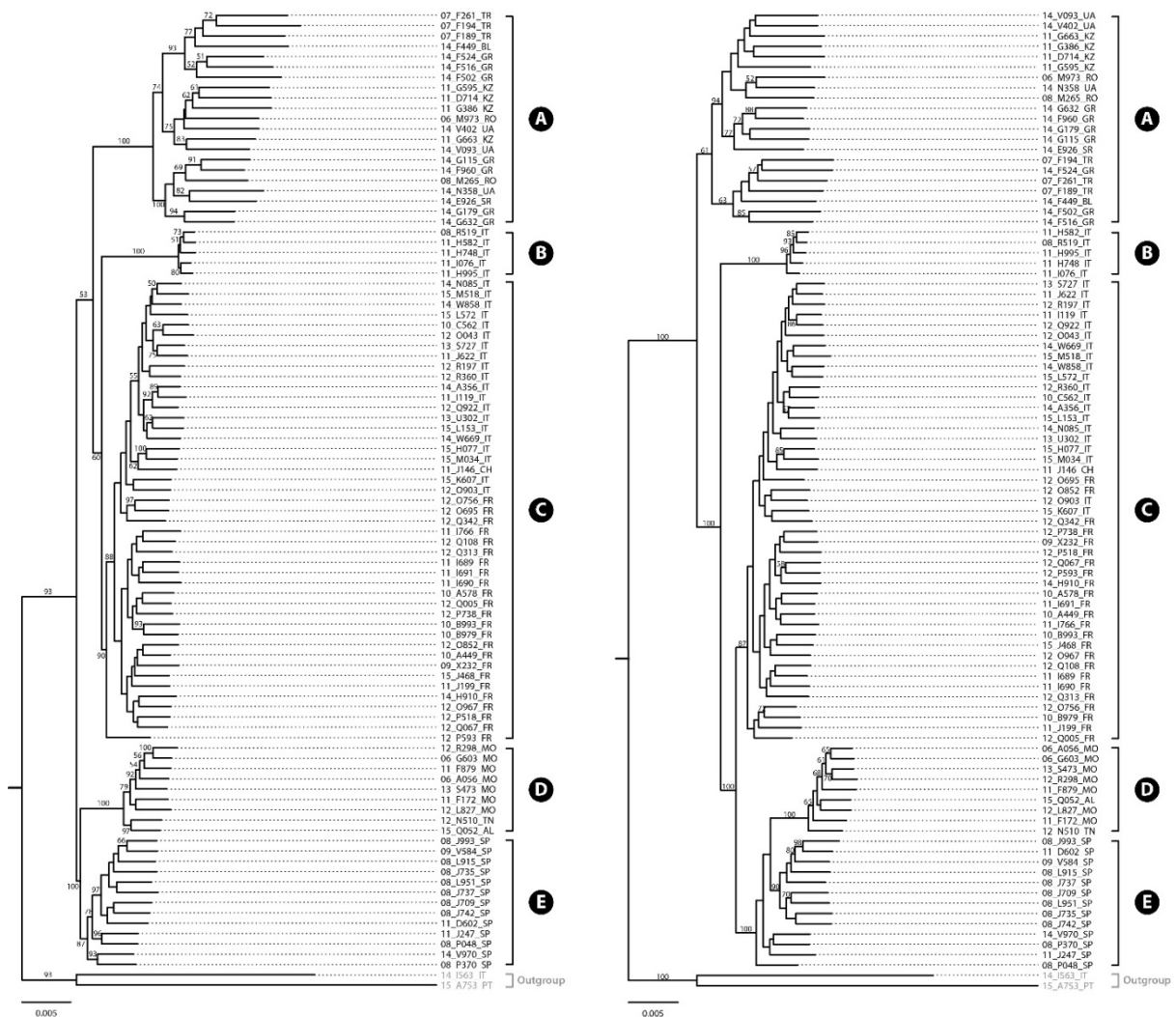
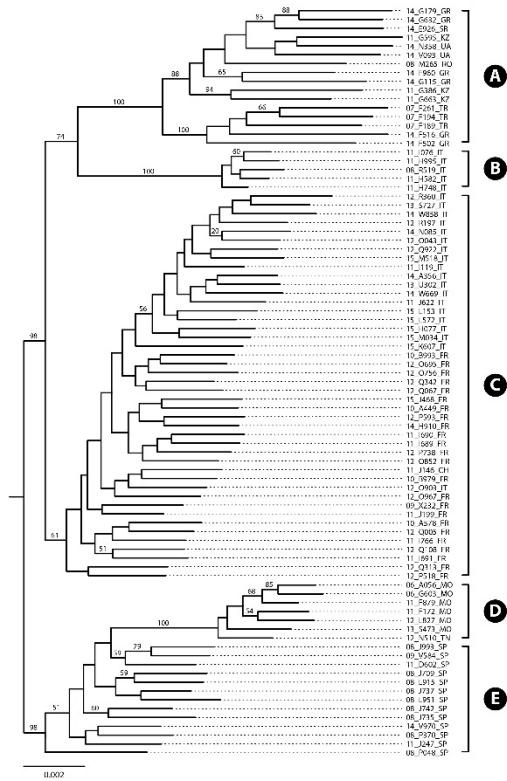
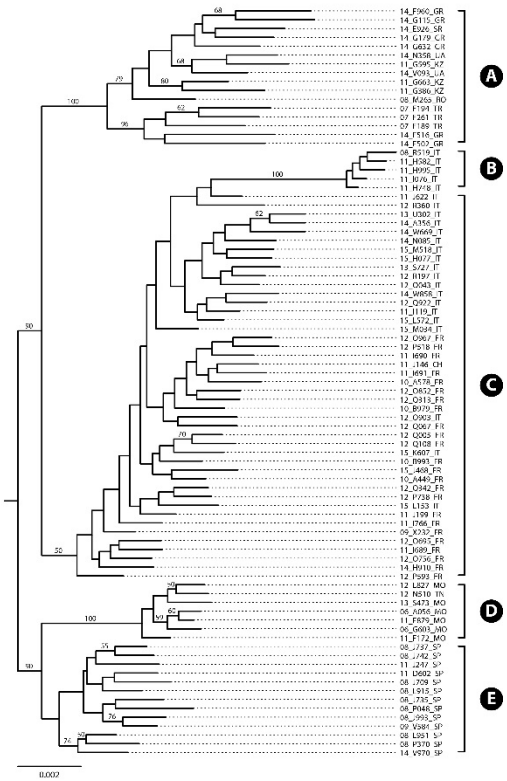


Fig. S3. ML trees inferred from (a) the *de novo* assembly data matrix and (b) the reference assembly data matrix against *Melitaea cinxia* genome (GCA_000716385) including two outgroup specimens. The phylogenetic trees were inferred with RAxML with 1,000 bootstrap replicates. Bootstrap values are indicated near branches. Only bootstrap values > 50% are shown. Branch lengths are proportional to the number of substitutions per site.

