

Likelihood-based inference of population history from low coverage *de novo* genome assemblies

Jack Hearn¹, Graham N. Stone¹, Lynsey Bunnefeld-McInnes¹, James A. Nicholls¹, Nick H. Barton², Konrad

¹ Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH9 3JT, United Kingdom

² Institute of Science and Technology, Am Campus 1, A-3400 Klosterneuburg, Austria

Keywords: Statistical phylogeography, population divergence, admixture, maximum likelihood

Proofs to be sent to:

Konrad Lohse

Institute of Evolutionary Biology,

University of Edinburgh,

Edinburgh EH9 3JT,

United Kingdom

Phone +44 (0)131 650 5508

Email: konrad.lohse@ed.ac.uk

1 **Abstract**

2 Short-read sequencing technologies have in principle made it feasible to draw detailed inferences about the
3 recent history of any organism. In practice, however, such inferences remain challenging due to the difficulty
4 of genome assembly in most organisms and the lack of statistical methods powerful enough to allow discrim-
5 ination among recent, non-equilibrium histories. We address both the assembly and inference challenges.
6 We develop a bioinformatic pipeline for generating outgroup-rooted alignments of orthologous sequence
7 blocks from *de novo* low-coverage short-read data for a small number of genomes, and show how such
8 sequence blocks can be used to fit explicit models of population divergence and admixture in a numerical
9 likelihood framework. To illustrate our approach, we reconstruct the Pleistocene history of an oak-feeding
10 insect (the oak gallwasp *Biorhiza pallida*) which, in common with many other taxa, was restricted during
11 Pleistocene ice ages to a longitudinal series of southern refugia spanning the Western Palaearctic. Our anal-
12 ysis of blocks sampled from a single genome from each of three major glacial refugia reveals support for an
13 unexpected history dominated by recent admixture. Despite the fact that 80% of lineages are affected by ad-
14 mixture during the last glacial cycle, we are able to infer the deeper divergence history of these populations.
15 These inferences are robust to variation in block length, mutation model, and the sampling location of indi-
16 vidual genomes within refugia. This combination of *de novo* assembly and numerical likelihood calculation
17 provides a powerful framework for estimating recent population history that can be applied to any organism
18 without the need for prior genetic resources.

19 **Introduction**

20 Short read sequencing technologies have made it affordable to sequence entire genomes and such data
21 are now being used routinely to infer population history in humans (Gutenkunst *et al.*, 2009; Pickrell &

22 Pritchard, 2012; Durand *et al.*, 2011) and a small number of model species (Kulathinal *et al.*, 2009; Pool
23 *et al.*, 2012). However, it has remained frustratingly difficult to use such data to infer population history in
24 species for which no prior reference genome exists (Pool *et al.*, 2010; McCormack *et al.*, 2013). Historical
25 inference not only matters for understanding the evolutionary past of a particular species or ecological com-
26 munity (Stone *et al.*, 2012), but is a fundamental pre-requisite for testing for past selection in sequence data
27 (Sousa & Hey, 2013) and in analysing patterns of phenotypic evolution among populations (Stone *et al.*,
28 2011).

29 Short-read data pose several technical challenges for historical inference. First, we lack genomic re-
30 sources for the great majority of organisms, necessitating some form of *de novo* assembly. Second, it is
31 still not cost-effective to obtain genome-level data for many individuals – the general sampling design of
32 population genomic analyses. Third, most methods available for inferring population history from genomic
33 data are either based on allele frequency information (Gutenkunst *et al.*, 2009; Pickrell & Pritchard, 2012;
34 Durand *et al.*, 2011) – which requires large samples and ignores the historical signal contained in the higher
35 moments of branch length distributions – or simply do not scale up to genomic datasets (but see François
36 *et al.*, 2008). And finally, given that population history is evolutionarily recent by definition, the information
37 contained even in whole genomes is fundamentally limited by the time-scales of mutation and genetic drift
38 (Sousa & Hey, 2013; Hey & Machado, 2003).

39 Given these difficulties, it is unsurprising that the few studies to have used high-throughput data in
40 non-model organisms to date have all incorporated a "genomic reduction" step, resulting in sequencing of
41 only a small proportion of the target genome (McCormack *et al.*, 2013; Arnold *et al.*, 2013). Restriction-
42 site-associated DNA (RAD) sequencing approaches, whereby only a few hundred bases on either side of
43 a particular set of restriction sites are sequenced (Davey & Blaxter, 2010), have been successfully used to
44 detect population structure (Emerson *et al.*, 2010). Similarly, McCormack *et al.* (2012) have developed a

45 protocol that uses restriction digest to generate a reduced representation library of longer loci. However,
46 while these methods drastically simplify the assembly challenge, they involve additional wet lab protocols
47 for selecting loci and/or multiplexing of genomic libraries. Furthermore, the data generated are not neces-
48 sarily ideal for inferring intra-specific history. For example, RAD data typically consist of large numbers
49 of unlinked SNPs and so lack the much more detailed information about population history which is con-
50 tained in the distribution of genealogical branch lengths and accessible only via longer sequences that span
51 multiple, linked polymorphic sites.

52 An alternative strategy to "genomic reduction" is to work with whole genomes, but limit the analysis
53 to just a few individuals. This has the great advantage that efficient likelihood methods able to deal with
54 genome-scale data already exist (Wang & Hey, 2010; Li & Durbin, 2011; Lohse *et al.*, 2011; Gronau *et al.*,
55 2011). Discrete population models in particular, although less realistic than spatially continuous models
56 (Barton *et al.*, 2010), have become a standard in population genomics (Harris & Nielsen, 2013; Li & Durbin,
57 2011; Green *et al.*, 2010; Lohse & Frantz, 2013) because they are tractable and easy to interpret. Li &
58 Durbin (2011) have developed a hidden Markov approach for inferring past changes in effective population
59 size from just a single diploid genome. Similarly, Harris & Nielsen (2013) use the length distribution of
60 homozygous tracts in pairwise alignments to fit more complex histories of divergence and admixture between
61 two populations. However, these methods are currently restricted to histories involving just one or two
62 populations and rely on long, phased sequence blocks and hence near-complete assemblies, which remain
63 challenging to obtain for most organisms. In contrast, multi-locus methods (Hey & Nielsen, 2004; Hey,
64 2010; Lohse *et al.*, 2011), which assume a set of loci or sequence blocks each of which is short enough
65 to ignore recombination within it, but long enough to contain multiple polymorphic sites, are intuitively
66 appropriate for the fragmented genome assemblies available at low cost using low coverage paired-end
67 sequencing data.

68 Using such *de novo* assemblies in a multilocus framework requires matching and aligning orthologous
69 sequences both between individuals and between in- and outgroup species. Perhaps more importantly, one
70 needs to show that neither the assembly itself, nor the filtering steps involved in aligning sequences across
71 individuals and species, lead to systematic biases that affect the population genetic analyses. Here we present
72 a pipeline for generation of outgroup-rooted sequence blocks from small numbers of low coverage genomes,
73 and show the resulting data to be representative of the mosaic of genealogies in the genome. We then use
74 a numerical likelihood approach to illustrate the signal inherent in such data by reconstructing the Western
75 Palearctic population history of an oak feeding insect, the oak apple gallwasp *Biorhiza pallida*.

76 A suite of detailed studies have addressed phylogeographic patterns in Western Palearctic oak gallwasp
77 communities, both for the gall inducers (Stone & Sunnucks, 1993; Rokas *et al.*, 2001, 2003; Stone *et al.*,
78 2007; Challis *et al.*, 2007) and their parasitoid enemies (Hayward & Stone, 2006; Lohse *et al.*, 2010, 2012;
79 Nicholls *et al.*, 2010a,b). Both groups show genetic structure compatible with three major Pleistocene refu-
80 gial areas (Iberia; Italy and the Balkans; Asia Minor and Iran, Fig. 1) that broadly parallel those for deciduous
81 oaks (Petit *et al.*, 2003). Most species in both groups show patterns compatible with westwards range ex-
82 pansion into Europe from Asia during or before the Pleistocene (the 'Out of Anatolia' hypothesis; Rokas
83 *et al.* (2003); Challis *et al.* (2007); Stone *et al.* (2009); see also Connord *et al.*), a pattern also supported by
84 a recent meta-analysis of 19 parasitoid and 12 gallwasp species (Stone *et al.*, 2012). The only exception to
85 this pattern to date has been *Biorhiza pallida*, for which mitochondrial and ITS nuclear sequence data show
86 evidence of a deep east-west divide (Rokas *et al.*, 2001). This raises the question of how general the 'Out of
87 Anatolia' pattern is for all three trophic elements of this community (Stone *et al.*, 2009, 2012). Here we use
88 *Biorhiza pallida* as a case study for phylogenomic inference, and ask whether genome-level data support the
89 anomalous pattern for this species within the oak gallwasp community.

