

Impact of past climatic changes and resource availability on the population demography of three food-specialist bees

SIMON DELLICOUR,*§ DENIS MICHEZ,† JEAN-YVES RASPLUS‡ and PATRICK MARDULYN*

*Evolutionary Biology and Ecology, Université Libre de Bruxelles, av. FD Roosevelt 50, 1050 Brussels, Belgium, †Laboratory of Zoology, Research Institute of Biosciences, University of Mons, Place du Parc 20, 7000 Mons, Belgium, ‡INRA, UMR CBGP (INRA/IRD/Cirad/Montpellier SupAgro), Montferrier-sur-Lez, France

Abstract

Past climate change is known to have strongly impacted current patterns of genetic variation of animals and plants in Europe. However, ecological factors also have the potential to influence demographic history and thus patterns of genetic variation. In this study, we investigated the impact of past climate, and also the potential impact of host plant species abundance, on intraspecific genetic variation in three codistributed and related specialized solitary bees of the genus *Melitta* with very similar life history traits and dispersal capacities. We sequenced five independent loci in samples collected from the three species. Our analyses revealed that the species associated with the most abundant host plant species (*Melitta leporina*) displays unusually high genetic variation, to an extent that is seldom reported in phylogeographic studies of animals and plants. This suggests a potential role of food resource abundance in determining current patterns of genetic variation in specialized herbivorous insects. Patterns of genetic variation in the two other species indicated lower overall levels of diversity, and that *M. nigricans* could have experienced a recent range expansion. Ecological niche modelling of the three *Melitta* species and their main host plant species suggested a strong reduction in range size during the last glacial maximum. Comparing observed sequence data with data simulated using spatially explicit models of coalescence suggests that *M. leporina* recovered a range and population size close to their current levels at the end of the last glaciation, and confirms recent range expansion as the most likely scenario for *M. nigricans*. Overall, this study illustrates that both demographic history and ecological factors may have contributed to shape current phylogeographic patterns.

Keywords: coalescent simulations, demographic history, food specialization, intraspecific diversity, phylogeography, phytophagous insects, population fragmentation

Received 8 April 2014; revision received 24 December 2014; accepted 15 January 2015

Introduction

Current patterns of genetic variation characterizing a species are known to reflect its evolutionary history. In particular, strong evidence suggests that past climate changes and geographic barriers had a major impact on

the genetic variation of most species in Europe (Hewitt 2004; Avise 2009). Genetic variation in Europe is usually interpreted in terms of demographic history, associated with available refugia during glacial episodes of the Quaternary, and to subsequent recolonization of Europe from these refugia at the onset of warmer periods (Taberlet *et al.* 1998; Hewitt 2004). However, ecological or life history traits could also play a major role in shaping current patterns of intraspecific variation but has been rarely investigated (but see e.g. Fine *et al.* 2013). One

Correspondence: Simon Dellicour,

E-mail: Simon.Dellicour@zoo.ox.ac.uk

§ Present address: Department of Zoology, University of Oxford, South Parks Road, Oxford, OX1 3PS, UK

such factor that could potentially impact genetic diversity and structure is the abundance of food resources. Widespread and abundant food resources would be expected to increase connectivity among populations and overall population size of the feeding species. This may translate into higher levels of genetic diversity compared to those of related species with similar history but smaller populations and more fragmented ranges. In other words, resource abundance may influence demographic parameters such as population sizes and migration rates, which are the proximal factors directly impacting genetic variation. Specialized herbivorous insects, especially those with a narrow diet, appear ideal to explore this hypothesis, because their distribution closely matches with those of their respective host plant species. Contrasting related specialist herbivores feeding on plants that vary in abundance should help to assess the hypothesized relationship between resource availability and (i) intraspecific genetic diversity and (ii) population structure.

We compared the genetic variability found in five independent genes (one mitochondrial and four nuclear loci; total of ~3750 bp) across the range of three species of specialist solitary bees codistributed in Europe. The three solitary bees, *Melitta leporina*, *M. nigricans* and *M. tricincta*, are sibling species (i.e. forming a clade in the genus *Melitta*; Dellicour *et al.* 2014a), are morphologically and physiologically similar (Celary 2006; Michez & Eardley 2007) and display comparable nesting behaviour (i.e. sharing nesting constraints; Celary 2006; Nilsson & Alves-dos-Santos 2009) and body size (i.e. similar dispersal and foraging abilities; Zurbuchen *et al.* 2010). On the other hand, they are characterized by different specialized diets (Michez *et al.* 2008). *Melitta leporina* is specialized on the Fabaceae plant family and visits common and largely abundant plants from different genera in Europe, mainly alfalfa and clovers (Michez *et al.* 2008). By comparison, *M. nigricans* and *M. tricincta* are specialized on resources that are much less abundant, including plants from the genus *Lythrum* (Lythraceae) and *Odontites* (Schrophulariaceae), respectively (Michez *et al.* 2008). Following Müller & Kuhlmann (2008), *M. leporina* can be considered as a broad oligolectic species (i.e. specialist species collecting pollen from more than one plant genus, but always from the same plant family), while *M. nigricans* and *tricincta* are narrow oligolectics (i.e. specialist species collecting pollen from host plant species always belonging to the same plant genus). Females of *M. leporina* have access to the pollen of a higher number of species, including crops (legumes are grown on 180 million Ha of earth's arable surface and account for 27% of the crop production; Graham & Vance 2003), which increases the relative abundance of its floral resources. The resulting difference in global host plant abundance

is at least partially reflected by their relative AOO (areas of occupancy index) in Europe (Table S1, Fig. S1, Supporting information): the AOO of *M. leporina* main host plant species (*M. sativa*, *T. pratense* and *T. repens*; ~340 000 km²) is more than two times larger than that of the host plant of *M. nigricans* (*L. salicaria*, AOO of ~150 000 km²) and up to approximately seven times larger than that of the main host plant of *M. tricincta* (*O. vernus*, AOO of ~50 000 km²). In addition, our field observations and previous studies suggest that when host plants from more than one species are present in the same region, those from *M. leporina* are locally much more abundant (e.g. Fortel *et al.* 2014). Bees are generally believed to be strongly associated with their host plants, and a shift in plant diet often translates into a new speciation (Michez *et al.* 2008; Müller & Kuhlmann 2008; Praz *et al.* 2008a,b). For the purpose of interpreting their current pattern of genetic variation, we have assumed that the plant diet of each *Melitta* species studied has remained stable over time, at least for the recent time period during which the observed patterns were formed. As floral specialization is the main ecological difference between the three *Melitta* species, they offer an interesting model to investigate the hypothesized relationship between resource availability and (i) intraspecific diversity and (ii) population structure. Following previous results obtained by Kelley *et al.* (2000), Packer *et al.* (2005) and Habel *et al.* (2009), we wished to test the prediction that the abundance of food resource is positively correlated with intraspecific diversity and negatively correlated with population structure (i.e. genetic differentiation among populations).

We therefore analysed and compared the phylogeographic patterns of the three *Melitta* species with the aim of investigating their respective demographic history, in relation to past climate changes and to the potential influence of differences in food resource abundance. More specifically, because past demography is known to impact patterns of genetic variation, we first (i) estimated and compared their current patterns of genetic variability, and (ii) inferred present and past species ranges by developing ecological niche models of both the insects and their host plants. We then derived historical hypotheses from the inferred species ranges and evaluated them using the acquired genetic variation data and spatially explicit models of coalescence.

Materials and methods

Sampling and sequencing

We sampled *M. leporina*, *M. nigricans* and *M. tricincta* in 57, 38 and 23 localities, respectively, across their geographic distribution (Fig. S2, Table S2, Supporting

information), described as Palearctic for *M. leporina* and West Palearctic for *M. nigricans* and *M. tricincta* (Michez & Eardley 2007). Genomic DNA of 506 haploid male individuals (208 *M. leporina*, 172 *M. nigricans* and 126 *M. tricincta*) were extracted using the Qiagen DNeasy® Blood & Tissue kit. Half a thorax per specimen was grounded in the Qiagen ATL buffer and incubated overnight with proteinase K at 56 °C. The remaining DNA-extraction steps were conducted as described in the manufacturer's protocol. We sequenced 454 samples for a ~900-base-pair (bp)-long fragment of the mitochondrial cytochrome oxidase I (COI) gene, 377 samples for a ~400-bp-long fragment of the wingless (WgL) gene, 368 samples for a ~750-bp-long fragment of the sodium-potassium adenosine triphosphatase (NaK) gene, 399 samples for a ~850-bp-long fragment of the RNA polymerase II (RNAP) gene and 425 samples for a ~850-bp-long fragment of the long-wavelength rhodopsin (opsin) gene, the last four genes being all protein-coding and of nuclear origin. All fragments were PCR-amplified with a TrueStart Hot Start Taq DNA polymerase, following the guidelines in the manufacturer's protocol (Fermentas International Inc.). The COI fragment was amplified (annealing temperature of 50.5 °C) using primers Sim (5'-AAT ATT TAT GAA TTC RGG RTC AGG-3') or Ron and Pat (Simon *et al.* 1994), the WgL fragment (annealing temperature of 63.5 °C) with primers Bee-wg-For1 or Bee-wg-For2 and Lep-wg2a-Rev (Almeida & Danforth 2009), the NaK fragment (annealing temperature of 66.5 °C) with primers NaKfor2 and NaKrev2 (Michez *et al.* 2009), the RNAP fragment (annealing temperature of 57 °C) with primers Polfor2a and Polrev2a (Danforth *et al.* 2006) and the opsin fragment (annealing temperature of 59 °C) with primers For3 and Rev4a (Danforth *et al.* 2004). See Table S3 (Supporting information) for detailed PCR conditions of each primer pair. All haplotype sequences are available from GenBank under Accession nos. KM922006-KM922543.

DNA sequence analyses

Sequences were aligned using the MUSCLE algorithm (Edgar 2004) implemented in CODONCODE ALIGNER (v. 3.7.1.1, CodonCode Corporation). These alignments were checked manually and pruned at both 5'- and 3'-ends. Small gaps identified in the opsin alignment were considered as missing data for the purpose of inferring a haplotype network. Median-joining networks (Bandelt *et al.* 1999) were inferred for each gene fragment using the software NETWORK 4.6.6 (available at <http://www.fluxus-engineering.com>) with $\epsilon = 0$. NETWORK was also used to identify hypervariable sites using the estimated maximum number of mutations at each site.