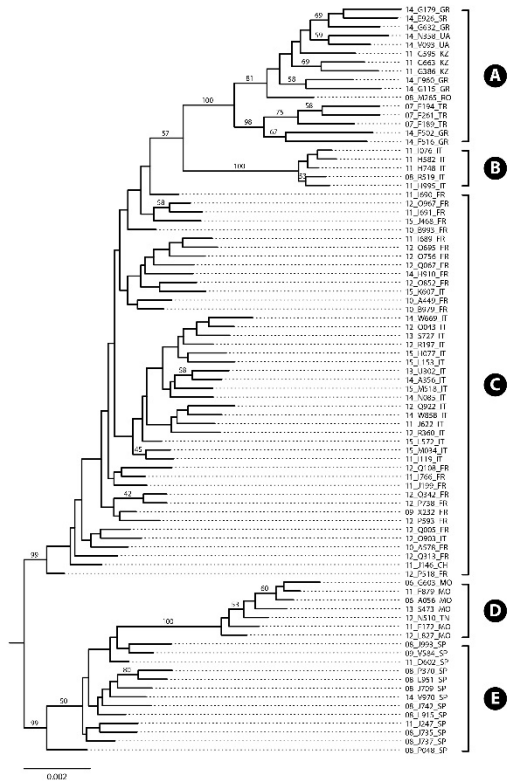
(a) missing 45%



(b) missing 20%



(c) missing 10%



(d) missing 5%

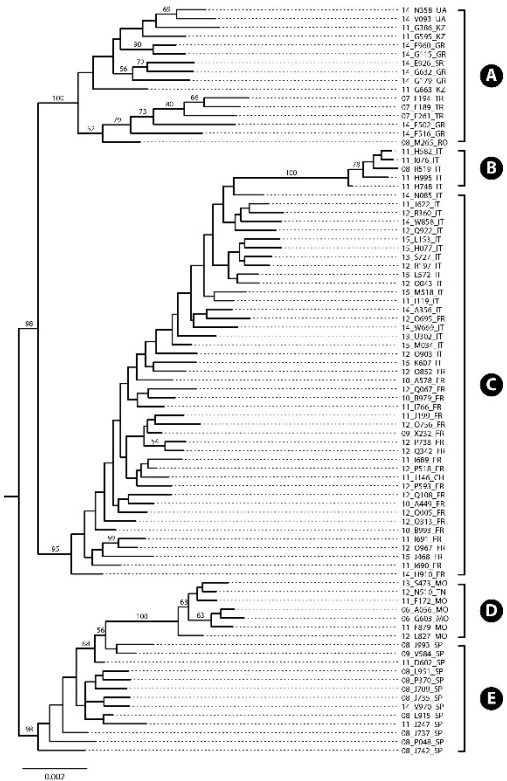


Fig. S4. ML trees based on (a) 45% of missing, (b) 20% of missing, (c) 10% of missing, and (d) 5% of missing ddRAD data.