90 We focus on three refugial areas detailed above, referred to hereafter as the Western, Central and Eastern

91 refuge (Fig. 1). Modelling the relationship between these populations as a series of instantaneous divergence
92 and admixture events (Fig. 2) enables us to test the longitudinal directionality of initial occupation of refugia
93 and of admixture between them during subsequent periods of range expansion, whilst taking incomplete
94 lineage sorting (which is expected) into account. Some recent studies (Stone *et al.*, 2012; Lohse *et al.*, 2010)
95 have explicitly fitted a model of (E(C,W)) population divergence. However, given the recent timescale and
96 the limited number of sequenced loci available in most species, it has so far rarely been possible to quantify
97 the relative contributions of incomplete lineage sorting, divergence and admixture to genetic diversity cur-
98 rently present in refugial populations. We show that this can be achieved using low coverage data for only
99 one individual from each refugial region.

100 An inherent feature of model-based analyses is that it is necessary to limit the space of models to be
101 explored. For example, it has been possible to fit a very specific divergence and admixture model in which
102 admixture occurs only from the most anciently diverged population and after the most recent population
103 split (see Fig. 2) to human and Neandertal genomic data (Green *et al.*, 2010; Durand *et al.*, 2011; Lohse &
104 Frantz, 2013) because, in this case, the order of population divergence was known *a priori*. However, such
105 prior information does not exist for *B. pallida* or indeed most species. Importantly, we have no reason to
106 assume that population relationships are dominated by divergence (so are tree-like) rather than admixture in
107 the first place (Pickrell & Pritchard, 2012). Thus, to be able to fit divergence and admixture in general, one
108 needs to search model space more broadly. To this end, we have extended existing coalescent theory for the
109 "Neandertal model" developed by Lohse & Frantz (2013) to all possible histories involving unidirectional
110 admixture of a fraction f of lineages to or from the most anciently diverged population (Fig. 2). We feel
111 our model set balances biologically realistic scenarios with computational tractability. It is important to note
112 that unlike D statistics (Green *et al.*, 2010), which are defined relative to the majority topology (assumed to
113 reflect population divergence), our framework can deal with histories that are dominated by admixture (i.e.

114 $f > 0.5$). This allowed us to use the *B. pallida* genomic data to compare the support for a large number of
115 models ($n = 32$) and make detailed inferences about the history of this species. In particular we investigate
116 i) how well the *B. pallida* data can be explained by an (E,(C,W)) divergence history as inferred previously
117 for many species in the Western Palearctic, and ii) whether its history is dominated by deep population
118 divergence or recent admixture.

119 Given that our likelihood scheme is restricted to minimal samples of individuals, it is crucial to test how
120 representative individual genomes are of the long-term population relationships captured by discrete popu-
121 lation models. Clearly, in spatially continuous populations, local genetic structure emerges as a consequence
122 of the limited dispersal ability of individuals (Barton *et al.*, 2010), and so any model that approximates a
123 population occupying a large area as a panmictic unit must break down over recent time-scales. To address
124 this, we repeated our analyses using two different individuals from each refugium.

125 **Materials and Methods**

126 **Sequencing and sampling**

127 DNA was extracted from individual wasps using the Qiagen DNeasy kit. Like most Hymenoptera, *Biorhiza*
128 *pallida* has haploid males and diploid females; haploid males were selected for genome sequencing because
129 there is no need to phase alleles and SNP calling and estimation of sequencing error rates are greatly sim-
130 plified. However, we stress that our method does not rely on haploid genomes and can easily be applied to
131 unphased diploid data (Lohse & Frantz, 2013) (see Discussion). Illumina 50 and 100 base-pair paired-end
132 libraries (Table S2) were prepared using the Illumina paired-end DNA sample preparation kit, the DNA
133 sheared was using the Covaris S2 instrument and size selection was carried out on a 2% agarose TAE gel
134 (fragments with an average insert size of 300bp were excised). These were then sequenced on the GAI,

135 GAIIX and HiSeq2000 platforms at the NERC GenePool facility in Edinburgh (CONFIRM THIS). Note
136 that the different read lengths used are simply a result of technological improvements during the course of
137 this work and not a necessary part of our sequencing strategy. Short read data are deposited at the ENA
138 Sequence Read Archive (<http://www.ebi.ac.uk/ena/data/view/ERP002280>).

139 Each of five male in-group individuals (2 West, 2 Central and 1 Eastern see Table S1 and S2 and Fig.
140 1) was sequenced to a modal coverage of 1.5 fold per individual, yielding a total modal coverage of 7.45
141 for the *B. pallida* genome across all individuals. In each of the West and Central refuges we sampled
142 replicate individuals (referred to as Wa/Wb and Ca/Cb respectively, see Fig. 1) from sites 400km apart.
143 This separation is well above the dispersal ability of an individual gallwasp (Stone & Sunnucks, 1993) and
144 was intended to incorporate any impact of within-refuge population genetic structure. Sampling of replicate
145 individuals was not possible in Iran. To polarize mutations as ancestral or derived, we sequenced two diploid
146 female individuals from a closely related outgroup species *Belizinella gibbera* (to a total coverage of 5.76
147 fold across individuals). A breakdown of reads per individual is given in tables S2 and S3.

148 Multi-locus inference methods (Hey & Nielsen, 2004; Lohse *et al.*, 2011) assume a large number of
149 sequence blocks that are i) sampled at random from the genome, ii) short enough to ignore recombina-
150 tion within them and iii) in linkage equilibrium within populations. We developed a simple bioinformatic
151 pipeline (Fig. 3) that generates out-group-rooted alignments meeting the above criteria for a small number
152 of individuals. In short, the pipeline consists of three steps:

153 **i) Assembly**

154 Initial read quality was assessed using *FastQC* (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>).
155 Reads were quality trimmed at Q20 using *sickle* and adapter-trimmed using *scythe* and *cutadapt*. After
156 quality filtering, reads from all individuals in each species were combined to create species-specific meta-

157 assemblies using the CLC *de novo* assembler v4.06 (Dryad repository doi: XXX). These were used as a
158 reference for downstream bioinformatics processing (Table 1). For *B. pallida* we included a small amount of
159 50-base read data from a pool of six related (sib or half-sib) females that were sequenced as part of another
160 project (Table S2).

161 Reads from each *B. pallida* individual were mapped back to the meta-assembly using *Stampy* (run with
162 a substitution rate of 0.0025 and the -baq and -sensitive options) to create BAM files. Although a high
163 proportion (92 – 97 % depending on the individual) of read-pairs mapped to the reference meta-assemblies,
164 the percentage of "properly paired mappings", in which both read pairs align to the same contig, was much
165 lower (37%) (see Table S3). This is expected, given the low coverage per individual and the short read
166 length. Many read-mapping failures with *Stampy* are known to be due to reads that overlap the ends of
167 contigs. This is related to assembly coverage, specifically of the reads used to generate the initial contigs
168 (and contig trimming/minimum coverage filtering implemented by assemblers) rather than the reads used to
169 call SNPs.

170 **ii) Filtering**

171 *RepeatScout* and *RepeatMasker* were used to *de novo* predict and mask repeat regions in both the *B. pal-*
172 *lida* and *B. gibbera* meta-assemblies. This removed 51% and 34% of the *B. pallida* and *B. gibbera* meta-
173 assemblies respectively. The data were further filtered according to the following three criteria. First, we
174 searched for orthologous regions shared amongst in- and outgroup meta-assemblies using a discontinuous,
175 reciprocal megaBlast search with an e-value cut-off of 10^{-20} to match contigs between the in- and outgroup
176 meta-assemblies (Altschul *et al.*, 1990). We only kept the reciprocally-best hit if it was at least 100 bit
177 scores better than the next best overlapping hit. To avoid penalising good unique hits with short overlaps,
178 we allowed for 15-base overhangs between best hits. The coordinates of the retained reciprocal blast hits

179 were used to create a BED file for each species and reads for each individual overlapping these regions
180 were extracted from the BAM alignments (Quinlan & Hall, 2010). VCF files were created for the ingroup
181 and outgroup using the sub-sampled BAM files. Second, any contigs in these VCF files that matched puta-
182 tive contaminant (bacteria, fungi) or mitochondrial DNA were removed. Finally, we removed contigs with
183 excessive coverage (>75 fold coverage for *B. pallida* and > 30 fold for *B. gibbera*, see Fig. S1) as they
184 were likely to indicate remaining unfiltered collapsed repeats whose sequences had been amalgamated dur-
185 ing assembly (Nagarajan & Pop, 2013). Together, these filtering steps reduced the number of contigs by
186 80.1% corresponding to 14.1% of the original *B. pallida* meta-assembly size (Table 1). Note that we are not
187 attempting to distinguish between coding and non-coding sequence at this stage (see Sensitivity analyses).
188 Raw variants (excluding indels) were called across individuals using *samtools mpileup* (Li *et al.*, 2009).