DNA sequences were analysed using BLAST (Zhang *et al.* 2000) to identify coding regions and the positions of introns. Translation of protein-coding sequences into amino acid sequences was performed in MACCLADE 4.08a (Maddison & Maddison 2000), while characterization of DNA polymorphisms and identification of nonsynonymous substitutions were conducted with DNASP 5.0 (Librado & Rozas 2009). Due to the high number of reticulations in the RNAP and opsin networks, we tested the possibility of recombination having occurred in all the DNA fragments, using the PHI test (Bruen *et al.* 2006) implemented in the software SplitsTree 4.12.3 (Huson & Bryant 2006). The location of the potential recombination events were estimated using (i) the method of Hudson & Kaplan (1985) implemented in the software DNASP (Librado & Rozas 2009) and (ii) all seven methods available in the RDP4 package (Maynard Smith 1992; Padidam *et al.* 1999; Gibbs *et al.* 2000; Martin & Rybicki 2000; Posada & Crandall 2001; Martin *et al.* 2005, 2010; Boni *et al.* 2007). Assuming homoplasy (i.e. convergent nucleotide substitutions) is responsible for reticulations in the RNAP and opsin networks (see Discussion), and for the purpose of analysing population fragmentation, we inferred new networks for these two loci, after deleting sites for which the estimated maximum number of substitutions was above 20. This limit was chosen to allow the RNAP and opsin networks to display a number of reticulations similar to those characterizing the networks of the two other nuclear loci. These two modified DNA sequence alignments are further referred as 'unsaturated' RNAP and opsin alignments.

Global genetic diversity and population structure

To analyse and compare current patterns of genetic diversity and population structure, we used SPADS 1.0 (Dellicour & Mardulyn 2014) and (i) estimated allelic richness (El Mousadik & Petit 1996) and nucleotide diversity (Nei & Li 1979) over the entire ranges of *M. nigricans* and *tricincta*, and over the European portion of the range of *M. leporina*, which corresponds to its co-distribution area with the two other species, for each locus separately, (ii) estimated the statistic $N_{ST} - G_{ST}$ (a measure of phylogeographic signal/structure; Pons & Petit 1996) over the same geographic ranges and (iii) performed a multilocus version of a SAMOVA (spatial analysis of molecular variance; Dupanloup *et al.* 2002) as well as a COI-only SAMOVA (as it was the only locus systematically associated with a significant phylogeographic signal; see Results). Allelic richness for a species is here calculated as the expected number of different haplotypes in a subsample of n sequences, n being the smallest sample size when comparing samples for each

species considered in the study (i.e. in this case, the number of sequences obtained for *M. tricincta*, the species for which we have the smallest sample size for each locus; $n = 98$ for COI, 86 for WgL, 77 for NaK, 76 for RNAP and 76 for opsin). It is important to note that allelic richness and nucleotide diversity are both corrected for unequal sample size, which makes possible their comparison among species for which we have different sample sizes. The statistical test for the difference between N_{ST} and G_{ST} (which highlights the extent of the phylogeographic signal) was based on 10 000 random permutations of haplotypes. $N_{ST} - G_{ST}$ was estimated for each locus separately and over all nuclear loci using multilocus weighted averages (Weir & Cockerham 1984): one multilocus weighted average based on the four original DNA sequence alignments and another calculated after the original RNAP and opsin alignments were replaced by their corresponding 'unsaturated' alignments. The SAMOVA method assigns populations to K groups based on geographic vicinity and sequence similarity. The most likely structure corresponds to the partition of populations maximizing among-group variation as measured by AMOVA (analysis of molecular variance) Φ_{CT} statistic (Excoffier *et al.* 1992). We performed ten independent runs of 10 000 simulated annealing steps for each K value varying from two to the number of sampled populations for each species. Compared to a SAMOVA based on a single locus, the multilocus version of this algorithm analyses all loci simultaneously using a multilocus weighted average Φ_{CT} to compare two successive iterations (Weir & Cockerham 1984; Dellicour & Mardulyn 2014).

Geographic distribution of genetic variation

To explore genetic variation across the species range, we created graphs displaying the geographic distribution of genetic diversity and genetic differentiation using an extension of a method developed by Miller (2005), based on an interpolation procedure (inverse distance-weighted interpolation; Watson & Philips 1985; Watson 1992), with R functions available in SPADS (Dellicour & Mardulyn 2014). For genetic diversity, interpolation is based on diversity values directly estimated at each sampling point, but for genetic differentiation, interpolation is based on distance values assigned at midpoints of each edge of a connectivity network built between the sampling points (e.g. a Delaunay triangulation). For each species, we generated three graphs based on different genetic diversity indices and two fragmentation graphs based on different interindividual distances. All diversity and distance indices used are independent from sample size, which varies among species (Fig. S2, Table S2, Supporting information). The three diversity

statistics were (i) allelic richness A_R of each population (the estimation of the expected number of different haplotypes in a subsample of n sequences for a given population, n being the smallest number of sequences obtained for a sampling site with more than one sampled sequence; El Mousadik & Petit 1996), (ii) nucleotide diversity π of each population (the average number of nucleotide differences per site between any two DNA sequences chosen randomly in a given population; Nei & Li 1979) and (iii) relative nucleotide diversity π_R of each population (the nucleotide diversity within this population divided by the nucleotide diversity within the group formed by all other populations; Mardulyn *et al.* 2009). For the estimation of these diversity statistics, only the populations with more than one sampled sequence per locus were considered. Note that for allelic richness, since $n = 2$ for all three species, interpolation surfaces are directly comparable. The two interindividual distances were (i) an interindividual distance based on allelic frequencies and treating the entire sequence as the locus (as defined by Miller 2005), and (ii) an interindividual distance based on distance DNA sequence mismatches averaged over the different loci (see SPADS manual for further details). All surfaces were generated using three values for distance weighting parameter a (1, 5 and 10), and fragmentation surfaces were both based on a Delaunay triangulation connectivity network. All interpolations were based on great circle geographic distances (i.e. distances on the surface of the earth) measured in kilometres and estimated using the R package 'fields' (Fields Development Team 2006). Furthermore, as advised by Miller *et al.* (2006), we performed the distance interpolations using residual distances derived from the linear regression of genetic vs. geographic distances (Manni *et al.* 2004). In the end, all generated surfaces were superimposed to a map of Europe (i.e. on the codistributed range of the three species).

Ecological niche modelling

For the purpose of investigating demographic history of each species and highlighting differences among them, current and last glacial maximum (LGM) distributions of the three *Melitta* species and their main host plant species (*Medicago sativa*, *Trifolium pratense* and *T. repens* for *M. leporina*, *Lythrum salicaria* for *M. nigricans* and *Odontites vernus* for *M. tricincta*) were inferred using the maximum entropy method implemented in MAXENT 3.3.3 (Phillips *et al.* 2006; Phillips & Dudík 2008). The model implemented in MAXENT minimizes relative entropy between two probability densities: one estimated from species occurrence data and one estimated from the landscape (Elith *et al.* 2011). The present distribution was estimated first and the result used to project the

species distribution on past climate layers using two distinct LGM models: (i) the CCSM (Community Climate System Model) and (ii) the MIROC (Model for Interdisciplinary Research On Climate). These inferences were based on ten pieces of bioclimatic data and 1923 pieces of occurrence data for *M. leporina*, 435 for *M. nigricans*, 594 for *M. tricineta* and, for the five host plant species, on all the GBIF records used for estimating AOO indices (see Introduction, Fig. S1 and Table S1, Supporting information). Nineteen bioclimatic variables (Bio1–Bio19) at a 2.5 arc-minutes (~5 km) resolution were initially extracted from the WorldClim database (WorldClim 1.4; Hijmans *et al.* 2005) for the current time (~1950–2000) and LGM models (CCSM and MIROC). Available bioclimatic variables on WorldClim for the CCSM and MIROC are provided by the PMIP2 database (Paleoclimate Modelling Intercomparison Project Phase II, Braconnot *et al.* 2007). Relationships among these bioclimatic variables were evaluated using Pearson's correlation coefficients. To avoid collinear variables (Pearson coefficient > 0.9), nine variables were discarded. The resulting set of ten selected variables included annual mean diurnal temperature range (Bio2), isothermality (Bio3), annual temperature range (Bio7), mean temperature of the wettest quarter (Bio8), mean temperature of the driest quarter (Bio9), mean temperature of warmest quarter (Bio10), mean temperature of coldest quarter (Bio11), annual precipitation (Bio12), precipitation seasonality (Bio15) and precipitation in the warmest quarter (Bio18). Ten replicates were performed for each analysis, from which we derived an average distribution. We used the default convergence threshold (10⁻⁵), 10 000 iterations and a 'random seed' to generate a random partition of our localities into training (90%) and test (10%) localities. These random partitions were used to test the model. The inferences were evaluated using the area under the ROC (receiver operating characteristic) curve. These AUC (area under the curve) values are commonly used to assess the MAXENT estimation performance (see Marske *et al.* 2009, 2011).

Comparison of spatially explicit demographic scenarios

Results obtained from ecological niche modelling suggest different levels of range size, range fragmentation and species abundance. In addition, an interpolation graph based on relative nucleotide diversity obtained for *M. nigricans* indicates a potential recent range expansion for this species (see Results). We then developed spatially explicit coalescence models using the framework implemented in PHYLOGEOSIM 1.0 (Dellicour *et al.* 2014b), formally translating the ambiguity in the results (ecological niche modelling and interpolation of relative nucleotide diversity) into alternative historical

hypotheses over the evolution of the three species ranges. We simulated DNA sequences following these models, with PHYLOGEOSIM, where the range at each generation is defined as a two-dimensional grid in which each cell is considered a single population. Going backward in time, at each generation (given these species are univoltine, 1 generation = 1 year), populations can host coalescence events between two or more gene copies and/or exchange gene copies with neighbour populations. Similarly to a method developed by Currat *et al.* (2004; see also Ray *et al.* 2010), PHYLOGEOSIM performs a preliminary forward simulation to estimate backward simulation parameters (backward migration rates and effective population sizes for each generation). Sequence length and a fixed number of mutations were specified for each simulated locus, matching the parameters of our real DNA sequence alignments. For opsin and RNAP, observed data included an extreme level of polymorphism, and multiple nucleotide substitutions at some sites have probably occurred. Therefore, we rather simulated the 'unsaturated' versions of the opsin and RNAP data sets (see above), for which we were able to specify a number of mutations, and compared them to their corresponding observed unsaturated data sets.