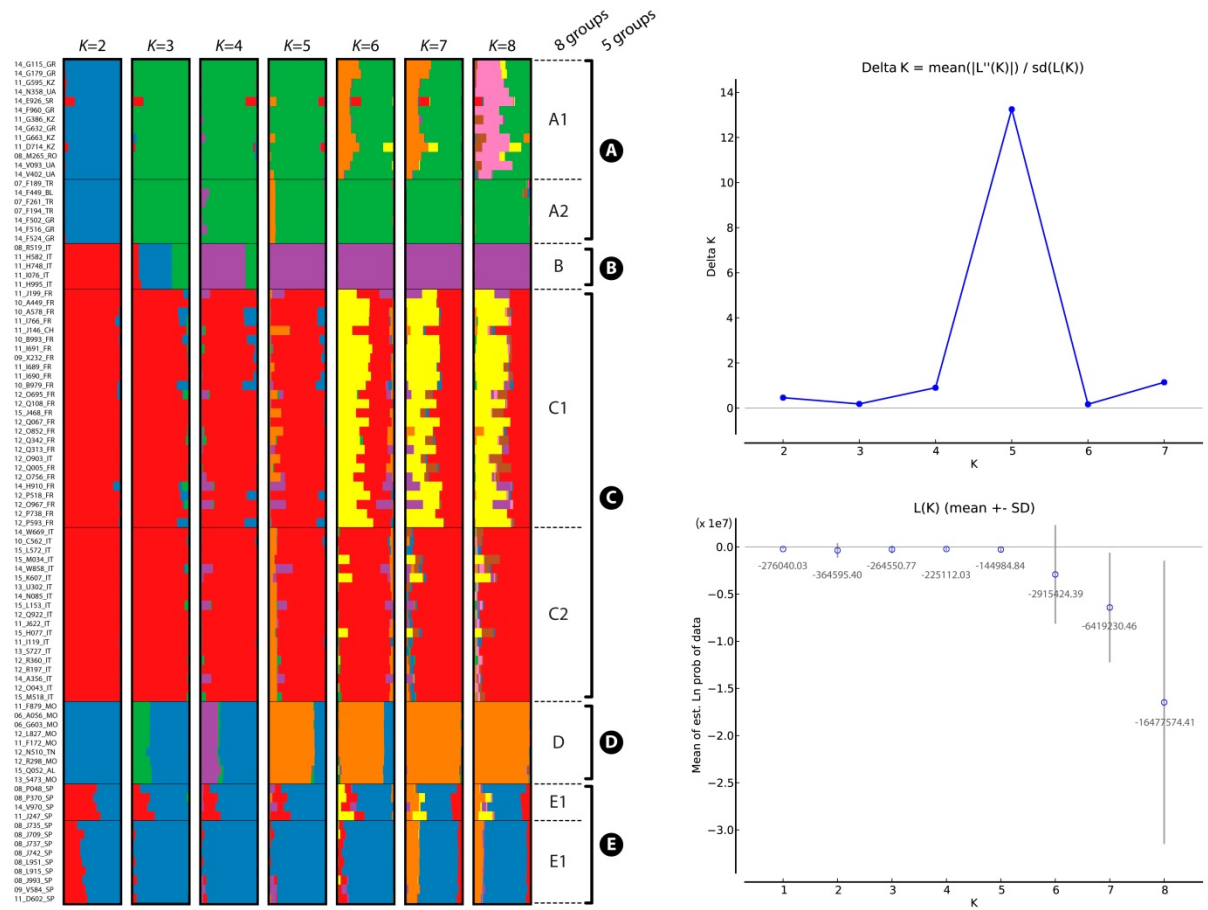


Fig. S5. Population genetic structure of the estimated delta K (ΔK) value for ten run replications of each from $K=2$ to $K=8$. The ΔK graph determined the maximum value at $K=5$. Mean probabilities $\ln P(K)$ and their standard deviation of posterior probability are shown. In the lower right panel, the numbers under the small circles represent mean estimated likelihood values for each K .

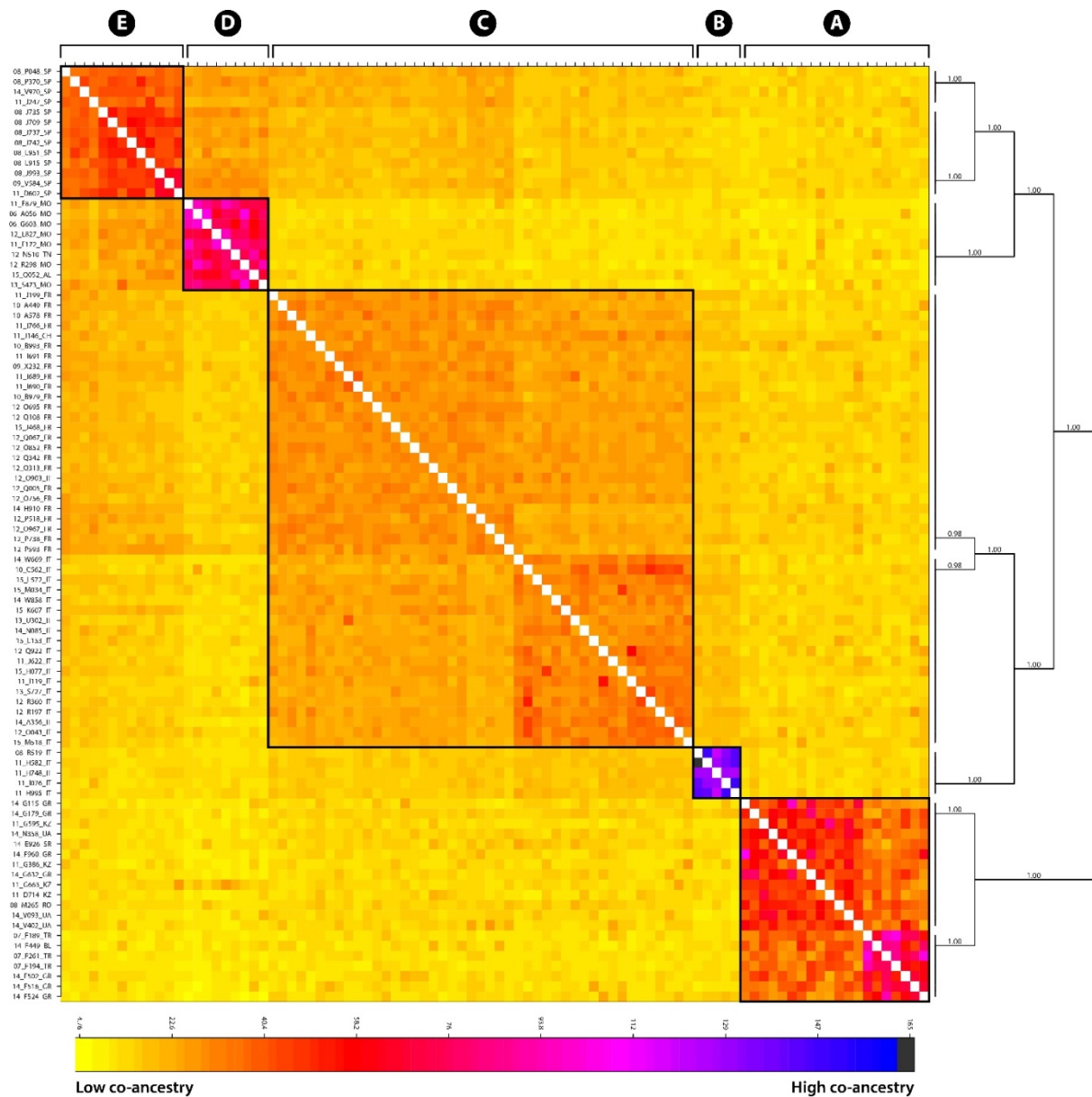


Fig. S6. Clustered FineRADstructure co-ancestry matrix for *M. didyma*. The highest levels of co-ancestry are evident among individuals from Sicily population (clade B), indicated by black, blue and purple colours. The lowest levels of co-ancestry sharing are indicated by yellow coloration. Individuals clustering into species/populations are indicated by clustering in the accompanying tree.