189 **iii) Generating individual consensus sequences and triplet alignments**

190 Consensus sequences were generated for each individual from the VCF files using a custom perl script
191 (available on Dryad XXX). For the in-group; (1) the reference base was called if no variant was present or
192 the variant did not reach a particular quality threshold, (2) an 'N' was coded if an individual had 0 coverage
193 at that position or was called heterozygous by *samtools* (indicating a sequencing error because we know the
194 individual to be haploid), or more than two alleles were present, violating the assumption of the infinite sites
195 model assumed in our likelihood analysis; (3) a SNP was called for sites that differed from the reference.
196 The script used to create consensus sequences from the VCF file has the option of specifying a quality score
197 filter for SNP calling. We explored Q0, Q10, Q20, and Q30 to assess differences in SNP frequencies at
198 different quality thresholds and selected Q0, as the frequencies did not change (see Fig. S2).

199 Because the outgroup was represented by two diploid individuals, '0/1' genotypes at a position could be
200 true heterozygotes. To avoid any impact of ancestral polymorphisms between in- and outgroup, positions

201 that were variable between the outgroup individuals (0/0 versus 1/1) were coded 'N'. We generated a single
202 outgroup consensus sequence, thus taking advantage of deeper sequencing by combining data from two
203 individuals. Finally, we generated outgroup-rooted triplet alignments consisting of a single individual from
204 each of the East (*E*), Center (*Ca* or *Cb*) and Western (*Wa* or *Wb*) refugia using *muscle* (Edgar, 2004). To
205 avoid including linked sequences as separate blocks and to increase the size of blocks, non-overlapping
206 alignments that mapped to the same contig in the *B. pallida* meta-assembly were concatenated into the
207 same block. Alignments were generated for all four possible West/Center/East combinations of *B. pallida*
208 individuals and the outgroup consensus sequence (Dryad repository doi: XXX).

209 For simplicity, the following analyses focus on two such triplet sets. One, referred to hereafter as dataset
210 *a*, comprised individuals *Wa*, *Ca* and *E* (Table S1 and Fig. 1). The other, referred to hereafter as dataset *b*,
211 comprised individuals *Wb*, *Cb* and *E*.

212 **Counting mutation types**

213 Given outgroup-rooting and assuming an infinite sites mutation model with only two allelic states per site,
214 each mutation can be unambiguously placed onto a genealogical branch. This means that the polymorphism
215 information can be condensed into a vector of mutation counts on branches (Patterson *et al.*, 2006). While
216 these counts of mutation types within sequence blocks constitute the input for the likelihood analyses (see
217 section below), their relative frequencies across all sites reveal the distribution of alternative genealogies
218 across the genome and hence the types of history that are plausible. As a check, we counted the three types
219 of shared derived mutations and the three singleton mutations before and after filtering alignments to contigs
220 $> 2kb$. Note also that the current implementation of our model does not allow for back mutations; given
221 the recent timescale of divergence the chance of a backmutation within the ingroup is slim. Back mutations
222 on the outgroup branch remain a possibility, and if present would perhaps slightly the estimated divergence

223 times.

224 **Maximum likelihood analyses of historical models**

225 To conduct a broad search of model space, we took a strict divergence model between three populations
226 as a starting point and considered all histories that involve a single unidirectional admixture event either to
227 or from the oldest population. We did not include models with bidirectional or multiple admixture events
228 because the additional parameters make the corresponding generating functions computationally intractable,
229 but also because these models are biologically unexpected: expansion out of refugia is expected to be a
230 unidirectional process. For a given order of population divergence, there are six possible models (Fig. 1),
231 each with five parameters: the time of the older split (T_2); the time of the more recent split (T_1); the time
232 of admixture, or gene flow, (T_{gf}) (all measured from the present); the admixture proportion (f) and the
233 effective population size (N_e). Again, for the sake of computational tractability, we assumed a single N_e
234 for both ancestral populations as well as the population receiving migrants. We assessed the support for all
235 six admixture scenarios as well as simpler, nested models that assume no admixture and divergence between
236 either three or two populations for each of the three possible orderings of population divergence (a total of 24
237 divergence and admixture models). We also quantified the support for a basal polytomy, a single panmictic
238 population and for distinct ancestral N_e values in the strict divergence models (to test whether the additional
239 parameter substantially improved model fit without the need to invoke admixture), giving 32 models in total.

240 The general method for calculating likelihoods is described in detail elsewhere (Lohse *et al.*, 2011; Lohse
241 & Frantz, 2013). In short, the probability of observing a particular mutational configuration in a sequence
242 block (which can be interpreted as the likelihood of the model) can be expressed in terms of a higher order
243 derivative of the generating function (GF) of genealogical branch lengths (Lohse *et al.*, 2011, eq. 1). Thus,
244 given the GF for a model, it is straightforward to tabulate the logarithm of the likelihood, $\ln L$, given all

245 observed mutational configurations. Assuming that loci are unlinked and hence statistically independent,
246 the joint $\ln L$ across loci is simply the sum of $\ln L$. We used *Mathematica* v.8 (Wolfram Research, 2010)
247 to tabulate $\ln L$ values and maximise the joint likelihood numerically (Supporting *Mathematica* notebook).
248 The GF conditional on a topology of a triplet genealogy has been previously derived for a divergence model
249 with recent admixture from the population involved in the older divergence event (scenarios D and E in
250 Fig. 2) (Lohse & Frantz, 2013). We used the general recursion for the GF (Lohse *et al.*, 2011, eq. 4)
251 (and *Mathematica* to solve equations) to find analogous expressions for the other four admixture scenarios
252 depicted in figure 1. Although their derivation is relatively straightforward, the resulting expressions are
253 cumbersome and given in the Supporting *Mathematica* notebook.

254 The accuracy of the likelihood method to estimate particular model parameters can be quantified using
255 the Fisher information (I), a measure of the sharpness of the $\ln L$ curve near the maximum (Edwards, 1972;
256 Lohse & Frantz, 2013). The average information about a parameter contained in a sequence block is given
257 by summing I over all possible mutational configurations weighted by their probability ((Lohse & Frantz,
258 2013), eqn. 3/6 check?). The expected information in a data set consisting of n sequence blocks is simply n
259 $\times E[I]$.

260 **Sampling blocks**

261 For each of the six models, we numerically computed the parameter values that maximized $\ln L$ across a
262 large number of sequence blocks of fixed length.

263 Because this calculation ignores statistical associations between blocks due to linkage and we lack infor-
264 mation about the relative position of contigs in the *B. pallida* genome, the number of blocks must be chosen
265 such that the probability that two blocks are physically linked by chance can be ignored. Assuming a genome
266 size of 1.75Gb for *B. pallida* (the average measured in oak gall wasps (Lima, 2012)) and sampling of blocks

267 by chance alone, the distance between neighbouring blocks is exponentially distributed with rate $n/1.75Gb$
268 (where n is the number of blocks). This implies that if we classify blocks separated from their nearest neigh-
269 bour by 20kb or more as being in linkage equilibrium and want to ensure that less than 5% of all blocks fall
270 below this threshold, we could in theory sample a maximum of $-(1.75Gb \times \text{Log}[0.95])/20kb \approx 4500$ blocks.

271 We chose a minimum length of 2kb for the inclusion of contigs in maximum likelihood analyses as this
272 length represented a good trade-off for obtaining blocks long enough to include enough polymorphic sites
273 for inference and short enough not to worry about linkage among contigs. Sub-sampling from the full set
274 of contigs with this length cut-off gave between 2231 and 2640 blocks (depending on the combination of
275 W/C/E individuals), roughly 10% of the contigs meeting the initial filtering requirements (Table 1). To be
276 able to compare likelihoods across datasets, we fixed the number of blocks to 2231 in all analyses.

277 We initially used the first 1kb of sequence from each aligned contig in the 2kb-filtered data and explored
278 the impact of block length by repeating the analysis with shorter (500bp) and longer (2kb) blocks.

279 We estimated the proportion of coding sequence in the filtered data by Blast-searching all aligned contigs
280 against a preliminary *B. pallida* transcriptome assembly (Dryad repository doi XXX). To incorporate mu-
281 tation rate heterogeneity, sequence blocks were partitioned according their predicted proportion of coding
282 sequences into 10 equally spaced bins. We used the average divergence between *B. pallida* and *B. gibbera*
283 to calibrate a relative mutation rate for blocks in each bin.

284 **Results**

285 Below, we first examine the counts of mutation types to draw qualitative inferences about the history of *B.*
286 *pallida*. We then describe how maximum likelihood can be used to distinguish quantitatively between alter-
287 native historical scenarios. Finally, we assess the sensitivity of these inferences to the mutation model, length
288 of sequence blocks, sampling location of individuals and our assumption of no intra-locus recombination.