We developed seven historical hypotheses (see Results). The first six (scenarios A–F) were developed from MAXENT outputs as follows: (i) resolution of the occurrence probability matrices returned by MAXENT was decreased by a factor 20, to prevent RAM limitations (occurring with extremely large inputs); this corresponds to an increase in cell size from 0.0416 to 0.8392 decimal degrees; (ii) occurrence probabilities were then multiplied by a constant value to generate matrices of maximal effective cell sizes (see Knowles & Alvarado-Serrano 2010 for a similar approach). As the order of magnitude of maximum effective sizes corresponding to a grid cell is difficult to estimate for these species, we compared three very different values for this constant in our simulations: 10 000, 100 000 and 1 000 000 for nuclear genes, and a third of these numbers for the mitochondrial gene COI (haplodiploid species resulting theoretically in a mitochondrial effective size three times lower than the nuclear effective size). We also compared different values of forward migration rates connecting adjacent cells on the grid ($m_f = 0.01, 0.001$ and 0.0001) and different values of reproduction rates ($t_R = 2$ and 5). See Table 1 for a summary of the different parameter values tested. Scenarios A–F each include a current interglacial layer, a LGM layer corresponding to the inferred distribution under CCSM or MIROC model for the three species and a previous interglacial layer that was assumed identical to the current interglacial (thus based on the occurrence probability matrices returned by MAXENT for the current time period).

Table 1 Combined P -values obtained from the comparison between real and simulated data sets with 1000 simulations per set of parameters and per locus, and a reproduction rate $t_R = 2$. $N_{e,cell}$ refers to maximum effective population size assigned to each cell on all two-dimensional grids of the corresponding model, and p_{MAXENT} to the probability of occurrence in each grid cell as inferred from $MAXENT$ results

Scenario	$N_{e,cell} = p_{MAXENT} \times 10^4$		$N_{e,cell} = p_{MAXENT} \times 10^5$		$N_{e,cell} = p_{MAXENT} \times 10^6$	
	$m_f = 0.01$	$m_f = 0.0001$	$m_f = 0.01$	$m_f = 0.0001$	$m_f = 0.01$	$m_f = 0.0001$
<i>Melitta leporina</i>						
A. LGM: CCSM <i>M. leporina</i>	<0.001	0.013	<0.001	<0.001	<0.001	<0.001
B. LGM: MIROC <i>M. leporina</i>	<0.001	0.125	0.003	<0.001	<0.001	<0.001
C. LGM: CCSM <i>M. nigricans</i> *	<0.001	0.016	0.021	<0.001	0.069	<0.001
<i>Melitta nigricans</i>						
A. LGM: CCSM <i>M. leporina</i> *	<0.001	<0.001	0.000	<0.001	<0.001	<0.001
C. LGM: CCSM <i>M. nigricans</i>	<0.001	0.025	0.001	0.001	<0.001	<0.001
D. LGM: MIROC <i>M. nigricans</i>	<0.001	0.013	0.000	<0.001	<0.001	<0.001
G. Recent range expansion	0.002	0.149	0.081	0.043	0.138	0.024
<i>Melitta tricincta</i>						
A. LGM: CCSM <i>M. leporina</i> *	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
E. LGM: CCSM <i>M. tricincta</i>	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
F. LGM: MIROC <i>M. tricincta</i>	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001

*Indicates scenarios for which interglacial spatial grids were built from $MAXENT$ results of another species. P -values in bold are significant (>0.05). See Table S11 (Supporting information) for combined P -values obtained with a reproduction rate $t_R = 5$

Although we initially considered all generated hypotheses (scenarios A-F) for each of the three investigated *Melitta* species, only a subset of these were actually tested (see Table 1) because some of these scenarios are very similar (e.g. ‘CCSM’ layers for *M. nigricans* and *tricincta*, see Results). Overall, we tested, for each species, two scenarios built using the past distributions inferred from different LGM models (CCSM and MIROC; i.e. scenarios A and B for *M. leporina*, C and D for *M. nigricans* and E and F for *M. tricincta*), as well as a third scenario that was initially built for one of the two other species and that clearly differed from the two-first scenarios tested (scenario C tested for *M. leporina*, involving a highly restricted range at the LGM, and scenario A for *M. nigricans* and *tricincta*, involving a large range at the LGM). For each species, we thus tested three scenarios involving notable differences in range size and probabilities of occurrence at the LGM.

The last hypothesis (scenario G) was specifically built to test the assumption of a recent range expansion for *M. nigricans*. This scenario was suggested by the geographic mapping of genetic diversity which revealed a lower overall diversity compared to the other two species, except for a restricted area associated with much higher relative nucleotide diversity (see Results). This small area characterized by high genetic variation may indicate the origin of a hypothesized range expansion. For this additional scenario, the layer for the current interglacial period was built as for the other scenarios, but the oldest glaciation layer was replaced by a more recent layer that allows the presence of this species only in a restricted portion of the range, considered as the area of origin for the expansion. The location of the area of origin was thus set according to the interpolation graph of relative nucleotide diversity for the species. Note that scenarios A–F also include an initial common layer with a restricted area of origin preceding the previous interglacial period (i.e. from the TMRCA, the time to the most recent common ancestor, to the beginning of the previous interglacial period). The location of this restricted area was arbitrarily set in the same place as the area of origin for the expansion set in scenario G. The purpose of this initial layer was to allow all gene copies to coalesce in a reasonably short time, which was necessary to simulate sequence data similar to our observed sequences. Indeed, failing to include this layer with a restricted range resulted in TMRCA values unrealistically large, because remaining gene copies at the end of the simulation take much too long to coalesce.

Each set of simulation parameters was used to simulate 1000 data sets under each tested species–hypothesis combination, and for each locus. To take stochastic variation associated with a forward simulation into account, we ran a new forward simulation for every 10

coalescence (backward) simulations. Additional information about coalescence simulations and the comparison of the simulated and observed sequence data through calculation of a series of summary statistics are given as Supplementary Information (Appendix S1 and Table S4, Supporting information).

Results

Genetic variation within and among species

Allele networks inferred from one mitochondrial (cytochrome oxidase I; Fig. 1) and four nuclear (wingless, sodium–potassium adenosine triphosphatase, RNA polymerase II, long-wavelength rhodopsin; Fig. 2) gene fragments highlight a stark contrast among the patterns of genetic variation characterizing the three species. Each locus displays a much higher genetic diversity for *M. leporina*. *M. tricincta* harbours systematically more alleles than *M. nigricans* for the four nuclear loci, although this pattern is reversed in the mitochondrial fragment. The observation of this pattern is confirmed by global diversity indices estimated from European samples (Table S5, Supporting information): allelic richness is always highest for *M. leporina* (and for *M. tricincta* for the NaK locus) and is also systematically

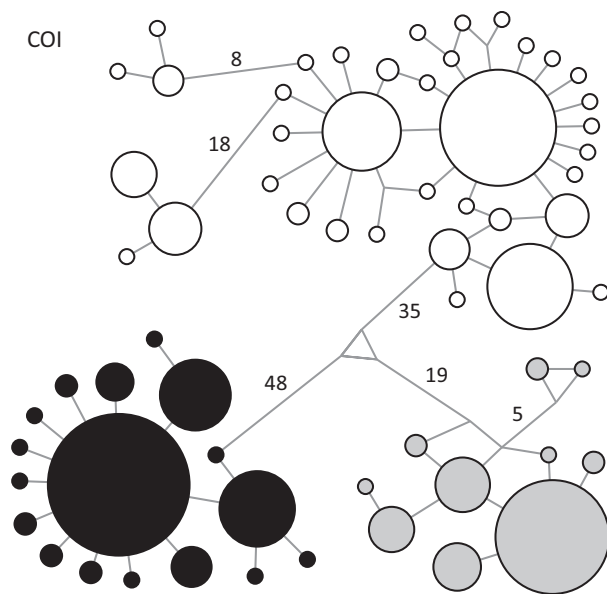


Fig. 1 Median-joining networks for the mitochondrial gene fragment COI. Each sequenced haplotype is represented by a circle, the size of which is proportional to its overall frequency. Each line in the network represents a single mutational change. In some cases, numbers are used to indicate a higher number of mutations. Haplotype colours correspond to the three *Melitta* species: white for *M. leporina*, black for *M. nigricans* and grey for *M. tricincta*.

higher for *M. tricincta* than for *M. nigricans* (except with COI, for which the reverse is true). Regarding nucleotide diversity, differences remain notable between *M. nigricans* and the two other species, but *M. leporina* and *tricincta* often display similar values.

Of particular note is the extremely high number of alleles and reticulations (revealing ambiguous, i.e. equally parsimonious, connections) observed for *M. leporina* in both the RNAP and opsin genes. Reticulations could result from multiple convergent mutations and/or intragene recombinations. Tests for the detection of recombination events provided ambiguous results, some of them suggesting the occurrence of a few recombination events, others detecting no recombination signal (Table S6, Supporting information). Removing the nucleotides associated with the highest number of mutations in the RNAP and opsin data sets resulted in more easily interpretable networks (Fig. 2; see Material and methods), displaying clear historical information related to the evolutionary relationships among alleles and species. For all genes, the majority of polymorphic sites are located on the third nucleotide positions of coding regions, and also within introns in the case of opsin. Only few substitutions between sampled sequences were identified as nonsynonymous (Table S7, Supporting information).