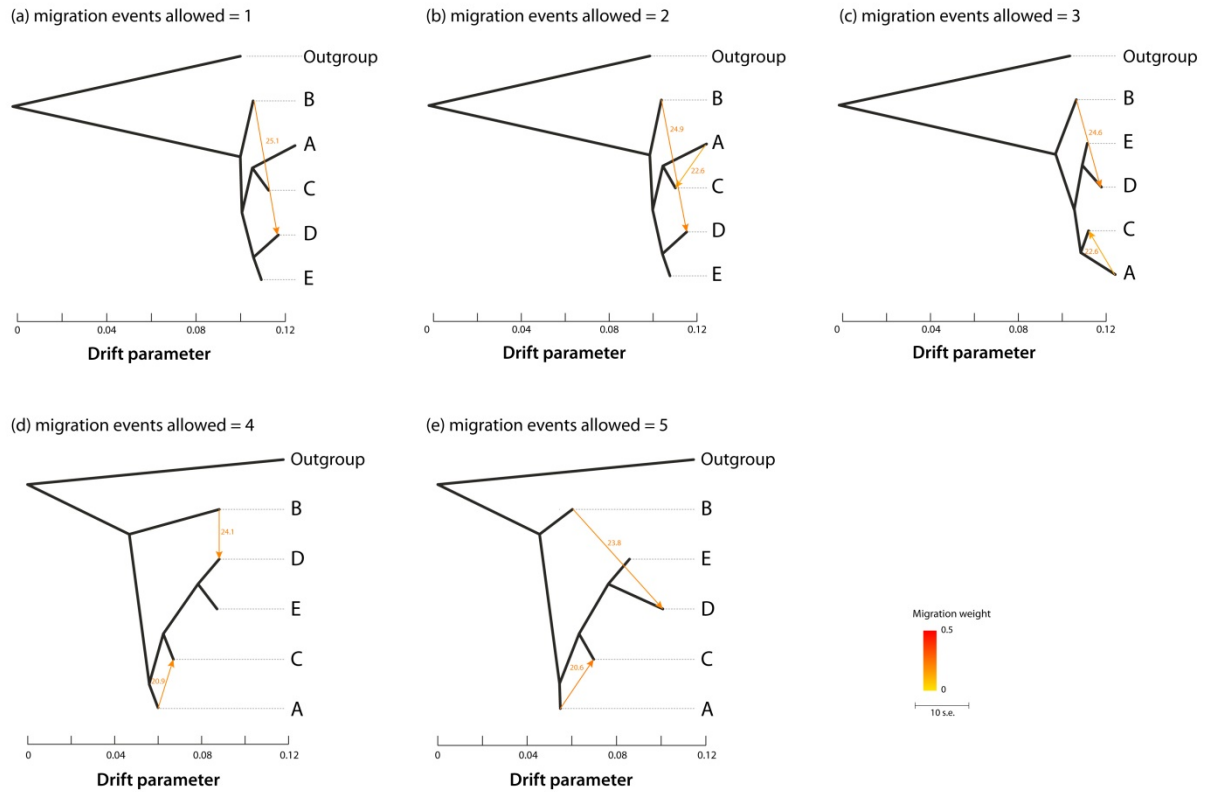


Fig. S7. Relatedness among *M. didyma* lineages recovered by *TreeMix*. The dataset included 25 specimens of *M. didyma* (five for each clade) and two specimens (*M. trivia* and *M. deione*) used as outgroup. The arrows (migration events) are coloured according to their weight. Only migration events that were statistically significant are shown. The migration weight represents the fraction of ancestry derived from the migration edge. The scale bar on the right shows the average standard error (s.e.) of the entries in the sample covariance matrix. Notable levels of gene flow are from B to D (Jackknife $P=0.00012$) and from A to C ($P=0.00030$).

Sample ID	BOLD process ID (COI)	COI-5P Accession	ddRADseq data	ddRADseq lineage	COI lineage	Wolbachia infection	BOLD process ID (Wolbachia)	WSP Accession	WSP allele	ftsZ Accession	ftsZ allele	Taxon	Collection Date	Country	Exact Site	Lat	Lon	Elevation (m)
CCDB-17951 A12	BPAL2292-14	KT874699										<i>M. mixta</i>	01-Jul-1992	Tajikistan	Peter I Range, Ganischou			
CCDB-17966 B04	BPAL2486-14	KT874711										<i>M. mixta</i>	12-Jul-2004	Tajikistan	Peter I Range			
CCDB-17966 D03	BPAL2509-14	KT874722										<i>M. mixta</i>	05-Jul-1994	Tajikistan	Pendzhikent Farob			
CCDB-17966 B06	BPAL2488-14	KT874731										<i>M. mixta</i>	19-Jun-2004	Tajikistan	Peter I Range			
CCDB-17949 E02	BPAL2235-13	KT874697										<i>M. montium</i>	22-Jun-2013	Israel	Mt. Hermon	33.3127	35.795	2056
CCDB-17949 E03	BPAL2236-13	KT874736										<i>M. montium</i>	22-Jun-2013	Israel	Mt. Hermon	33.3127	35.795	2056
		FJ462273										<i>M. perse</i>		Iran	Ardabil			
		AF187796										<i>M. perse</i>		Lebanon	Les Cedres, Mohafazat Beharre			
		HQ004810										<i>M. phoebe</i>		Romania				
52025		KT792906										<i>M. phoebe</i>	5 July 2014	Russia	Samara obl., Kinel district	53°0'13.8	50°9'24.6	
160017		KT792907										<i>M. phoebe</i>	27 June 2014	Russia	Samara obl., Bachilova polyana	53°4'28.6	49°6'61.8	
		FJ462281										<i>M. saxatilis</i>		Iran	Tehran			
FJ462254	GBLN1854-09	FJ462254										<i>M. sutschana</i>		China	Hebei			
AF187805	GBLN0099-06	AF187805										<i>M. sutschana</i>		Russia	Chita region, Kyra			
CCDB-17966 G01	BPAL2543-14	KT874696										<i>M. sutschana</i>	18-Jun-1995	Russia	Chita reg., Talacha			
CCDB-17966 G03	BPAL2545-14	KT874701										<i>M. sutschana</i>	18-Jun-1995	Russia	Chita reg., Talacha			
CCDB-17966 G02	BPAL2544-14	KT874717										<i>M. sutschana</i>	18-Jun-1995	Russia	Chita reg., Talacha			
	BPAL2309-14	KY086105										<i>M. sutschana</i>		Russia	Primorskij Kray			
CCDB-17949 C05	BPAL2214-13	KT874693										<i>M. telona</i>	24-Apr-2013	Israel	Odem Forest	33.174	35.7522	960
CCDB-17968 E09	BPAL2717-14	KT874723										<i>M. telona</i>	03-Jul-2014	Israel	Hermon			
CCDB-17949 G08	BPAL2265-13	KT874724										<i>M. telona</i>	23-Jun-2013	Israel	Avivim	33.0825	35.4594	616
CCDB-17949 C06	BPAL2215-13	KT874729										<i>M. telona</i>	24-Apr-2013	Israel	Odem Forest	33.174	35.7522	960
CCDB-17949 G09	BPAL2266-13	KT874747										<i>M. telona</i>	23-Jun-2013	Israel	Avivim	33.0825	35.4594	616
CCDB-17949 G10	BPAL2267-13	KT874749										<i>M. telona</i>	23-Jun-2013	Israel	Avivim	33.0825	35.4594	616
		HQ004812										<i>M. trivla</i>		Romania				
		HQ004813										<i>M. trivla</i>		Romania				
		HQ004814										<i>M. trivla</i>		Romania				
		HQ004815										<i>M. trivla</i>		Romania				
		HQ004816										<i>M. trivla</i>		Romania				
		HQ004817										<i>M. trivla</i>		Romania				
		HQ004818										<i>M. trivla</i>		Romania				