289 **Counting mutation types**

290 The full *a* datasets comprised 84,822 aligned contigs > 300 bases long with an N50 value of 803 bases (see
291 Table 1). We recovered a total of 171,694 polymorphic sites in the in-group, corresponding to an average
292 per site diversity (as measured by Watterson's θ_W) of 0.188 % (Table 2). Average divergence between
293 the outgroup and the Eastern individual was 4%. If population divergence were to take place in the order
294 (E,(C,W)) without admixture, we expect derived sites shared by Central and Western individuals (*C/W*) to
295 be more common than both derived sites shared by Central and Eastern individuals (*C/E*) and sites shared
296 by Western and Eastern individuals (*W/E*). Likewise *C/E* and *W/E* sites, which correspond to internal
297 branches of genealogies that are incongruent with the population history, are expected to occur at equal
298 frequency (Hudson, 1983; Tajima, 1983). Analogously, under null models of a polytomic split or a single
299 panmictic population, all three types of shared derived sites are equally likely. Contrary to these simple
300 models, we found that *C/E* sites were more frequent (9.6 %) than *W/E* sites (5.1%), which in turn were
301 more frequent than *W/C* sites (2.8 %) (see top two rows of Table 2, CHECK THESE FIGURES!). This
302 double asymmetry suggests that simple divergence models without gene flow provide a poor fit to the data.
303 If we assume that the majority class of informative sites corresponds to the order of population divergence,
304 then these results imply that the Western population diverged from the common ancestor of the Central and
305 Eastern populations before these in turn diverged. Under this model, the observed excess of *W/E* sites
306 relative to *W/C* sites could arise as a consequence of gene flow between Western and Eastern refugia (Fig.
307 1) after the more recent *C/E* split (Durand *et al.*, 2011; Lohse & Frantz, 2013).

308 **Maximum likelihood analyses of historical models**

309 Comparing the three possible histories of strict divergence, a population tree topology (W,(C,E)) had highest
310 support ($\ln L$), as expected from the frequencies of shared derived sites. Allowing for different values of N_e

311 in the two ancestral populations did improve the fit of the strict divergence model. However, 8-9 of the 18
312 models involving admixture (depending on which of the a or b datasets is considered) had greater support
313 than models of strict divergence (Table 3). The best supported history assumes a (W,(C,E)) population tree
314 topology with substantial admixture ($f = 0.76 - 0.83$ across a and b datasets) (Table 4) from the Eastern
315 into the Western refuge shortly after the split between Center and East populations (model B in Fig. 2). The
316 observed number of blocks showing each mutational configuration fits the number expected (Table S4) under
317 this best model well. Reassuringly, the alignments for the two sets of individuals *a* and *b* yielded the same
318 ranking of models and gave very similar parameter estimates with broadly overlapping 95 % C.I. (see Tables
319 4 and S5). Interestingly, however, the estimated admixture proportion *f* was slightly higher in both triplet
320 analyses involving the individual from southern rather than northern Spain (*Wb*, Fig. 1) (see Discussion).

321 Because our models are not nested, we cannot use likelihood ratio tests to test for significance. Each
322 admixture model also contains the same number of parameters, so comparisons based on Akaike information
323 criterion (AIC) reduce to comparisons based on change in log likelihood ($\Delta \ln L$). To assess our confidence in
324 our ability to identify the best model, we conducted a simulation study to quantify the power of our method.
325 Briefly, we simulated 100 replicate 2231 loci datasets using Hudson's ms program (Hudson, 2002) and the
326 ML parameter estimates for our 1kb *WaCaE* data (Table 4). Ninety-nine out of 100 replicates identified the
327 same best model as obtained for our observed data (see Table 3). The second best models were 61% model
328 A, 32% model C and 6% model E (Fig. 2, all (W,(C,E)) topology). Furthermore, the parameter estimates
329 for the intervals between T_{gf} and T_1 were tiny, such that none predicted the observed asymmetry in the
330 number of blocks specifying (C,(E,W)) or (E,(C,W)) topologies. This result underlines the fact that support
331 for model B comes from both mutation counts and configurations.

332 To provide an order of magnitude calibration for the inferred history, we applied a direct, genome-
333 wide estimate of the effective neutral mutation rate of 3.5×10^{-9} per site and generation as measured in

334 *Drosophila melanogaster* (Keightley *et al.*, 2009). To account for the bias towards conserved sequence in
335 our 2kb filtered data, we scaled the *D. melanogaster* rate by the ratio of per site diversity in the filtered and
336 unfiltered data (0.47 and 0.54 for *a* and *b* data respectively, see θ_W in Table 2). Assuming that *B. pallida* has
337 two generations per year (Csóka *et al.*, 2005; Atkinson *et al.*, 2003) this calibration gives effective population
338 sizes between 39,000 – 52,000 (Table 4). The time of admixture and the more recent split (t_{gf} , t_1) both date
339 to the last glacial period (Weichselian, 12-110kya), whereas the maximum likelihood estimate for the oldest
340 split (t_2) falls in the previous (Saalian, 130-200kya) glacial period (Table 4).

341 We also calculated $\Delta \ln L$ as parameters move away from their maximum likelihood estimates (Fig. 4).
342 and used Fisher Information to quantify how informative our data are about a particular model parameter,
343 and hence how accurate one can expect parameter estimates to be. We found that, for our best supported
344 model and the 1kb *WaCaE* dataset there is less information associated with T_2 than with the other three
345 parameters (Table S6). With 2231 loci, we expect a standard deviation (SD) of 0.143 in estimates of T_2 but
346 0.0274 in estimates of T_1 (see also Table 4).

347 **SNPs vs. blocks**

348 To assess what information, if any, is gained by using sequence blocks instead of SNPs for inference, for
349 each parameter we calculated the expected Fisher information in a single block as a function of θ (setting all
350 other parameters to their maximum likelihood estimates from the best supported model for the 1kb *WaCaE*
351 dataset). We also calculated the expected information in a single SNP for each parameter for the same model.
352 We found that even blocks containing only one SNP on average ($\theta = 0.266$) are more informative across all
353 parameters than a single SNP (Fig. S3).

354 One can use the generating function framework to obtain expected branch lengths, the sum of which gives
355 the expected total tree length. Our likelihood approach can then be used with the observed SNP frequencies

356 (Table 2, third row) to find maximum likelihood parameter values without recourse to any blocking scheme.
357 Reassuringly, when applied to the generating function for model B, very similar parameter estimates are
358 obtained, namely: $T_{gf} = 1.095$, $T_1 = 1.095$ and $T_2 = 3.34$. Note, that the interval between T_{gf} and T_1 is
359 estimated at zero.

360 **Sensitivity analyses**

361 **i) Length filtering**

362 Filtering contigs by length could result in various biases that might affect inference. For example, more
363 conserved and/or structurally complex regions of the genome are expected to assemble better and align with
364 fewer errors, and so should be represented by longer contigs. To quantify this effect, we correlated contig
365 length against per site divergence. As expected, longer contigs were on average less diverged between
366 ingroup and outgroup (Fig. S4) (Kendall's $\tau = -0.0419$, $p < 10^{-6}$). Consistent with this, the average
367 per site diversity (θ_W) in the 2kb filtered *a* data was about half of that in the unfiltered data (Table 2).
368 This confirms that length filtering does indeed enrich for conserved sequences. However, for the purpose of
369 estimating population history, any overall bias in absolute diversity can be incorporated by a simple rescaling
370 of the mutation rate (see below). In contrast, in order to justify treating the 2kb-filtered data as a random
371 sample of genealogies, we need to show that length filtering does not affect the relative frequencies of
372 mutational types (i.e. the polarized allele frequency spectrum normalized by the proportion of polymorphic
373 sites). To test this, we obtained a random sample of putatively unlinked SNPs before and after filtering the
374 *a* data for contigs > 2kb by selecting one SNP at random from each sequence block. In the length-filtered
375 data, all 2231 blocks were included. In the full data, SNPs were drawn from a random sample of 4500
376 sequence blocks to minimize linkage effects. Reassuringly, we found no significant difference between the
377 filtered and unfiltered data in the relative frequencies of the three types of shared derived mutations (the

378 most informative site types) (Table 2) ($\chi^2 = 1.96, p = 0.38$). However, there was a significant (but slight)
379 excess of singleton mutations compared to shared derived sites in the 2kb data ($\chi^2 = 9.3, p = 0.0023$). This
380 may be either due to assembly or alignment bias or purifying selection (which is likely to be stronger in the
381 2kb-filtered data) (Charlesworth *et al.*, 1993).

382 **ii) Mutational heterogeneity**

383 The likelihood method ignores mutational heterogeneity between blocks. This assumption may be problem-
384 atic given that the *B. pallida* data consists of a mix of coding and non-coding sequence. There was no hit to
385 the transcriptome of *B. pallida* for 50% of all contigs and, across all sites, the proportion of coding sequence
386 was 70%. This, together with the increased GC content in the filtered *a* data (Table 1), clearly showed that
387 our filtering strategy enriched for coding sequence. To incorporate mutational heterogeneity, we partitioned
388 blocks by their predicted proportion of coding sequence (see Methods) and scaled the effective neutral muta-
389 tion rate of each bin using the within bin divergence (per site) relative to the total divergence across all sites.
390 This drastically improved model fit (i.e. increased $\ln L$) (see Table 4), but had no impact on the ranking
391 of alternative models or parameter estimates under the best supported model. However, we did find that
392 incorporating mutational heterogeneity led to a slight reduction in both divergence time and N_e estimates
393 (see Table 4).