While the mitochondrial gene COI appears strongly differentiated among species (Fig. 1, see also Figs S3 and S4, Supporting information), forming three potentially monophyletic groups of alleles separated by a large number of mutations, this level of differentiation is much less pronounced in the nuclear loci. Except for opsin, rooting of allele networks for nuclear genes does not result in three monophyletic species. Even in the case of opsin, the species are separated from each other by a single mutation. Incomplete lineage sorting among the three species is the most likely explanation to account for this pattern in the three other nuclear loci. For WgL and NaK, some alleles are even shared by two species (Fig. 2). Likewise, when focusing on intraspecies variation, the mitochondrial data set is the only one associated with a strong and significant phylogeographic signal for each of the three species (i.e. significant $N_{ST} - G_{ST}$ values, Table S8, Supporting information). The association between genealogy and geography can be further examined in this network with Figs S3 and S4 (Supporting information). For most nuclear loci, the phylogeographic signal for each species is not significant (Table S8, Supporting information). The stronger phylogeographic signal in the mitochondrial fragment is probably related to its associated lower effective population size. With no clear maximum Φ_{CT} value, the multilocus SAMOVA (Fig. S5, Supporting information) did not highlight a particular population partition. This

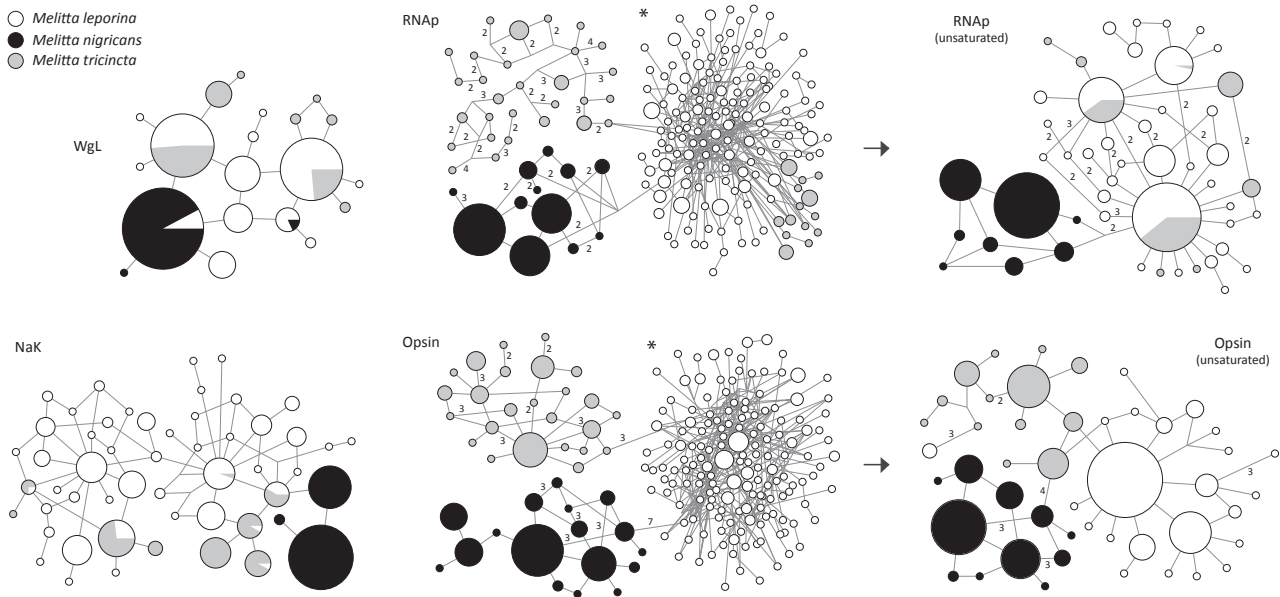


Fig. 2 Median-joining networks for the four nuclear gene fragments (WgL, NaK, RNAP and opsin) used in this study. Haplotype colours correspond to the three *Melitta* species: white for *M. leporina*, black for *M. nigricans* and grey for *M. tricincta*. The 'unsaturated' RNAP and opsin networks were inferred using new DNA alignments free of polymorphic sites with an estimated maximum number of mutations higher than 20 (see the text for further detail). (*) Due to the important number of haplotype connections, numbers of mutations higher than one are not indicated among *M. leporina* haplotypes in RNAP and opsin networks.

is likely to be related to the lack of phylogeographic signal found in nuclear fragments. On the contrary, SAMOVAS based on the mitochondrial gene only identified K values (number of partitions) associated with high Φ_{CT} in most cases, except in the case of the European samples of *M. leporina*. Partitions associated with highest Φ_{CT} values in this analysis are detailed in Table S9 (Supporting information).

Geographic distribution of genetic variation

Interpolation graphs generated with parameter $a = 5$ are presented in Fig. 3 for the population distance based on allelic frequencies ($IID1$), the population diversity based on allelic richness A_R and relative nucleotide diversity π_R , and in Fig. S6 (Supporting information) for the population distance based on DNA sequence mismatches ($IID2$) and the population diversity based on nucleotide diversity π (see also Figs S7 and S8, Supporting information for equivalent graphs generated with $a = 1$ and 10). Global patterns remain similar with different a , but the level of spatial variation within each surface logically increases with a . While large differences in allelic richness were highlighted among the three species at the global level, interpolation surfaces (Figs 3B and S6, Supporting information) confirm this trend at the local level for both allelic richness and nucleotide diversity: *M. leporina* displays the

highest local genetic diversity, followed by *M. tricincta* and then by *M. nigricans*. Interestingly, *M. nigricans* is the only species whose distribution map includes an area where the relative nucleotide diversity π_R (Fig. 3C), a statistic comparing the genetic diversity present in one locality to that of the rest of the distribution, is significantly higher than on the rest of the range, which could suggest that this area has been the starting point for a range expansion of the species. While *M. leporina* displays the highest level of genetic diversity, *M. tricincta* presents the highest level of population fragmentation, when comparing interpolation graphs of interindividual distances based on allelic frequencies (Fig. 3A), although this contrast is much lower when taking genetic distances into account (confirming the low or absent phylogeographic signal found with the nuclear loci; Fig. S6, Supporting information). Note that these graphs (Figs 3A and S6, Supporting information) focus on interindividual distances among adjacent populations and are complementary to population or phylogeographic structure analyses focusing on the entire range (i.e. SAMOVA, $N_{ST} - G_{ST}$).

Ecological niche modelling

Distributions estimated with MAXENT (Figs 4 and S9, Supporting information) were associated with good AUC (area under the curve) values (all values > 0.823 ;

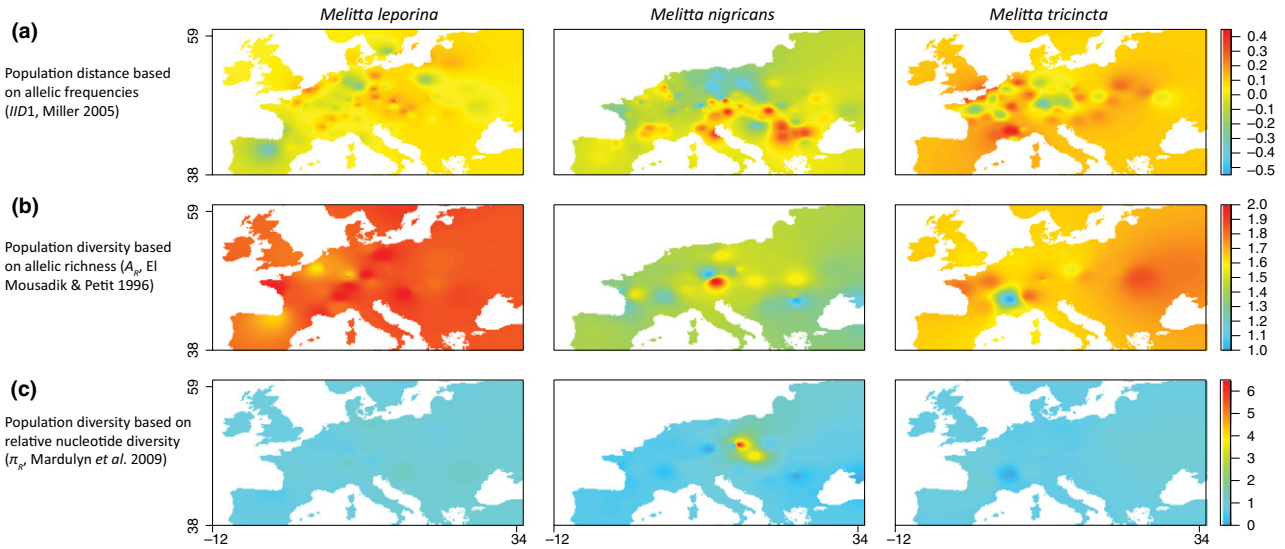


Fig. 3 Interindividual distance and population diversity interpolating graphs generated with a distance weighting parameter $a = 5$ (see Fig. S6 for the population distance based on DNA sequence mismatches and for the population diversity based on nucleotide diversity). Graphs based on interindividual distances were generated with the GDisPAL function and those based on population diversity with the GDivPAL function (see text). Interindividual distances were computed with the unsaturated rNAP and opsin data sets. See also Figs S7 and S8 (Supporting information) for equivalent graphs generated with $a = 1$ and 10.