Table S2. Specimens of *Melitaea didyma* analysed in this study and a summary of the ddRAD data in *de novo* and reference assemblies. Individuals selected for BFD* analyses assuming 5 taxa are marked with *, while individuals selected for analyses assuming 8 taxa are marked with †. Specimens selected for *TreeMix* and D-statistics as representatives for each clade are marked with ^d in Sample ID.

Clade	Sample ID	Population	Total reads (x10 ⁶)	<i>de novo</i> assembly				Reference assembly				
				Clusters at 90%	Mean depth	Retained loci	Recovered loci	Mapped reads	Clusters total	Clusters depth	Reads consensus	Recovered loci in assembly
Clade A	06_M973_RO	Romania	2.28	15274	33.2	5337	524	59642	5845	6.24	1651	350
	07_F189_TR * †	Turkey	4.03	24358	55.4	9850	1700	209962	12829	115.8	4437	882
	07_F194_TR * †	Turkey	2.48	34519	33.5	11310	1674	247576	16262	53.7	4911	961
	07_F261_TR * † ^d	Turkey	6.02	50392	88.8	14263	1995	244540	27353	26.2	5800	1093
	08_M265_RO * †	Romania	2.31	25899	53.8	9138	1664	157324	14989	32.5	4638	886
	11_D714_KZ	Kazakhstan	1.18	17098	16.1	6667	995	68142	9082	22.3	2914	516
	11_G386_KZ * † ^d	Kazakhstan	3.48	47146	18.4	15385	2118	382712	24238	58.8	6136	1138
	11_G595_KZ * † ^d	Kazakhstan	4.06	47318	17.8	15437	2074	398166	23400	54.3	5927	1067
	11_G663_KZ * † ^d	Kazakhstan	3.3	56864	25.8	17012	2318	369878	25657	45.5	5648	1065
	14_E926_SR	Serbia	5.74	25307	18.2	7796	1287	181210	13618	32.5	3276	614
	14_F449_BL	Bulgaria	1.23	16678	24.9	4827	613	123670	9597	53.1	1618	230
	14_F502_GR * †	Greece	1.95	27476	26.7	9099	1361	182892	15929	45	4072	676
	14_F516_GR * †	Greece	1.48	20478	21.3	7460	1213	172892	13050	50	3709	668
	14_F524_GR	Greece	2.5	16681	29.8	5929	729	124156	9460	62.3	1838	273
	14_F960_GR	Greece	1.88	36357	9.4	8913	1390	195102	20758	33.9	3576	668
	14_G115_GR	Greece	2.21	28807	13.8	8359	1320	194692	18825	30.6	3585	592
14_G179_GR	Greece	3.27	31195	34.1	9274	1715	236990	19871	22.5	4217	829	

Clade	Sample ID	Population	Total reads (x10 ⁶)	<i>de novo</i> assembly				Reference assembly				
				Clusters at 90%	Mean depth	Retained loci	Recovered loci	Mapped reads	Clusters total	Clusters depth	Reads consensus	Recovered loci in assembly
	14_G632_GR ^{***d}	Greece	3.08	58007	14.7	14931	2546	386768	29190	44.7	6029	1152
	14_N358_UA	Ukraine	1.12	31258	11.7	12066	1983	175170	15627	38.2	5076	1039
	14_V093_UA	Ukraine	3.3	32518	40.1	11458	1870	257138	17099	79.4	4724	949
	14_V402_UA	Ukraine	0.12	8381	4.9	4113	814	30494	5878	8.4	1405	274
Clade B	08_R519_IT ^{***d}	Sicily, Italy	1.36	21986	22.4	8988	2121	147270	10043	62.8	3750	907
	11_H582_IT ^{***d}	Sicily, Italy	2.32	23857	60.4	9413	2239	165896	14130	31.6	4453	1149
	11_H748_IT ^{***d}	Sicily, Italy	1.39	19218	32.2	7754	1843	124888	10450	42.1	3408	808
	11_H995_IT ^{***d}	Sicily, Italy	1.80	40248	17.6	12444	2599	253268	22922	36.0	5456	1310
	11_I076_IT ^{***d}	Sicily, Italy	1.19	26924	12.4	8569	1833	165730	14848	32.2	3940	937
Clade C	09_X232_FR ^{***}	France	1.93	40312	14.7	11869	2476	279780	20151	37.8	5330	1298
	10_A449_FR ^{***}	France	0.91	22729	16.8	9457	2197	134320	13104	28.6	4777	1225
	10_A578_FR	France	1.29	29083	16.1	9186	1762	191166	19295	25	4398	1003
	10_B979_FR ^{***}	France	2.71	40840	17.7	12316	2402	310712	21028	49.7	5629	1429
	10_B993_FR ^{***}	France	2.62	36188	33.3	11311	2372	235530	18688	43.7	5078	1244
	10_C562_IT	Italy	1.77	13311	75.9	4480	725	87592	8896	42.7	1667	273
	11_I119_IT ^{***}	Italy	0.95	24702	10.5	8622	1798	128742	15079	22.8	4260	1034
	11_I689_FR ^{***d}	France	2.65	39531	24	12549	2834	295120	24756	36.3	6170	1561
	11_I690_FR	France	2.71	35575	26.4	11326	2524	289072	20227	45.2	5449	1370
	11_I691_FR	France	3.46	45263	37.9	12475	2422	272110	23106	46.5	5419	1372