394 **iii) Intra-locus recombination**

395 To investigate the robustness of the maximum likelihood estimates to the assumption of no recombination
396 within blocks, we repeated the analysis with shorter (500b) and longer (2kb) blocks, both sub-sampled
397 from each contig in the 2kb-filtered data. In both cases, the relative ranking of models was unaffected and
398 parameter estimates were similar to those obtained in the initial 1kb analysis (Tables S7 and S8). This
399 suggests that undetected recombination within blocks has a minor effect on our results.

400 To investigate this further, we also conducted a simulation study to assess whether ignoring intra-locus
401 recombination biases model choice or parameter values. As starting points, we took the ML param-
402 eter estimates from the 1kb *WaCaE* dataset (Table 4) and the estimated recombination rate for *Nasonia*
403 of 1.5 cM/Mb (or crossovers/generation/bp $\times 10^{-8}$) from Niehuis (2010). Using Hudson's *ms* (Hudson,
404 2002), we simulated triplet datasets using seven recombination rates (0, 0.015, 0.15, 0.3, 0.33, 1.5, 3.3, 7.5
405 crossovers/generation/bp $\times 10^{-8}$). If the Keightley *et al.* (2009) mutation rate calibration is assumed, these
406 correspond to r/μ ratios of 0, 0.0454, 0.454, 0.907, 1, 4.5, 10 and 22.7. For each parameter combination, we
407 simulated 1,000,000 loci in order to obtain expected mutational configurations. We removed any loci that
408 failed the four gamete test (removing between 0.0338 and 3.47 % of loci), although the remaining loci will
409 still include undetected recombination events. We then parsed the polymorphism information into vectors of
410 mutation counts. We found that across all r/μ ratios, the best supported model matches that recovered from
411 the observed data. $\Delta \ln L$ values between this and the second, third and fourth-ranked models are similar to
412 that for the observed data (Fig. S6). As the r/μ ratio increases, parameter estimates become more biased: θ
413 and f decrease, while the splitting and admixture times increase (Fig. S6).

414 Discussion

415 We show how outgroup-rooted alignments of thousands of orthologous sequence blocks can be generated
416 for multiple individuals using low-coverage (< 2 fold per individual) genomic data and standard *de novo*
417 assembly tools. Although the requirement for orthologous sequences in in- and outgroup, the filtering against
418 repetitive sequences and short contigs enrich for coding and otherwise selectively constrained sequence –
419 in the case of *B. pallida* – the allele frequency spectrum is little affected. This suggests that the resulting
420 data provide a representative sample of neutral variation in the genome which, if analysed in a multi-locus
421 framework, is highly informative about recent history.

422 **Admixture dominates the history of *Biorhiza pallida***

423 The model we fit to *B. pallida* of (W,(C,E)) population divergence with strong East to West admixture differs
424 qualitatively from previous population genomic inferences of divergence with admixture (Green *et al.*, 2010;
425 Lohse & Frantz, 2013) in two ways. Firstly, admixture is from the more recently diverged population (E) into
426 the older population (W), and hence in the opposite direction to that observed in three-population analyses
427 of our own Neandertal ancestry (Green *et al.*, 2010; Durand *et al.*, 2011). Secondly, the history of *B. pallida*
428 is dominated by admixture rather than by divergence. Despite this, the majority class of shared derived
429 sites is still *C/E*, and so concordant with the order of population divergence (W,(E,C)). This is a peculiar
430 consequence of the direction of admixture: going backwards in time, *W* lineages that trace back to the *E*
431 population via admixture only spend a short time in the *E* population before they trace back to the ancestral
432 *C/E* population.

433 Both the order of population divergence and the direction of admixture are unexpected. First, our infer-
434 ence of initial divergence of the Western refuge contrasts with a previous meta-analysis of 12 oak gallwasps
435 (including *B. pallida*) and 19 associated parasitoid species (Stone *et al.*, 2012), as well as a multi-locus
436 study that compared the history of four oak gall parasitoid species (Lohse *et al.*, 2012). Both studies found
437 a general signature of (E,(C,W)) divergence on a community scale, but had insufficient power to resolve
438 the order of population divergence in individual species (or to fit additional admixture parameters). Inter-
439 estingly, however, the deep split of the Iberian population from other refugia we infer here for *B. pallida*
440 is compatible with the mitochondrial genealogy reconstructed by Rokas *et al.* (2001). Second, the history
441 of *B. pallida* involves substantial admixture from the Middle East into Iberia without affecting the Balkans.
442 One plausible route for such admixture that would not pass through the central refuge is westwards migra-
443 tion into Iberia across North Africa, possibly via southern Italy and Sicily. Striking floristic links between
444 Iberia and Asia Minor have been found across a range of plant taxa (Davis & Hedge, 1971), including oaks

445 (Lumaret *et al.*, 2002), and there is genetic evidence that Iberia was colonised from North Africa during
446 the Pleistocene by some animal taxa (Griswold & Baker, 2002; Habel *et al.*, 2008). This scenario is also
447 compatible with our finding of a higher admixture fraction for the sample from Southern Iberia (*Wb*) com-
448 pared to Central (*Wa*) Iberia, since the *Wb* sample would be closer to the putative origin of North African
449 immigrants. Similarly, the genetic similarity of extant populations of oak gallwasps (Rokas *et al.*, 2003) and
450 their parasitoids (Nicholls *et al.*, 2010b) in Morocco and Spain suggests that the Strait of Gibraltar presents
451 little or no barrier to gene flow. Given that we lack molecular calibrations for Hymenoptera in general and
452 gallwasps in particular, our absolute time estimates are tentative at best. Nevertheless, it is clear that the
453 divergence and admixture between refugial populations of *B. pallida* is recent, encompassing no more than
454 two or three glacial cycles.

455 **Sampling the genome and the limits of power**

456 While in the past, most statistical analyses of phylogeographic scenarios were limited in power by the num-
457 ber of available loci (Carstens *et al.*, 2009; Lohse *et al.*, 2012), the massive replication of sequence blocks
458 afforded by short-read sequencing overcomes this and – in the case of *B. pallida* – allowed us to reliably
459 identify the best fitting history among a set of alternative divergence and admixture scenarios.

460 However, despite increasing the number of loci by several orders of magnitude, the difference in support
461 we find for some alternative models (Table 3) is still relatively modest, suggesting that the power to dis-
462 tinguish more complex models is limited. For example, it would be hard to distinguish multiple admixture
463 events from a single event or a model of continuous migration (Hey & Nielsen, 2004). It is worth reiter-
464 ating that the lack of linkage information for the *B. pallida* assembly imposes a limitation on the number
465 of blocks we were able to include in the maximum likelihood analyses. The final analysis only included
466 2.2Mb of sequence, a mere 0.13 % of the genome, and most of the assembled genome remained unused. If

467 one had complete linkage information, i.e. if the relative position of blocks was known, one could sample
468 blocks at fixed intervals (Lohse & Frantz, 2013), which would increase the number of blocks that can safely
469 be taken as unlinked by an order of magnitude. Alternatively, one could ignore linkage between blocks
470 altogether when obtaining point estimates of parameters (which are unaffected) and use a simple scaling
471 factor to adjust confidence intervals and $\Delta \ln L$ (see Lohse & Frantz, 2013). However, the gain in power one
472 could expect is limited. In general, increasing the number of independently segregating blocks by a factor k
473 increases the accuracy of parameter estimates by $\sim \sqrt{k}$. Instead, it is the recent time-scale of the *B. pallida*
474 history that sets an inherent limit to the complexity of models among which one can hope to discriminate
475 using a multi-locus approach. If only a small number of mutations have occurred during the history of inter-
476 est (as is the case in *B. pallida* where most 1kb blocks contain two or fewer mutations), there are only a few
477 mutational configurations that are observed at appreciable frequency (Table S4).

478 Given this mutational limitation, it is clear that increasing the number of individuals sampled from within
479 each population would also only slightly improve inference: most ancestral lineages would coalesce rapidly,
480 i.e. the vast majority of genealogical branches added by larger samples would be unresolved, and so would
481 not give much extra information. Very large samples of a long non-recombining sequence can be informative
482 (Kong *et al.*, 2011), but mainly about even more recent population history than the timescale considered
483 here. Sampling individuals a further distance apart would give extra information, but also requires more
484 complex models, involving multiple parameters for separation times and admixture rates. In general, these
485 considerations suggest that there will be an upper limit to the signal contained in even an extremely large
486 number of short, unlinked sequence blocks. Nevertheless, even short blocks containing on average only one
487 SNP contain more information than single SNPs (Fig. S4).