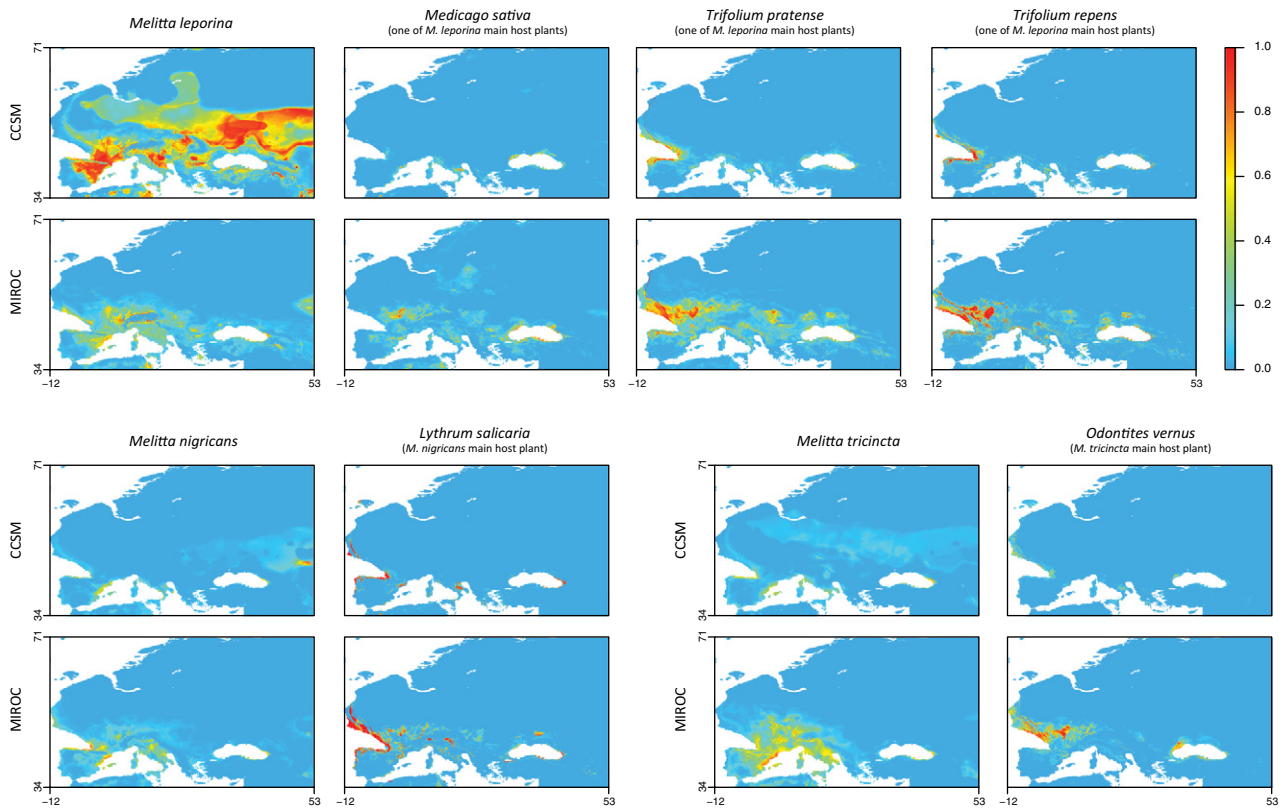


Fig. 4 Last glacial maximum (LGM) ranges of the three *Melitta* species and their main host plant species inferred using MAXENT and based on estimated LGM palaeoclimatic data ('CCSM' and 'MIROC' models). Probabilities of occurrence are shown with a colour scale.

Table S10, Supporting information), indicating that all models performed significantly better than random predictions (Phillips *et al.* 2006). As displayed in Fig. 4, ecological niche modelling of past species distributions mainly revealed that the outcome is highly dependent on the LGM model (CCSM or MIROC) used for the inference. For example, in the case of *M. leporina*, the LGM distribution inferred with CCSM predicts a much larger range with globally higher probability of occurrence, while for *M. tricineta*, it is the LGM distribution inferred with MIROC that suggests larger range and occurrence probabilities. Similar differences can be observed among inferred LGM distributions for the main host plant species. To understand the origin of these differences, we mapped and compared the relative values of each selected bioclimatic variable between the two LGM models (Fig. S10, Supporting information). This comparison highlighted the predicted values for the annual temperature range (Bio7) and the isothermality (Bio3) as the main contributors to the difference between models. Overall, comparison of LGM predictions for the three *Melitta* species indicates that the range of *M. leporina* would have been larger than that of the two other species during the last glaciation. On the other hand, no clear difference in LGM range size appears among the main host plant species.

Comparison of spatially explicit demographic scenarios

The seven tested historical hypotheses are displayed on Figure 5. Comparisons between real and simulated data sets are summarized in Table 1 and lead to different results for each of the three species. For *M. leporina*, scenarios B and C, associated with a stronger reduction in the range during the last glaciation, are clearly favoured (two *P*-values > 0.1) over scenario A (all *P*-values < 0.015). For *M. nigricans*, scenario G (recent range expansion) is clearly favoured over all others (two *P*-values > 0.1). However, the range expansion hypothesis was initially derived directly from the observation of the genetic variation, that is the same data that were compared to the simulated data to test the hypothesis. A certain level of circularity is therefore inherent to this approach, and the confidence in the result should be lower than when all considered hypotheses were derived entirely from external sources, that is, here, from ecological niche modelling alone. On the other hand, none of the tested scenarios for *M. tricineta* appear compatible with the data (all *P*-values < 0.001). Note that in this context, each *P*-value represents the estimated probability that the simulated scenario has generated the observed corresponding pattern of genetic variation.

Discussion

The unusually high level of polymorphism and reticulation observed for *Melitta leporina* in RNAP and opsin networks could potentially be explained by three hypotheses: balancing selection acting on the two loci, unusually high intragene recombination and an unusually large population size that led to the maintenance of a high level of ancestral polymorphism in that species. It is usually considered that genes subject to balancing selection are quite rare in a genome. Although such genes may occur more commonly than previously thought (Burgess 2013; Leffler *et al.* 2013; Delph & Kelly 2014), the likelihood that two of four randomly selected nuclear loci are under such selection seems fairly low. While recombination can occur within any nuclear gene fragment, its probability of occurrence is usually small. To generate the high number of reticulations displayed by the RNAP and opsin allele networks through recombination only, clearly more than one recombination event is needed (14–19 recombination events according to the method of Hudson & Kaplan 1985; see Table S6, Supporting information), which seems also highly unlikely in such a short DNA fragment, even more so in independent (again, randomly selected) loci. On the contrary, the fact that such extremely high diversity is displayed by two independent loci favours the demographic hypothesis, and supports an unusually large historical population size for *M. leporina*. Even though such a strong contrast in diversity is not found in the other loci, perhaps because they are characterized by lower mutation rates, genetic patterns of diversity in these other loci also unambiguously demonstrate *M. leporina* to be the most polymorphic species. In fact, after deleting the sites associated with the highest number of substitutions in RNAP and opsin, located mainly in third codon positions or introns, the networks of all four nuclear loci become highly similar.

Whether we favour the balancing selection or demographic hypotheses, the unusual RNAP and opsin networks are at least partially explained by a number of mutations so large that multiple substitutions have occurred at many sites in the two sequences. While similarly 'saturated' historical signals are often observed in phylogenetic studies, especially those comparing distantly related species, this saturation is surprising at the intraspecies level. To the best of our knowledge, such extreme patterns are seldom reported in the phylogeographic literature and therefore appear highly unusual. Networks similar to those presented in Fig 2 were found only in a few cases involving renowned hyper-variable regions, including, for example, networks inferred from (i) the control region of the mitochondrial genome (e.g. Bamshad *et al.* 2001; de Bruyn *et al.* 2009;

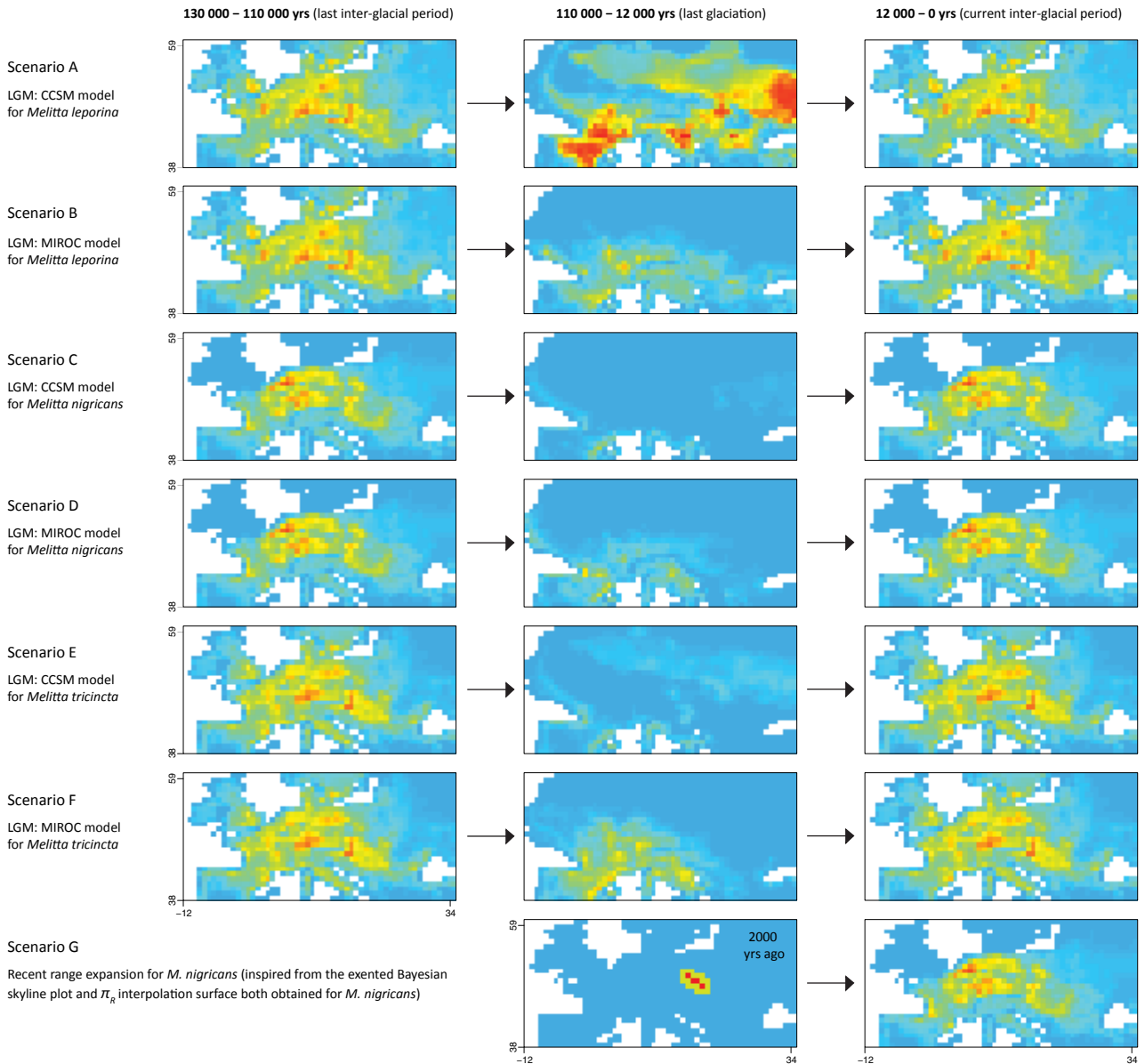


Fig. 5 Spatially explicit models built to test alternative evolutionary hypotheses with PHYLOGEOSIM 1.0. Historical scenarios are described by successive geographic layers (the oldest on the left). The colour scale indicates probabilities of occurrence inferred with MAXENT, used to define matrices of maximum effective sizes in the model (see text). Effective size matrices were generated by multiplying MAXENT probabilities by 10 000, 100 000 and 1 000 000.