Clade	Sample ID	Population	Total reads (x10 ⁶)	<i>de novo</i> assembly				Reference assembly				
				Clusters at 90%	Mean depth	Retained loci	Recovered loci	Mapped reads	Clusters total	Clusters depth	Reads consensus	Recovered loci in assembly
	11_I766_FR	France	3.11	49303	30	12876	2331	383782	23303	33	5081	1201
	11_J146_CH ^d	Switzerland	3.6	49038	42.1	13873	2804	309550	23328	42.3	5552	1394
	11_J199_FR	France	1.74	28745	30.1	10340	2252	188032	15267	35.5	4729	1111
	11_J622_IT * ⁺	Italy	3.62	38588	38.7	11996	2410	283910	19980	88.2	4652	1129
	12_O043_IT * ⁺	Italy	5.95	130628	16.4	26196	2717	706154	64475	34.3	5970	1324
	12_O695_FR	France	1.87	34320	20.6	11375	2683	230602	21511	32.3	5665	1501
	12_O756_FR	France	3.8	41203	55.8	11642	2331	220854	23501	26.7	5033	1266
	12_O852_FR ^d	France	3.06	49632	20.6	13744	2914	373276	26048	40.9	6389	1618
	12_O903_IT	Italy	1.3	26338	22.4	9553	2103	154660	15999	26.3	4552	1083
	12_O967_FR	France	2.31	33119	32.7	10920	2155	187056	17017	32.9	5009	1271
	12_P518_FR	France	2.63	38884	18.9	11638	2281	314420	17962	79.3	4889	1163
	12_P593_FR	France	1.64	25433	30	9273	1973	147616	15834	31.4	4447	1071
	12_P738_FR	France	1.52	22080	18.9	7857	1650	171090	15531	35.1	4083	939
	12_Q005_FR	France	1.49	21618	31.3	7789	1806	173616	16320	45.8	4201	1018
	12_Q067_FR	France	1.2	22320	17.8	8259	1788	168776	16449	31.1	4214	930
	12_Q108_FR	France	1.82	28594	22.4	9660	2284	216026	21778	30	5176	1299
	12_Q313_FR	France	1	23929	13.2	8850	1970	154020	16189	27.1	4599	1139
	12_Q342_FR	France	2.71	25984	44.5	9477	2154	230218	15571	34	4725	1225
	12_Q922_IT	Italy	3.36	22750	23.6	7917	1584	184080	14917	47.9	3816	821

Clade	Sample ID	Population	Total reads (x10 ⁶)	<i>de novo</i> assembly				Reference assembly				
				Clusters at 90%	Mean depth	Retained loci	Recovered loci	Mapped reads	Clusters total	Clusters depth	Reads consensus	Recovered loci in assembly
	12_R197_IT ^{*+}	Italy	2.23	33055	25.6	12192	2722	247562	19296	44.3	5692	1454
	12_R360_IT	Italy	3.62	28497	39.6	8777	1536	168472	12917	40.7	3416	750
	13_S727_IT ^{*+}	Italy	2.19	26052	35.2	9885	2180	175724	14622	51.2	4248	1010
	13_U302_IT	Italy	5.01	42209	48.4	11070	1699	195064	19359	33.9	3465	719
	14_A356_IT	Italy	2.33	28158	43.3	9196	1814	182500	16071	44.3	4031	909
	14_H910_FR ^d	France	2.32	44470	19.8	13097	2940	270434	27814	25.5	5913	1449
	14_N085_IT	Italy	1.29	20529	16.6	7421	1664	145756	12681	41.2	3746	883
	14_W669_IT	Italy	2.61	30856	35.9	10658	2121	217732	16638	49.9	4496	1081
	14_W858_IT	Italy	2.84	22778	62.2	8825	1848	175536	11864	61.4	3479	827
	15_H077_IT	Italy	2.24	31959	26	10311	2007	210362	17707	39.1	4603	1081
	15_J468_FR	France	1.35	26273	12.4	9251	2118	170442	18686	28.8	4817	1133
	15_K607_IT	Italy	2.03	28905	23.9	10192	2245	216260	17364	44.3	5020	1252
	15_L153_IT	Italy	2.05	26591	31.4	9543	1969	183886	14037	48.1	4442	1060
	15_L572_IT ^d	Italy	3.53	48735	41.9	13524	2818	328446	24536	49.1	5787	1457
	15_M034_IT	Italy	3.24	44554	28.7	13189	2786	365898	26581	42.3	6022	1479
	15_M518_IT	Italy	2.64	31180	37.1	10033	2073	233116	17417	66	4756	1165
Clade D	06_A056_MO	Morocco	1.19	19869	21.4	6923	1472	125092	11244	40.8	3272	739
	06_G603_MO ^{*+d}	Morocco	2.83	29637	33.2	7816	1502	264250	15270	82.6	3334	700
	11_F172_MO ^{*+d}	Morocco	3.09	36496	19.7	12968	2281	360028	24329	64.1	5824	1361