488 In contrast, we would have far more information if we could analyse the full linear sequence and explic-
489 itly use linkage information. In *B. pallida*, a total of 3.5% of the genome would be usable after filtering for

490 unique orthologous sequence, but allowing an arbitrary degree of linkage; ultimately, of course, we could
491 use the whole genome in such an analysis. The gain does not come primarily from the sheer volume of data;
492 rather, we gain extra information from the lengths of sequence blocks. For example, the length of block
493 that shares the same genealogy within a population is inversely proportional to its coalescence time, and the
494 length of unrecombined, introgressed blocks of genome decreases with the time since introgression. Thus,
495 recombination gives an additional time-scale, beyond that provided by mutation, as used here. Barton *et al.*
496 (2013) show that in a two-dimensional continuum, the distribution of block lengths shared between genomes
497 allows inference of both dispersal rate and neighbourhood size, whereas samples of allele frequencies do not
498 give information about dispersal rate. Li & Durbin (2011) use the distribution of heterozygous SNPs to infer
499 ancestral population size through time, whilst Harris & Nielsen (2013) use this information to infer complex
500 migration histories. However, a full statistical analysis that takes into account the linear structure of the
501 genetic map not only remains extremely challenging analytically, but also requires much better assemblies
502 or linkage maps than can currently be achieved for most organisms in practice. Nevertheless, even without
503 such whole-genome data, correlations between linked loci can be informative and it will be interesting to see
504 to what extent including this information in the present maximum likelihood framework improves inference.

505 In the meanwhile, the combination of *de novo* assembly and numerical likelihood computation we de-
506 velop here provides a level of resolution far beyond that of traditional phylogeographic analyses of a few
507 loci. The fact that our bioinformatic pipeline yielded sufficient data (and resolution to distinguish between
508 models) in an oak gallwasp, the group with the largest known genomes in the Hymenoptera (Lima, 2012),
509 should encourage those working on other non-model species species and ecological communities (Stone
510 *et al.*, 2012). Furthermore, our sensitivity analyses suggest that population historical inferences based on
511 large numbers of blocks and few individuals are robust in two fundamental ways. Firstly, and despite the
512 fact that undetected recombination can bias multi-locus analyses (Strasburg & Rieseberg, 2009), neither

513 model selection nor parameter estimates are much affected by the length of sequence block. The slight bi-
514 ases observed at high r/μ ratios are in line with expectations: as recombination scrambles histories across
515 blocks of sequence the variance in branch lengths across loci is artificially decreased, leading to underesti-
516 mations of θ and overestimations of divergence times (Wall, 2003). Incidentally, the fact that small (500bp)
517 blocks with less than two mutations on average are sufficient to distinguish these models also implies that
518 the unresolved phase in diploid genomes would not be an issue when applying this framework (Lohse &
519 Frantz, 2013). The chance of multiple heterozygous sites within such short blocks and within an individual
520 is negligible. Secondly, the fact that we recover essentially the same population history using individuals
521 sampled many dispersal distances apart highlights that simple, discrete population models can be a useful
522 approximation to recent, intra specific histories.

523 **Acknowledgements**

524 We would like to thank Majide Tavakoli, Julja Ernst, Pablo Fuentes-Utrilla and Rachel Atkinson for col-
525 lecting specimens, Marian Thomson for help with the DNA extractions, the NERC Biomolecular Anal-
526 ysis Facility (NBAF) node in Edinburgh (The GenePool) for library preparation and Illumina sequenc-
527 ing and John Davey and Mark Blaxter for advice on bioinformatic strategy. This work was funded by
528 NERC grants to G Stone, J. Nicholls, K Lohse and N Barton (NE/J010499, NBAF375, NE/E014453/1 and
529 NER/B/S2003/00856). K Lohse was supported by a UK NERC fellowship (NE/I020288/1).

530 **References**

531 Altschul, S.F., Gish, W., Miller, W., Myers, E.G. & Lipman, D.J. (1990). Basic local alignment search tool.
532 *Journal of Molecular Biology*, 215(3), 403 – 410.

- 533 Arnold, B., Corbett-Detig, R.B., Hartl, D. & Bomblies, K. (2013). Radseq underestimates diversity and
534 introduces genealogical biases due to nonrandom haplotype sampling. *Molecular Ecology*, 22(11), 3179–
535 3190. ISSN 1365-294X. doi:10.1111/mec.12276.
- 536 Atkinson, R.J., Brown, G.S. & Stone, G.N. (2003). Skewed sex ratios and multiple founding in galls of the
537 oak apple gall wasp *Biorhiza pallida*. *Ecological Entomology*, 28(1), 14–24.
- 538 Barton, N.H., Etheridge, A.M., Kelleher, J. & Véber, A. (2013). Inference for the spatial lambda-Fleming-
539 Viot process. *Theoretical Population Biology*, page in press.
- 540 Barton, N.H., Kelleher, J. & Etheridge, A.M. (2010). A new model for extinction and recolonisation in two
541 dimensions: Quantifying phylogeography. *Evolution*, 64(9), 2701–2715.
- 542 Carstens, B.C., Stoute, H.N. & Reid, N.M. (2009). An information-theoretical approach to phylogeography.
543 *Molecular Ecology*, 18(20), 4270–4282.
- 544 Challis, R.J., Mutun, S., Nieves-Aldrey, J.L., Preuss, S., Rokas, A., Aebi, A., Sadeghi, E., Tavakoli, M. &
545 Stone, G.N. (2007). Longitudinal range expansion and cryptic eastern species in the western palaeartic
546 oak gallwasp *Andricus coriarius*. *Molecular Ecology*, 16(10), 2003–2014.
- 547 Charlesworth, B., Morgan, M.T. & Charlesworth, D. (1993). The effect of deleterious mutations on neutral
548 molecular variation. *Genetics*, 134(4), 1289–303.
- 549 Connord, C., Gurevitch, J. & Fady, B. (????). Large-scale longitudinal gradients of genetic diversity: a
550 meta-analysis across six phyla in the mediterranean basin. *Ecology and Evolution*, 2, 2600–2614.
- 551 Csóka, G., Stone, G.N. & Melika, G. (2005). The biology, ecology and evolution of gall-inducing Cynipidae.
552 In C. Raman, W. Schaefer & T.M. Withers, editors, *Biology, ecology and evolution of gall inducing insects*,
553 pages 573–642. Science Publisher, Enfield, New Hampshire.

- 554 Davey, J.W. & Blaxter, M.L. (2010). Radseq: next-generation population genetics. *Briefings in Functional*
555 *Genomics*, 9(5-6), 416–423. doi:10.1093/bfgp/elq031.
- 556 Davis, P.H. & Hedge, I. (1971). Floristic links between NW Africa and SW Asia. *Annales Naturhistorisches*
557 *Museum Wien*, 75(1), 43–57.
- 558 Durand, E.Y., Patterson, N., Reich, D. & M., S. (2011). Testing for ancient admixture between closely
559 related populations. *Mol Biol Ecol*, 28(8), 2239–2252.
- 560 Edgar, R. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
561 *Acids Res.*, 32(5), 1792–1797.
- 562 Edwards, A.W.F. (1972). *Likelihood*. Cambridge University Press, Cambridge.
- 563 Emerson, K.J., Merz, C.R., Catchen, J.M., Hohenlohe, P.A., Cresko, W.A., Bradshaw, W.E. & Holzapfel,
564 C.M. (2010). Resolving postglacial phylogeography using high-throughput sequencing. *Proceedings of*
565 *the National Academy of Sciences*, 107(37), 16196–16200. doi:10.1073/pnas.1006538107.
- 566 François, O., Blum, M.G.B., Jakobsson, M. & Rosenberg, N.A. (2008). Demographic history of european
567 populations of *Arabidopsis thaliana*. *PLoS Genet*, 4(5), e1000075. doi:10.1371/journal.pgen.1000075.
- 568 Green, R.E., Krause, J., Briggs, A.W., Maricic, T., Stenzel, U., Kircher, M., Patterson, N., Li, H., Zhai, W.,
569 Fritz, M.H.Y., Hansen, N.F., Durand, E.Y., Malaspinas, A.S., Jensen, J.D., Marques-Bonet, T., Alkan, C.,
570 Prufer, K., Meyer, M., Burbano, H.A., Good, J.M., Schultz, R., Aximu-Petri, A., Butthof, A., Hober, B.,
571 Hoffner, B., Siegemund, M., Weihmann, A., Nusbaum, C., Lander, E.S., Russ, C., Novod, N., Affourtit,
572 J., Egholm, M., Verna, C., Rudan, P., Brajkovic, D., Kucan, Z., Gusic, I., Doronichev, V.B., Golovanova,
573 L.V., Lalueza-Fox, C., de la Rasilla, M., Fortea, J., Rosas, A., Schmitz, R.W., Johnson, P.L.F., Eichler,
574 E.E., Falush, D., Birney, E., Mullikin, J.C., Slatkin, M., Nielsen, R., Kelso, J., Lachmann, M., Reich, D.
575 & Paabo, S. (2010). A draft sequence of the Neanderthal genome. *Science*, 328(5979), 710–722.