Zimmermann *et al.* 2011), (ii) human Y-chromosome microsatellite haplotypes (e.g. Sengupta *et al.* 2006; Shou *et al.* 2010; Kim *et al.* 2011) and (iii) the human MHC (major histocompatibility complex) locus, in which high diversity is maintained through balancing selection (e.g. Andrés *et al.* 2010).

The initial prediction that higher genetic diversity should be observed in the species with the most abundant plant resources (i.e. *M. leporina*) is verified separately with all three loci, and concordance among

independent loci is a strong signal in favour of a historical explanation. Because the three studied species are sibling and characterized by a highly similar life history and morphology (i.e. size, which was found to be correlated with dispersal distance in other bees; Gathmann & Tscharntke 2002; Araújo *et al.* 2004; Greenleaf *et al.* 2007), local and/or global differences in genetic diversity among them are likely to reflect, at least partially, differences in their population sizes. Genetic variation in *M. leporina* therefore suggests the maintenance of a

larger effective population size over a long period of time, while genetic variation in *M. nigricans* rather suggests a recent and dramatic size increase from a small effective ancestral population. A large effective population size maintained over a long period of time in *M. leporina* is expected to result in a longer coalescence process leading backwards in time to the most recent common ancestor of alleles at a locus, allowing a high number of mutations, and possibly an unusually high number of recombination events, to occur at each locus. Thus, the demographic hypothesis could in fact be compatible with the occurrence of recombination, and the unusual allele networks highlighted for RNAp and opsin could be explained by a combination of convergent nucleotide substitutions and recombination events.

The higher level of genetic diversity and low population structure characterizing the range of *M. leporina* could be related to its access to abundant plant resources, both in terms of number of species and in terms of abundance of each host plant species. While it is true that the range of *M. leporina* is larger than that of its two sister species, this difference alone appears insufficient to explain the extent of the difference in genetic diversity between them, at least if we assume a linear relationship between species range and population size. Indeed, the range of *M. leporina* is approximately three times larger than that of the other species, but its estimated current effective size is approximately 10 times greater than that of *M. tricineta*. The widespread nature of the host plants of *M. leporina*, together with the general ability of bees to disperse (e.g. Delli-cour *et al.* 2014c), could thus result in a large network of connected populations encompassing the entire range of the species. Availability of host plants would have persisted over a long period of time. While ecological niche modelling estimates suggested that host plant ranges of all *Melitta* species were strongly reduced at the LGM, interglacial levels of plant resource abundance may have recovered earlier for *M. leporina* than for the two other *Melitta* species, or, if the different host plants have recovered their distribution simultaneously after the LGM, the *M. leporina* bee populations themselves may have recovered more quickly than those of the two other *Melitta* species. In other words, differences in genetic diversity among species could in fact reflect differential population dynamics after the last glaciation. The two other bees, although phylogenetically and ecologically very close to and codistributed with *M. leporina*, are associated with less common plants and currently inhabit a more fragmented distribution. If host plant abundance and coverage influenced genetic variation, we would expect these species to be characterized by lower diversity and stronger population structure. These expectations are clearly met

for *M. tricineta*, associated with the least abundant resources, as shown in Figs 3 and S6 (Supporting information). The pattern of genetic variation displayed by *M. nigricans* is less straightforward to explain. Although also displaying less genetic diversity, the species does not show a clear pattern of stronger population differentiation/structure, when compared to *M. leporina*. Because Fig. 3 highlights a restricted region of high genetic diversity lying partially over Poland and Czech Republic, it was *a priori* tempting to hypothesize that a recent increase in effective size associated with a range expansion occurred, originating from this location.

These historical hypotheses, suggested by our observations of genetic variation across the range of the three *Melitta* species, were further tested with simulations. Sequence data were simulated according to each hypothesis using a spatially explicit model of coalescence and compared to observed sequence data. For *M. leporina*, the best supported hypotheses are the ones that indeed implement a strong reduction in its range and population size at the LGM, as suggested by the MIROC model estimation for that species, and a relatively quick recovery 12 000 years ago, in which it reached range conditions and population sizes similar to those we see today. The large reduction in its range at the LGM is further supported by the distribution estimates of the main host plant species of this bee for the same time period, which also suggest highly restricted distributions, assuming *M. leporina* could not survive without its current host plant species. On the other hand, the best supported hypothesis for *M. nigricans* involves a much more recent range expansion/population size increase that occurred around 2000 years ago. In that species, the combination of a less available plant resource with a recent history of range expansion from a small portion of the current range may explain the lower genetic diversity and population structure observed. Finally, for *M. tricineta*, while no tested model conditions generated simulated data close to our observed data, the lower level of host plant abundance correlates with lower genetic diversity (compared with *M. leporina*). Lower genetic variation encountered in that species could then be associated with the lower abundance of its plant resource alone.

While it is clear that plant resources of *M. leporina* are currently much more abundant than those of the two other bees, no direct evidence is available to estimate their relative abundance during the last 12 000 years, as pollen diagrams (i.e. pollen records in deposits) and vegetation of insect-pollinated plants are poorly related (Djamali *et al.* 2009). However, we could hypothesize a constant relative dominance of the legume family over the genera *Lythrum* and *Odontites* during this time frame. In fact, it is worth-noting that the domestication

of many Fabaceae species in the Mediterranean basin began approximately 8000 years ago (Graham & Vance 2003). *M. leporina* may thus have benefitted from this domestication, because it may have increased the availability of its floral resources at the time.

Previous work had already suggested that generalist-feeding herbivorous insects tend to display higher genetic variation than specialized species (Kelley *et al.* 2000; Packer *et al.* 2005; Habel *et al.* 2009). Our data further suggest that even for specialist-feeding species, host plant abundance could be a crucial factor influencing population size and fragmentation, and thus genetic variation, across their range. More importantly, this study illustrates that ecological and historical factors are somewhat related and can both potentially influence current patterns of intraspecific genetic variation. *M. leporina* and *M. tricineta* appear to have both experienced a strong range contraction at the LGM, yet the two species are characterized by different patterns of genetic variation, which we can only assume are related to the different availability of their plant resources. *M. tricineta* and *M. nigricans* seem to have benefited from quite similar abundance of plant resources during their recent history, yet their patterns of population structure appear widely divergent, as *M. nigricans* has probably gone through a relatively recent range expansion, possibly due to a different response to past climate changes, or to another unidentified cause. The influence of past climate change should thus be considered along with that of important ecological factors when interpreting phylogeographic data. In some phylogeographic studies, deciphering the impact of past climate change and of other ecological factors may in fact become challenging, as it may be difficult to identify the important ecological or life history traits that played an important role for the focal species, as well as to disentangle the respective effects of all these factors. While the development of ecological niche modelling has helped us to integrate many climatic factors in our interpretation of genetic variation, this approach is limited by the availability of predictive variables for the LGM. Indeed, predictions inferred for the last glaciation are often based on only a few climatic variables measuring temperature and/or humidity across the range, as was done in the present study. While these variables are definitely important, other abiotic variables, and even important biotic variables, such as ecological interactions among species, are seldom taken into account. In this study, for example, the estimated distribution of *M. leporina* host plants at the LGM seems much more restricted than that estimated for the insect based on climatic variables alone. Taking the estimated host plants distribution at the LGM into account would result in a much smaller estimated range for the insect. Even when

the climatic variables included in the ecological niche model are the main determinants of the current pattern of genetic variation, the resolution of bioclimatic data available for the estimation may be insufficient (Dellacour *et al.* 2014d). Outputs of ecological niche modelling should therefore be interpreted with caution, even though they are interesting to identify relevant demographic scenarios that can be tested with genetic variation data (e.g. Carstens & Richards 2007).

Overall, the results presented here suggest that climate alone cannot explain the observed patterns of genetic variation, as these differ among the three co-distributed bee species investigated, which have presumably been subject to the same climatic conditions. While codistributed species in Europe have often displayed divergent phylogeographic patterns (e.g. Taberlet *et al.* 1998), the differences highlighted here among three related bee species, which share most of their traits (with the notable exception of host plant species), are compatible with an important influence of food resource abundance, as suggested in a few previous studies (Kelley *et al.* 2000; Packer *et al.* 2005; Habel *et al.* 2009). Further studies comparing closely related species that differ by their resource abundance or other potentially important ecological or life history traits are needed before strong general conclusions over the influence of these factors can be inferred. These results have also practical implications for conservation biology, as the abundance of the host plant could significantly influence the genetic diversity, and thus the potential for long-term survival, of specialist insect herbivores.

Acknowledgements

We thank three anonymous reviewers and editor for providing helpful comments and advices. We are grateful to A. Lachaud, A.Q. Zhang, B. Cederberg, C. Philippe, C. Saure, D.W. Baldock, E. Dufrêne, E. Lamas Delgado, F. Mayer, F. Vyghen, H.K. Schwenninger, H. Özbek, I. Raemakers, J. Gellhaus, J.S. Ascher, J. Straka, L.A. Nilsson, L. Crépin, L. Fortel, M. Augbert, M. Terzo, M. Vanderplanck, N.J. Vereecken, O. Berg, O. Dayez, P. Bogusch, S. Bellens, S. Gadoum, S.W. Droege, T. De Meulemeester, T. Lecocq, T. Levchenko, V.G. Radchenko and W. Celary for their invaluable help in collecting specimens. We wish to thank G. Genson for her technical assistance in the laboratory, Steven Bachman for his help regarding the IUCN AOO indices as well as I. Meusnier, H. Darras and F. Massonet for their useful comments on this study. We are also grateful to S. Hill for her careful proofreading of the final version of our manuscript. This research project was funded by the Belgian *Fonds pour la Recherche Scientifique* (FRS-FNRS; FRFC 2.4613.10) and by the Royal Academy of Belgium (Agathon De Potter funds). Computational resources have been provided by the High Performance Computing Centre cofunded by ULB and VUB (HPC cluster 'Hydra'). S.D. was supported by a grant from the *Fonds pour la Recherche dans l'Industrie et*

l'Agriculture (FRIA) and by an award from the Fonds David and Alice Van Buuren.