Clade	Sample ID	Population	Total reads (x10 ⁶)	<i>de novo</i> assembly				Reference assembly				
				Clusters at 90%	Mean depth	Retained loci	Recovered loci	Mapped reads	Clusters total	Clusters depth	Reads consensus	Recovered loci in assembly
	11_F879_MO ^{***d}	Morocco	3.14	39557	21.1	13140	2272	388136	24631	55.9	5805	1331
	12_L827_MO ^{***d}	Morocco	2.45	39441	16.6	13409	2250	351008	20860	51.2	5547	1190
	12_N510_TN ^{***d}	Tunisia	1.90	38983	18.9	12627	2395	248060	22387	37.5	5724	1373
	12_R298_MO	Morocco	1.76	15547	27.9	6232	984	145658	10280	55.9	2822	523
	13_S473_MO	Morocco	1.90	21259	21.7	8169	1497	166364	11948	69.8	3680	849
	15_Q052_AL	Algeria	0.26	10845	9.5	5074	1103	56032	8201	12.7	2614	555
Clade E	08_J709_SP ^{***d}	Spain	2.13	57522	12.3	13094	2529	368606	24180	40.0	5683	1366
	08_J735_SP ^{***d}	Spain	4.40	93052	12.8	19474	2328	584200	32936	38.4	5090	1100
	08_J737_SP ^{**}	Spain	2.67	40706	17.4	12095	2235	359580	22230	54.6	5355	1287
	08_J742_SP ^{**}	Spain	2.21	45260	13.8	13269	2241	319584	22549	41.3	5313	1293
	08_J993_SP ^{***d}	Spain	2.59	44521	19.5	12803	2371	372706	20395	57.7	5112	1153
	08_L915_SP	Spain	2.36	32009	20.1	10440	1970	288670	17769	50.5	4825	1158
	08_L951_SP	Spain	2.05	27384	19.2	9816	1874	257652	17286	51.4	5034	1209
	08_P048_SP ⁺	Spain	1.81	24840	34.1	9532	2113	178930	15001	40.5	5037	1247
	08_P370_SP ^{+d}	Spain	2.51	45280	23.6	12985	2286	304924	22921	35.8	5266	1262
	09_V584_SP	Spain	2.64	42804	17.9	12371	2167	300940	22653	44.2	5492	1364
	11_D602_SP	Spain	2.11	32095	20.1	10542	1861	241174	21479	37.9	4827	1117
	11_J247_SP ⁺	Spain	3.07	41516	45.6	11579	2259	235334	18293	27.5	4688	1139
	14_V970_SP ^{+d}	Spain	2.32	41624	21.9	13118	2590	292324	19658	53.7	5506	1367

Clade	Sample ID	Population	Total reads (x10 ⁶)	<i>de novo</i> assembly				Reference assembly				
				Clusters at 90%	Mean depth	Retained loci	Recovered loci	Mapped reads	Clusters total	Clusters depth	Reads consensus	Recovered loci in assembly
Outgroup	14_I563_IT * + ^d	Italy	0.99	45229	63.2	12810	844	30955	4981	73.82	1940	131
	15_A753_PT * + ^d	Portugal	2.34	41414	90.4	10153	4057	113910	15916	13.74	5315	200
AVERAGE			2.42	34168	28.3	10634	2003	238205	18636	43.5	4623	1059

Table S3. Summary of ddRAD and mtDNA data sets.

Data source	ddRAD			mtDNA
Data matrix	ddRAD_mt_m4	ddRAD_dn_out	ddRAD_ref	mtDNA_COI
Outgroup	NA	14_I563_IT, 15_A753_PT	14_I563_IT, 15_A753_PT	14_I563_IT, 15_A753_PT
Assembly method	de novo – reference ^a	de novo	Reference ^b	NA
Number of taxa	92 ^c	95	95	93
Number of loci	22,353	22,342	14,525	1
Alignment length (bp)	3,489,654	3,487,180	2,548,000	658
SNPs	143,201	144,004	116,942	141
Informative sites	46,371	46,096	41,362	88
Missing data (%)	90.9	91.3	89.7	0

^a *Melitaea cinxia* mitochondrion genome (CM002851) was used as reference genome.

^b *Melitaea cinxia* whole genome sequences (GCA_000716385) were used as reference.

^c Specimen “RV-06-M973” was removed due to the low number of recovered loci in the final data matrix.

Table S4. Primers and PCR protocols used for the amplification of COI, wsp and ftsZ.

COI	
Reagents	µl per reaction
Buffer 5X	5
MgCl ₂ (25 mM)	2
dNTPs (10mM)	0.5
LepF1 (10µM)	0.5
LepR1 (10µM)	0.5
H ₂ O	14.4
Taq polymerase (Promega)	0.1
DNA extraction	2

Total volume = 25

PCR program			
Step	Temperature	Duration	Cycles
1	92°C	60 s	
2	92°C	15 s	
3	48°C	45 s	Steps 2-4, 5X
4	62°C	150 s	
5	92°C	15 s	
6	52°C	45 s	Steps 5-7, 35X
7	62°C	150 s	
8	62°C	420 s	

wsp	
Reagents	µl per reaction
Buffer 10X	1
MgCl ₂ (25 mM)	1
dNTPs (10mM)	0.2
wsp81_F (10µM)	0.5
wsp691_R (10µM)	0.5
H ₂ O	5.7
AmpliTaQ Gold (Applied Biosystems)	0.1
DNA extraction	1

Total volume = 10

PCR program			
Step	Temperature	Duration	Cycles
1	95°C	300 s	
2	95°C	30 s	
3	55°C	30 s	Steps 2-4, 38X
4	72°C	120 s	
5	72°C	120 s	

ftsZ	
Reagents	µl per reaction
Buffer 10X	1
MgCl ₂ (25 mM)	1
dNTPs (10mM)	0.2
ftsZ_F (10µM)	0.5
ftsZ_R (10µM)	0.5
H ₂ O	5.7
AmpliTaQ Gold (Applied Biosystems)	0.1
DNA extraction	1

Total volume = 10

PCR program			
Step	Temperature	Duration	Cycles
1	95°C	300 s	
2	95°C	30 s	
3	54°C	30 s	Steps 2-4, 38X
4	72°C	120 s	
5	72°C	120 s	

Primer name	Primer sequence (5' - 3')	Direction	Marker	Reference
LepF1	ATTCAACCAATCATAAAGATATTGG	Forward	COI	Hebert et al. 2004
LepR1	TAAACTTCTGGATGTCCAAAAAATCA	Reverse	COI	Hebert et al. 2004
wsp81_F	TGGTCCAATAAGTGATGAAGAAAC	Forward	wsp	Baldo et al. 2006
wsp691_R	AAAAATTAAACGCTACTCCA	Reverse	wsp	Baldo et al. 2006
ftsZ_F	ATYATGGARCATATAAARGATAG	Forward	ftsZ	Baldo et al. 2006
ftsZ_R	TCRAGYAATGGATRRGATAT	Reverse	ftsZ	Baldo et al. 2006

References

Baldo L, Dunning Hotopp JC, Jolley KA, Bordenstein SR, Biber SA, Choudhury RR, Hayashi C, Maiden MC, Tettelin H, Werren JH. 2006. Multilocus sequence typing system for the endosymbiont *Wolbachia pipientis*. *Applied and Environmental Microbiology* **72**, 7098–7110. (doi:10.1128/AEM.00731-06)

Hebert PDN, Penton EH, Burns JM, Janzen DH, Hallwachs W. 2004. Ten species in one: DNA barcoding reveals cryptic species in the neotropical skipper butterfly *Astraptes fulgerator*. *Proc. Natl. Acad. Sci. USA* **101**, 14812–14817. (doi:10.1073/pnas.0406166101)

Table S5. Minimum p-distance matrix between ten COI lineages of *M. didyma*. Minimum and maximum values are in blue and red, respectively.