- 576 Griswold, C.K. & Baker, A.J. (2002). Time to the most recent common ancestor and divergence times of
577 populations of common chaffinches (*Fringilla coelebs* in europe and north africa: Insights into pleistocene
578 refugia and curent levels of migration. *Evolution*, 56(1), 143–153.
- 579 Gronau, I., Hubisz, M., Gulko, B., Danko, C. & Siepel, A. (2011). Bayesian inference of ancient human
580 demography from individual genome sequences. *Nature Genetics*, page 43.
- 581 Gutenkunst, R.N., Hernandez, R.D., Williamson, S.H. & D., B.C. (2009). Inferring the joint demo-
582 graphic history of multiple populations from multidimensional SNP frequency data. *PLoSGenetics*, 5(10),
583 e1000695.
- 584 Habel, J.C., Meyer, M., El Mousadik, A. & Schmitt, T. (2008). Africa goes europe: The complete phylo-
585 geography of the marbled white butterfly species complex *Melanargia galathea/M. lachesis* (lepidoptera:
586 Satyridae). *Organisms Diversity & Evolution*, 8(2), 121–129.
- 587 Harris, K. & Nielsen, R. (2013). Inferring demographic history from a spectrum of shared haplotype lengths.
588 *PLoS Genetics*, page *in press*.
- 589 Hayward, A. & Stone, G.N. (2006). Comparative phylogeography across two trophic levels: the oak gall
590 wasp *Andricus kollari* and its chalcid parasitoid *Megastigmus stigmatizans*. *Molecular Ecology*, 15(2),
591 479–489.
- 592 Hey, J. (2010). Isolation with migration models for more than two populations. *Molecular Biology and*
593 *Evolution*, 27, 905–920.
- 594 Hey, J. & Machado, C.A. (2003). The study of structured populations - new hope for a difficult and divided
595 science. *Nature Reviews Genetics*, 4(7), 535–543.

596 Hey, J. & Nielsen, R. (2004). Multilocus methods for estimating population sizes, migration rates and
597 divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*.
598 *Genetics*, 167(2), 747–760.

599 Hudson, R.R. (1983). Testing the constant-rate neutral allele model with protein sequence data. *Evolution*,
600 37, 203–217.

601 Hudson, R.R. (2002). Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioin-*
602 *formatics*, 18, 337–338.

603 Keightley, P.D., Trivedi, U., Thomson, M., Oliver, F., Kumar, S. & Blaxter, M.L. (2009). Analysis of the
604 genome sequences of three drosophila melanogaster spontaneous mutation accumulation lines. *Genome*
605 *Research*, 19(7), 1195–1201.

606 Kong, Q.P., Sun, C., Wang, H.W., Zhaio, M., Wang, W.Z., Zhong, L., Hao, X.D., H. Pan, S.Y.W., Cheng,
607 Y.T., Zhu, C.L., Wu, S.F., Liu, L.N., Jin, J.Q., Yao, Y.G. & Zhang, Y.P. (2011). Large-scale mtDNA
608 screening reveals a surprising matrilineal complexity in East Asia and its implications to the peopling of
609 the region. *MBE*, 28, 513–522.

610 Kulathinal, R.J., Stevison, L.S. & Noor, M.A.F. (2009). The genomics of speciation in *Drosophila*: Diversity,
611 divergence, and introgression estimated using low- coverage genome sequencing. *PLoSGenetics*, 5(7),
612 e1000550.

613 Li, H., Wysoker A, Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G. & Durbin, R. (2009). The
614 sequence alignment/map SAMformat. *Bioinformatics*, pages 2078–2079.

615 Li, H. & Durbin, R. (2011). Inference of human population history from individual whole-genome se-
616 quences. *Nature*, 475(7357), 493–6.

- 617 Lima, J. (2012). *Species Richness and Genome Size Diversity in Hymenoptera with Different Developmental*
618 *Strategies: A DNA Barcoding Enabled Study*. Ph.D. thesis, University of Guelph.
- 619 Lohse, K., Barton, N.H., Melika, N. & Stone, G.N. (2012). A likelihood-based comparison of population
620 histories in a parasitoid guild. *Molecular Ecology*, 49(3), 832–842.
- 621 Lohse, K. & Frantz, L. (2013). Maximum likelihood evidence for Neandertal admixture in Eurasian popu-
622 lations from three genomes. *MBE*, page *in review*.
- 623 Lohse, K., Harrison, R.J. & Barton, N.H. (2011). A general method for calculating likelihoods under the
624 coalescent process. *Genetics*, 58(189), 977–987.
- 625 Lohse, K., Sharanowski, B. & Stone, G.N. (2010). Quantifying the population history of the oak gall
626 parasitoid *C. fungosa*. *Evolution*, 58(4), 439–442.
- 627 Lumaret, R., Mir, C., Michaud, H. & Raynal, V. (2002). Phylogeographical variation of chloroplast DNA in
628 holm oak (*Quercus ilex* l.). *Molecular Ecology*, 11(11), 2327–2336.
- 629 McCormack, J.E., Hird, S.M., Zellmer, A.J., Carstens, B.C. & Brumfield, R.T. (2013). Applications of next-
630 generation sequencing to phylogeography and phylogenetics. *Molecular Phylogenetics and Evolution*,
631 66, 526–538.
- 632 McCormack, J.E., Maley, J.M., Hird, S.M., Derryberry, E.P., Graves, G.R. & Brumfield, R.T. (2012). Next-
633 generation sequencing reveals phylogeographic structure and a species tree for recent bird divergences.
634 *Molecular Phylogenetics and Evolution*, 62(1), 397–406.
- 635 Nagarajan, N. & Pop, M. (2013). Sequence assembly demystified. *Nature Reviews Genetics*, 14(157-167).
- 636 Nicholls, J.A., Preuss, S., Hayward, A., Melika, G., Csóka, G., Nieves-Aldrey, J.L., Askew, R.R., Tavakoli,

637 M., Schönrogge, K. & Stone, G.N. (2010a). Concordant phylogeography and cryptic speciation in two
638 western Palaearctic oak gall parasitoid species complexes. *Molecular Ecology*, 19, 592–609.

639 Nicholls, J., Fuentes-Utrilla, P., Hayward, A., Melika, G., Csoka, G., Nieves-Aldrey, J.L., Pujade-Villar, J.,
640 Tavakoli, M., Schonrogge, K. & Stone, G. (2010b). Community impacts of anthropogenic disturbance:
641 natural enemies exploit multiple routes in pursuit of invading herbivore hosts. *BMC Evolutionary Biology*,
642 10(1), 322.

643 Niehuis, Oliver Gibson, J.D.R.M.S.P.B.A.K.T.J.A.K.D.C.A.K.K.D.D.B.L.W.v.d.Z.L.S.D.M.W.J.H.G.J.
644 (2010). Recombination and its impact on the genome of the haplodiploid parasitoid wasp nasonia. *PLoS*
645 *ONE*, 5(1), e8597.

646 Patterson, N., Richter, D.J., Gnerre, S., Lander, E.S. & Reich, D. (2006). Genetic evidence for complex
647 speciation of humans and chimpanzees. *Nature*, 441(7097), 1103–1108.

648 Petit, R.J., Csaikl, U.M., Bordacs, S. & Burg, K. Coart, E.e.a. (2003). Chloroplast dna variation in european
649 white oaks phylogeography and patterns of diversity based on data from over 2600 populations. *Forest*
650 *Ecology and Management*, 176, 595–599.

651 Pickrell, J.K. & Pritchard, J.K. (2012). Inference of population splits and mixtures from genome-wide allele
652 frequency data. *PLoS Genet*, 8(11), e1002967. doi:10.1371/journal.pgen.1002967.

653 Pool, J.E., Corbett-Detig, R.B., Sugino, R.P., Stevens, K.A., Cardeno, C.M., Crepeau, M.W., Duchon, P.,
654 Emerson, J.J., Saelao, P., Begun, D.J. & Langley, C.H. (2012). Population genomics of sub-saharan
655 *Drosophila melanogaster*: African diversity and non-african admixture. *PLoS Genet*, 8(12), e1003080.
656 doi:10.1371/journal.pgen.1003080.

657 Pool, J.E., Hellmann, I., Jensen, J.D. & Nielsen, R. (2010). Population genetic inference from genomic
658 sequence variation. *Genome Research*, 20(3), 291–300. doi:10.1101/gr.079509.108.