References

- Almeida EAB, Danforth BN (2009) Phylogeny of colletid bees (Hymenoptera: Colletidae) inferred from four nuclear genes. *Molecular Phylogenetics and Evolution*, **50**, 290–309.
- Andrés AM, Dennis MY, Kretzschmar WW *et al.* (2010) Balancing selection maintains a form of ERAP2 that undergoes nonsense-mediated decay and affects antigen presentation. *PLoS Genetics*, **6**, e1001157.
- Araújo ED, Costa M, Chaud-Netto J, Fowler HG (2004) Body size and flight distance in stingless bees (Hymenoptera: Meliponini): inference of flight range and possible ecological implications. *Brazilian Journal of Biology*, **64**, 563–568.
- Avise JC (2009) Phylogeography: retrospect and prospect. *Journal of Biogeography*, **36**, 3–15.
- Bachman S, Moat J, Hill AW, de laTorre J, Scott B (2011) Supporting red list threat assessments with GeoCAT: Geospatial conservation assessment tool. *ZooKeys*, **150**, 117–126.
- Bamshad M, Kivisild T, Watkins WS *et al.* (2001) Genetic evidence on the origins of Indian caste populations. *Genome Research*, **11**, 994–1004.
- Bandelt HJ, Forster P, Röhl A (1999) Median-joining networks for inferring intraspecific phylogenies. *Molecular Biology and Evolution*, **16**, 37–48.
- Boni MF, Posada D, Feldman MW (2007) An exact nonparametric method for inferring mosaic structure in sequence triplets. *Genetics*, **176**, 1035–1047.
- Braconnot P, Otto-Bliesner B, Harrison S *et al.* (2007) Results of PMIP2 coupled simulations of the Mid-Holocene and last glacial maximum – Part 1: experiments and large-scale features. *Climate of the Past*, **3**, 261–277.
- Bruen TC, Philippe H, Bryant D (2006) A simple and robust statistical test for detecting the presence of recombination. *Genetics*, **172**, 2665–2681.
- de Bruyn M, Hall BL, Chauke LF, Baroni C, Koch PL, Hoelzel AR (2009) Rapid response of a marine mammal species to holocene climate and habitat change. *PLoS Genetics*, **5**, e1000554.
- Burgess DJ (2013) Evolution: a matter of balance. *Nature Reviews Genetics*, **14**, 240–241.
- Carstens BC, Richards CL (2007) Integrating coalescent and ecological niche modeling in comparative phylogeography. *Evolution*, **61**, 1439–1454.
- Celary WD (2006) Biology of the solitary ground-nesting bee *Melitta leporina* (Panzer, 1799) (Hymenoptera: Apoidea: Melittidae). *Journal of the Kansas Entomological Society*, **79**, 136–145.
- Curat M, Ray N, Excoffier L (2004) SPLATCHE: a program to simulate genetic diversity taking into account environmental heterogeneity. *Molecular Ecology Notes*, **4**, 139–142.
- Danforth BN, Brady SG, Sipes SD, Pearson A (2004) Single-copy nuclear genes recover Cretaceous-age divergences in bees. *Systematic Biology*, **53**, 309–326.
- Danforth BN, Fang J, Sipes S (2006) Analysis of family-level relationships in bees (Hymenoptera: Apiformes) using 28S and two previously unexplored nuclear genes: CAD and RNA polymerase II. *Molecular Phylogenetics and Evolution*, **39**, 358–372.
- Dellicour S, Mardulyn P (2014) SPADS 1.0: a toolbox to perform spatial analyses on DNA sequence datasets. *Molecular Ecology Resources*, **14**, 647–651.
- Dellicour S, Lecocq T, Kuhlmann M, Mardulyn P, Michez D (2014a) Molecular phylogeny, biogeography, and host plant shifts in the bee genus *Melitta* (Hymenoptera: Anthophila). *Molecular Phylogenetics and Evolution*, **70**, 412–441.
- Dellicour S, Kastally C, Hardy OJ, Mardulyn P (2014b) Comparing phylogeographic hypotheses by simulating DNA sequences under a spatially explicit model of coalescence. *Molecular Biology and Evolution*, **31**, 3359–3372.
- Dellicour S, Mardulyn P, Hardy OJ, Hardy C, Roberts SPM, Vereecken NJ (2014c) Inferring the mode of colonisation of a rapid range expansion from multi-locus DNA sequence variation. *Journal of Evolutionary Biology*, **27**, 116–1329.
- Dellicour S, Fearnley S, Lombal A *et al.* (2014d) Inferring the past and present connectivity across the range of a North American leaf beetle: combining ecological-niche modeling and a geographically explicit model of coalescence. *Evolution*, **68**, 2371–2385.
- Delph LF, Kelly JK (2014) On the importance of balancing selection in plants. *New Phytologist*, **201**, 45–56.
- Djamali M, de Beaulieu JL, Campagne P *et al.* (2009) Modern pollen rain-vegetation relationships along a forest-steppe transect in the Golestan National Park, NE Iran. *Review of Palaeobotany and Palynology*, **153**, 272–281.
- Dupanloup I, Schneider S, Excoffier L (2002) A simulated annealing approach to define the genetic structure of populations. *Molecular Ecology*, **11**, 2571–2581.
- Edgar RC (2004) MUSCLE: Multiple sequence alignment with improved accuracy and speed. pp. 728–729.
- El Mousadik A, Petit RJ (1996) High level of genetic differentiation for allelic richness among populations of the argan tree [*Argania spinosa* (L.) Skeels] endemic to Morocco. *Theoretical and Applied Genetics*, **92**, 832–839.
- Elith J, Phillips SJ, Hastie T, Dudík M, Chee YE, Yates CJ (2011) A statistical explanation of MaxEnt for ecologists. *Diversity and Distributions*, **17**, 43–57.
- Excoffier L, Smouse PE, Quattro JM (1992) Analysis of molecular variance inferred from metric distances among DNA haplotypes: application to human mitochondrial DNA restriction data. *Genetics*, **131**, 479–491.
- Fields Development Team (2006). *fields: Tools for Spatial Data*. National Center for Atmospheric Research, Boulder, Colorado. <http://www.cgd.ucar.edu/Software/Fields>
- Fine PVA, Zapata F, Daly DC *et al.* (2013) The importance of environmental heterogeneity and spatial distance in generating phylogeographic structure in edaphic specialist and generalist tree species of *Protium* (Burseraceae) across the Amazon Basin. *Journal of Biogeography*, **40**, 646–661.
- Fortel L, Henry M, Guilbaud L *et al.* (2014) Decreasing abundance, increasing diversity and changing structure of the wild bee community (Hymenoptera: Anthophila) along an urbanization gradient. *PLoS ONE*, **9**, e104679.
- Gathmann A, Tschardt T (2002) Foraging ranges of solitary bees. *Journal of Animal Ecology*, **71**, 757–764.
- Gibbs MJ, Armstrong JS, Gibbs AJ (2000) Sister-Scanning: a Monte Carlo procedure for assessing signals in recombinant sequences. *Bioinformatics*, **16**, 573–582.