	Minimum p-distance (%)									
	L1	L2	L3	L4	L5	L6	L7	L8	L9	L10
L1	n/a									
L2	2.43	n/a								
L3	2.43	2.43	n/a							
L4	2.89	3.04	1.82	n/a						
L5	2.28	3.04	2.43	2.43	n/a					
L6	3.34	3.34	3.34	3.65	2.89	n/a				
L7	3.34	3.34	3.04	3.65	2.74	1.52	n/a			
L8	4.26	5.02	4.41	3.8	3.8	3.95	3.8	n/a		
L9	3.65	4.26	3.50	3.19	2.58	2.74	2.13	1.22	n/a	
L10	5.93	7.45	6.53	6.23	5.62	6.08	6.69	7.14	6.38	n/a

Table S6. Pairwise F_{ST} values among all clades of *M. didyma* using unlinked SNP data. An asterisk denotes significant statistical support from 1,000 permutations ($p < 0.05$).

	A	B	C	D	E	Outgroup
A	–					
B	0.79959 *	–				
C	0.63294 *	0.95894 *	–			
D	0.68593 *	1.00000 *	0.88916 *	–		
E	0.67690 *	0.87873 *	0.65320	0.65689 *	–	
Outgroup	0.61531 *	1.00000 *	0.94708 *	1.00000 *	0.87171 *	–

Table S7. Species delimitation scenarios for *Melitaea didyma*. a) Results assuming five species. b) Results assuming eight species. The table shows the different species delimitation models for the group evaluated with the BFD* method and the associated results. Each row indicates a different species delimitation model and two samples ('14_I563_IT' and '15_A753_PT') were included in all tests as outgroup taxa (O). The best delimitation scenario is shown in bold.

a)

Scenarios	Description	n species (w/o outgroup)	MLE	Bayes factor	Rank
Current taxonomy	(A,B,C,D,E),(O)	1	-3783.7	NA	9
Split 2.1	(A,B,C),(D,E),(O)	2	-3718.1	-131.2	7
Split 2.2	(A),(B,C,D,E),(O)	2	-3750.6	-66.2	8
Split 2.3	(A,C,D,E),(B),(O)	2	-3686.4	-194.5	6
Split 3.1	(A),(B,C),(D,E),(O)	3	-3661.3	-244.9	3
Split 3.2	(A,B),(C),(D,E),(O)	3	-3678.8	-209.7	5
Split 3.3	(A,B,C),(D),(E),(O)	3	-3667.4	-232.6	4
Split 4	(A),(B,C),(D),(E),(O)	4	-3612.1	-343.3	2
Split 5	(A),(B),(C),(D),(E),(O)	5	-3551.0	-465.3	1

b)

Scenarios	Description	n species (w/o outgroup)	MLE	Bayes factor	Rank
Current taxonomy	(A,B,C,D,E),(O)	1	-224334.8	NA	10
Split 2.1	(A,B,C),(D,E),(O)	2	-213984.6	-20700.5	7
Split 2.2	(A),(B,C,D,E),(O)	2	-215945.4	-16778.8	9
Split 2.3	(A,C,D,E),(B),(O)	2	-215330.0	-18009.7	8
Split 3.1	(A),(B,C),(D,E),(O)	3	-207846.9	-32975.8	4
Split 3.2	(A,B),(C),(D,E),(O)	3	-208976.7	-30716.1	5
Split 3.3	(A,B,C),(D),(E),(O)	3	-212030.2	-24609.2	6
Split 4	(A),(B,C),(D),(E),(O)	4	-205842.8	-36984.1	3
Split 5	(A),(B),(C),(D),(E),(O)	5	-199276.4	-50116.7	2
Split 8	(A₁),(A₂),(B),(C₁),(C₂),(D),(E₁),(E₂),(O)	8	-198600.0	-51469.5	1

Table S8. Tests of admixture using five clades and 25 specimens of *M. didyma* (five for each clade), as well as two specimens (*M. trivialis* and *M. deione*) used as outgroup. P1, P2, and P3 indicate species/population used in a given topology position when testing for admixture using Patterson’s D-statistics. Each test was repeated over all possible four-sample replicates (n), with a range of Z-scores reported, and the number of significant replicates shown (nSig). Outgroup, not shown in the table, consists of two individuals (‘14_I563_IT’ and ‘15_A753_PT’). n loci used is the number of loci analyzed in each test.

Test	P1	P2	P3	Range Z	nSig/n	nSig/n (%)	n loci used
1	A	A	B	(0.2 – 5.4)	2/49	4.1	104
2	A	A	C	(0.0 – 9.5)	6/49	12.2	110
3	A	A	D	(0.0 – 11.8)	12/49	24.5	106
4	A	A	E	(0.1 – 3.4)	0/49	0	108
5	B	B	A	(0.2 – 8.3)	4/49	8.2	114
6	B	B	C	(0.0 – 23.5)	8/49	16.3	124
7	B	B	D	(0.0 – 42.4)	9/49	18.4	117
8	B	B	E	(0.0 – 5.0)	2/49	4.1	118
9	C	C	A	(0.0 – 12.6)	10/49	20.4	117
10	C	C	B	(0.0 – 45.5)	12/49	18.4	122
11	C	C	D	(0.0 – 29.0)	12/49	18.4	120
12	C	C	E	(0.1 – 4.9)	2/49	4.1	125
13	D	D	A	(0.0 – 10.5)	9/49	18.4	117
14	D	D	B	(0.0 – 91.2)	10/49	20.4	118
15	D	D	C	(0.0 – 12.0)	7/49	14.3	124
16	D	D	E	(0.0 – 2.2)	0/49	0	128
17	E	E	A	(0.0 – 5.8)	5/49	10.2	114
18	E	E	B	(0.0 – 12.3)	9/49	18.4	114
19	E	E	C	(0.0 – 9.1)	13/49	25.5	123
20	E	E	D	(0.0 – 3.6)	0/49	0	122