- Graham PH, Vance CP (2003) Legumes: importance and constraints to greater use. *Plant Physiology*, **131**(3), 872–877.
- Greenleaf SS, Williams NM, Winfree R, Kremen C (2007) Bee foraging ranges and their relationship to body size. *Oecologia*, **153**, 589–596.
- Habel JC, Meyer M, Schmitt T (2009) The genetic consequence of differing ecological demands of a generalist and a specialist butterfly species. *Biodiversity and Conservation*, **18**, 1895–1908.
- Hewitt GM (2004) Genetic consequences of climatic oscillations in the Quaternary. *Philosophical Transactions of the Royal Society of London Series B-Biological Sciences*, **359**, 183–195.
- Hijmans RJ, Cameron SE, Parra JL, Jones PG, Jarvis A (2005) Very high resolution interpolated climate surfaces for global land areas. *International Journal of Climatology*, **25**, 1965–1978.
- Hudson RR, Kaplan NL (1985) Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*, **111**, 147–164.
- Huson DH, Bryant D (2006) Application of phylogenetic networks in evolutionary studies. *Molecular Biology and Evolution*, **23**, 254–267.
- Kelley ST, Farrell BD, Mitton JB (2000) Effects of specialization on genetic differentiation in sister species of bark beetles. *Heredity*, **84**, 218–227.
- Kim SH, Kim KC, Shin DJ *et al.* (2011) High frequencies of Y-chromosome haplogroup O2b-SRY465 lineages in Korea: a genetic perspective on the peopling of Korea. *Investigative Genetics*, **2**, 10.
- Knowles L, Alvarado-Serrano DF (2010) Exploring the population genetic consequences of the colonization process with spatio-temporally explicit models: insights from coupled ecological, demographic and genetic models in montane grasshoppers. *Molecular Ecology*, **19**(17), 3727–3745.
- Leffler EM, Gao Z, Pfeifer S *et al.* (2013) Multiple instances of ancient balancing selection shared between humans and chimpanzees. *Science*, **340**, 1578–1582.
- Librado P, Rozas J (2009) DnaSP v5: A software for comprehensive analysis of DNA polymorphism data. *Bioinformatics*, **25**, 1451–1452.
- Maddison DR, Maddison WP (2000) *MacClade 4: Analysis of Phylogeny and Character Evolution. Version 4.0*. Sinauer Associates, Sunderland, Massachusetts.
- Manni F, Guerard E, Heyer E (2004) Geographic patterns of (genetic, morphologic, linguistic) variation: how barriers can be detected by using Monmonier's algorithm. *Human Biology*, **76**, 173–190.
- Mardulyn P, Mikhailov YE, Pasteels JM (2009) Testing phylogeographic hypotheses in a Euro-Siberian cold-adapted leaf beetle with coalescent simulations. *Evolution*, **63**, 2717–2729.
- Marske KA, Leschen RAB, Barker GM, Buckley TR (2009) Phylogeography and ecological niche modelling implicate coastal refugia and trans-alpine dispersal of a New Zealand fungus beetle. *Molecular Ecology*, **18**, 5126–5142.
- Marske KA, Leschen RAB, Buckley TR (2011) Reconciling phylogeography and ecological niche models for New Zealand beetles: looking beyond glacial refugia. *Molecular Phylogenetics and Evolution*, **59**, 89–102.
- Martin DP, Rybicki E (2000) RDP: detection of recombination amongst aligned sequences. *Bioinformatics*, **16**, 562–563.
- Martin DP, Posada D, Crandall KA, Williamson C (2005) A modified bootscan algorithm for automated identification of recombinant sequences and recombination breakpoints. *AIDS Research and Human Retroviruses*, **21**, 98–102.
- Martin DP, Lemey P, Lott M, Moulton V, Posada D, Lefevre P (2010) RDP3: a flexible and fast computer program for analyzing recombination. *Bioinformatics*, **26**, 2462–2463.
- Maynard Smith J (1992) Analyzing the mosaic structure of genes. *Journal of Molecular Evolution*, **34**, 126–129.
- Michez D, Eardley C (2007) Monographic revision of the bee genus *Melitta* Kirby 1802 (Hymenoptera: Apoidea: Melittidae). *Annales de la Société entomologique de France (n.s.)*, **43**, 379–440.
- Michez D, Patiny S, Rasmont P, Timmermann K, Vereecken NJ (2008) Phylogeny and host plant evolution in Melittidae s.l. (Hymenoptera: Apoidea). *Apidologie*, **39**, 146–162.
- Michez D, Patiny S, Danforth BN (2009) Phylogeny of the bee family Melittidae (Hymenoptera: Anthophila) based on combined molecular and morphological data. *Systematic Entomology*, **34**, 574–597.
- Miller MP (2005) Alleles In Space (AIS): computer software for the joint analysis of interindividual spatial and genetic information. *Journal of Heredity*, **96**, 722–724.
- Miller MP, Bellinger MR, Forsman ED, Haig SM (2006) Effects of historical climate change, habitat connectivity, and vicariance on genetic structure and diversity across the range of the red tree vole (*Phenacomys longicaudus*) in the Pacific Northwestern United States. *Molecular Ecology*, **15**, 145–159.
- Müller A, Kuhlmann M (2008) Pollen hosts of western palaeartic bees of the genus *Colletes* (Hymenoptera: Colletidae): the Asteraceae paradox. *Biological Journal of the Linnean Society*, **95**, 719–733.
- Nei M, Li WH (1979) Mathematical model for studying genetic variation in terms of restriction endonucleases. *Proceedings of the National Academy of Sciences of the United States of America*, **76**, 5269–5273.
- Nilsson LA, Alves-dos-Santos I (2009) The oligolectic solitary bee *Melitta tricincta* Kirby, 1802 (Sw. rödtöppebi) in Sweden (Hymenoptera, Apoidea, Melittidae). *Entomologisk Tidskrift*, **130**, 85–98.
- Packer L, Zayed A, Grixti JC *et al.* (2005) Conservation genetics of potentially endangered mutualisms: reduced levels of genetic variation in specialist versus generalist bees. *Conservation Biology*, **19**, 195–202.
- Padidam M, Sawyer S, Fauquet CM (1999) Possible emergence of new geminiviruses by frequent recombination. *Virology*, **265**, 218–225.
- Phillips SJ, Dudík M (2008) Modeling of species distributions with Maxent: new extensions and a comprehensive evaluation. *Ecography*, **31**, 161–175.
- Phillips SJ, Anderson RP, Schapire RE (2006) Maximum entropy modeling of species geographic distributions. *Ecological Modelling*, **190**, 231–259.
- Pons O, Petit RJ (1996) Measuring and testing genetic differentiation with ordered versus unordered alleles. *Genetics*, **144**, 1237–1245.
- Posada D, Crandall KA (2001) Evaluation of methods for detecting recombination from DNA sequences: computer simulations. *Proceedings of the National Academy of Sciences of the United States of America*, **98**, 13757–13762.

- Praz CJ, Müller A, Dorn S (2008a) Host recognition in a pollen-specialist bee: evidence for a genetic basis. *Apidologie*, **39**, 547–557.
- Praz CJ, Müller A, Dorn S (2008b) Specialized bees fail to develop on non-host pollen: do plants chemically protect their pollen? *Ecology*, **89**, 795–804.
- Ray N, Currat M, Foll M, Excoffier L (2010) SPLATCHE2: a spatially explicit simulation framework for complex demography, genetic admixture and recombination. *Bioinformatics*, **26**, 2993–2994.
- Sengupta S, Zhivotovsky LA, King R *et al.* (2006) Polarity and temporality of high-resolution Y-chromosome distributions in India identify both indigenous and exogenous expansions and reveal minor genetic influence of Central Asian pastoralists. *American Journal of Human Genetics*, **78**, 202–221.
- Shou WH, Qiao EF, Wei CY *et al.* (2010) Y-chromosome distributions among populations in Northwest China identify significant contribution from Central Asian pastoralists and lesser influence of western Eurasians. *American Journal of Human Genetics*, **55**, 314–322.
- Simon C, Frati F, Beckenbach A, Crespi B, Liu H, Flook P (1994) Evolution, weighting, and phylogenetic utility of mitochondrial gene sequences and a compilation of conserved polymerase chain reaction primers. *Annals of the Entomological Society of America*, **87**, 651–701.
- Taberlet P, Fumagalli L, Wust-Saucy AG, Cosson JF (1998) Comparative phylogeography and postglacial colonization routes in Europe. *Molecular Ecology*, **7**, 453–464.
- Watson DF (1992) *Contouring: A Guide to the Analysis and Display of Spatial Data*. Pergamon Press, New York, New York.
- Watson DF, Philips GM (1985) A refinement of inverse distance weighted interpolation. *Geo-processing*, **2**, 315–327.
- Weir BS, Cockerham CC (1984) Estimating F-statistics for the analysis of population structure. *Evolution*, **38**, 1358–1370.
- Zhang Z, Schwartz S, Wagner L, Miller W (2000) A greedy algorithm for aligning DNA sequences. *Journal of Computational Biology*, **7**, 203–214.
- Zimmermann B, Röck A, Huber G, Krämer T, Schneider PM, Parson W (2011) Application of a west Eurasian-specific filter for quasi-median network analysis: sharpening the blade for mtDNA error detection. *Forensic Science International: Genetics*, **5**, 133–137.
- Zurbuchen A, Landert L, Klaiber J, Müller A, Hein S, Dorn S (2010) Maximum foraging ranges in solitary bees: only few individuals have the capability to cover long foraging distances. *Biological Conservation*, **143**, 669–676.

S.D., D.M. and P.M. designed the study. S.D. and D.M. collected specimens. S.D. performed the experiments and analysed the data. J.Y.R. contributed to reagents and materials. S.D., D.M. and P.M. interpreted the results. S.D. wrote the initial draft of the manuscript, and all authors contributed substantially to revisions.

Data accessibility

DNA sequences are deposited in GenBank under Accession nos. KM922006–KM922543. DNA sequence alignments for all individuals and loci, detailed sampling information and occurrence data used for niche modelling are deposited in Dryad Digital Repository (doi:10.5061/dryad.0g2f5).

Supporting information

Additional supporting information may be found in the online version of this article.

Table S1 Number of distinct European GBIF records (Global Biodiversity Information Facility, www.gbif.org, visited on February 2014) and AOO index (areas of occupancy index calculated from European GBIF records with a cell width of 2 km, corresponding to the standard reference used for IUCN red list assessments; www.iucnredlist.org, see also Bachman *et al.* 2011) for the main host plants (Michez *et al.* 2008) of the three *Melitta* species examined.

Table S2 Sampling localities and distribution of haplotypes for the five loci used in this study.

Table S3 Detailed PCR conditions.

Table S4 List, brief description and reference of the different summary statistics computed on each simulated data set.

Table S5 Global allelic richness A_R (El Mousadik & Petit 1996) and global nucleotide diversity π (Nei & Li 1979) estimated on the overall sampling for each locus (and also considering only the European sampling of *M. leporina*).

Table S6 Results of the recombination tests: PHI test P -values (Bruen *et al.* 2006) and minimum numbers of recombination events (in parentheses) inferred with the method of Hudson & Kaplan (1985).

Table S7 Number of polymorphic sites per codon position or for introns, and number of identified non-synonymous mutations, calculated for each locus separately.

Table S8 $N_{ST} - G_{ST}$ estimated on sampled populations.

Table S9 Partition of populations obtained for highest Φ_{CT} values with SAMOVA based on the mitochondrial gene COI.

Table S10 AUC (area under the curve) scores for each replicate run performed with MAXENT.

Table S11 Combined P -values obtained from the comparison between real and simulated data sets with 1000 simulations per set of parameters and per locus and a reproduction rate $t_R = 5$.

Appendix S1 Spatially explicit simulations and comparison between simulated and observed data.

Fig. S1 Red contours of the European area used to select GBIF records (Global Biodiversity Information Facility; www.gbif.org) of host plants for AOO index computations (areas of occupancy, Table S1).

Fig. S2 Sampling localities of *Melitta leporina*, *Melitta nigricans* and *Melitta tricincta*.

Fig. S3 Median-joining networks for the COI gene fragment.

Fig. S4 Maps showing COI haplotypes distributions for the three *Melitta* species.

Fig. S5 Evolution of the combined Φ_{CT} statistic regarding the number of clusters 5 in multi-locus SAMOVA analyses and in SAMOVA analysis only based on the mitochondrial COI locus.

Fig. S6 Inter-individual distance and population diversity interpolating surfaces generated with a distance weighting parameter $a = 5$.

Fig. S7 Inter-individual distance and population diversity interpolating surfaces generated with a distance weighting parameter $a = 1$.

Fig. S8 Inter-individual distance and population diversity interpolating surfaces generated with a distance weighting parameter $a = 10$.

Fig. S9 Inferred current distributions of the three *Melitta* and their main host plant species inferred using MAXENT.

Fig. S10 Comparison of the CCSM and MIROC LGM models of the PMIP2 database for the 10 variables used in MAXENT analysis.