- 659 Quinlan, A.R. & Hall, I.M. (2010). Bedtools: a flexible suite of utilities for comparing genomic features.
660 *Bioinformatics*, 26(6), 841–842. doi:10.1093/bioinformatics/btq033.
- 661 Rokas, A., Atkinson, R., Brown, G., West, S.A. & Stone, G.N. (2001). Understanding patterns of genetic
662 diversity in the oak gallwasp *Biorhiza pallida*: demographic history or a *Wolbachia* selective sweep?
663 *Heredity*, 87, 294–305.
- 664 Rokas, A., Atkinson, R.J., Webster, L., Csóka, G. & Stone, G.N. (2003). Out of Anatolia: longitudinal
665 gradients in genetic diversity support an eastern origin for a circum-mediterranean oak gallwasp *Andricus*
666 *quercustozae*. *Molecular Ecology*, 12(8), 2153–2174.
- 667 Sousa, V. & Hey, J. (2013). Understanding the origin of species with genome-scale data: modelling gene
668 flow. *Nat Rev Genet*, 1(6), 404–414.
- 669 Stone, G.N., Challis, R.J., Atkinson, R.J., Csóka, G., Hayward, A., Melika, G., Mutun, S., Preuss, S., Rokas,
670 A., Sadeghi, E. & Schönrogge, K. (2007). The phylogeographical clade trade: tracing the impact of
671 human-mediated dispersal on the colonization of northern Europe by the oak gallwasp *Andricus kollari*.
672 *Molecular Ecology*, 16, 2768–2781.
- 673 Stone, G.N., Hernandez-Lopez, A., Nicholls, J.A., di Pierro, E., Pujade-Villar, J., Melika, G., Cook, J.M. &
674 Abbot, P. (2009). Extreme host plant conservatism during at least 20 million years of host plant pursuit
675 by oak gallwasps. *Evolution*, 63(4), 854–869.
- 676 Stone, G.N., Lohse, K., Nicholls, J.A., Fuentes-Utrilla, P., Sinclair, F., Schönrogge, K., Csóka, G., Melika,
677 G., Nieves-Aldrey, J.L., Pujade-Villar, J., Tavakoli, M., Askew, R.R. & Hickerson, M.J. (2012). Recon-
678 structing community assembly in time and space reveals enemy escape in a western palaeartic insect
679 community. *Current Biology*, 22(6), 531–537.

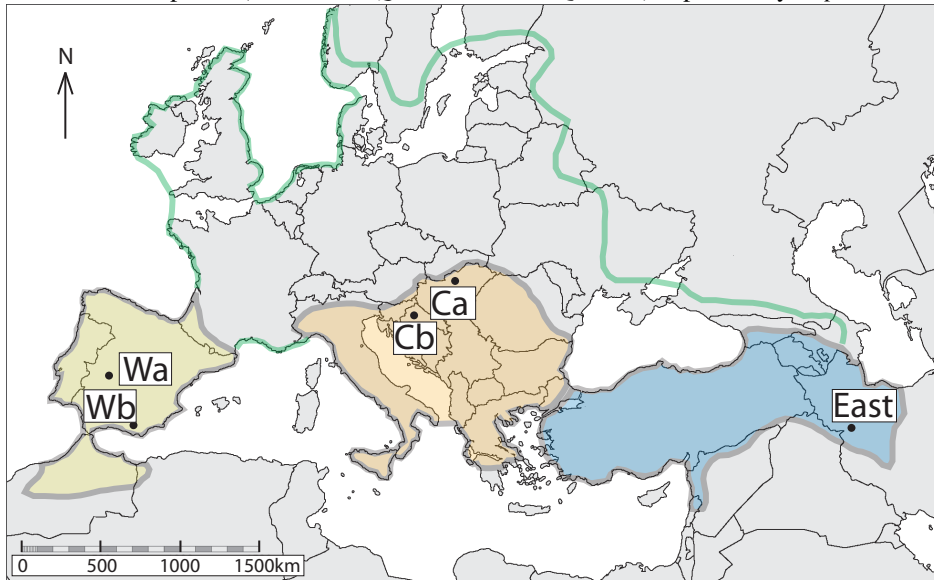
- 680 Stone, G.N. & Sunnucks, P. (1993). Genetic consequences of an invasion through a patchy environment - the
681 cynipid gallwasp *Andricus quercuscalicis* (Hymenoptera: Cynipidae). *Molecular Ecology*, 2(4), 251–268.
- 682 Stone, G., Nee, S. & Felsenstein, J. (2011). Controlling for non-independence in comparative analysis of
683 patterns across populations within species. *Phil. Trans. Roy. Soc. Series B Biol Sci.*, 366, 1410–1424.
- 684 Strasburg, J.L. & Rieseberg, L.H. (2009). How robust are isolation with migration analyses to violations of
685 the IM model? A simulation study. *Molecular Biology and Evolution*, 27(2), 297–310.
- 686 Tajima, F. (1983). Evolutionary relationships of DNA sequences in finite populations. *Genetics*, 105(2),
687 437–460.
- 688 Wall, J.D. (2003). Estimating ancestral population sizes and divergence times. *Genetics*, 163(1), 395–404.
- 689 Wang, Y. & Hey, J. (2010). Estimating divergence parameters with small samples from a large number of
690 loci. *Genetics*, 184, 363–373.
- 691 Wolfram Research, I. (2010). *Mathematica, Version 8.0*. Wolfram Research, Inc., Champaign, Illinois.

692 **Data Accessibility**

- 693 • Raw Illumina reads for *B. pallida* and *B. gibbera*: ENA Sequence Read Archive (<http://www.ebi.ac.uk/ena/data/view/ERP002>).
- 694 • Meta-assemblies for *B. pallida* and *B. gibbera*: Dryad doi: XXX
- 695 • *B. pallida* transcriptome: Dryad doi: XXX
- 696 • Final triplet alignments and summary files for Mathematica: Dryad doi: XXX.

697 **Figures**

Figure 1: Sampling locations of the five *B. pallida* individuals (Wa, Wb, Ca, Cb and E, see Table S1) used for genome sequencing and population genomic analyses. Refugial regions are colour coded as follows: W, West, in green; C, Centre, in orange; and E, East, in blue. The green line shows the current postglacial extent of the oak host plants (white oaks, *Quercus* section *Quercus*) exploited by *B. pallida*.



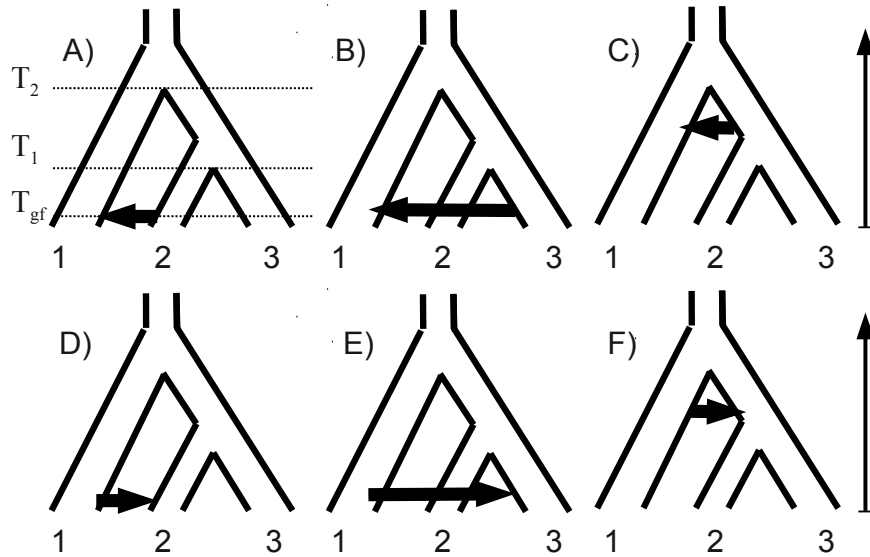


Figure 2: The six possible models for divergence histories involving unidirectional admixture to or from the older population (labelled: 1). Given the three possible orders of population divergence, there are 18 admixture models in total. Divergence times (T_1 and T_2) and the time of admixture (T_{gf}) are measured back in time from the present.

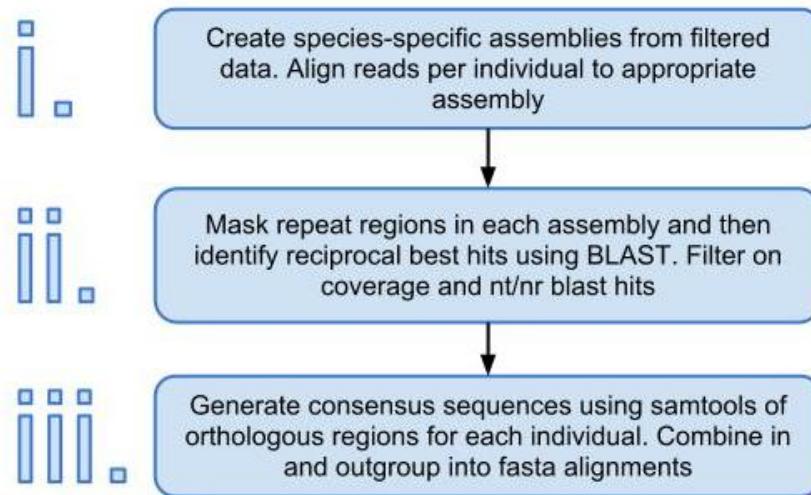


Figure 3: Assembly and filtering steps used to generate population genomic datasets in *B. pallida*.

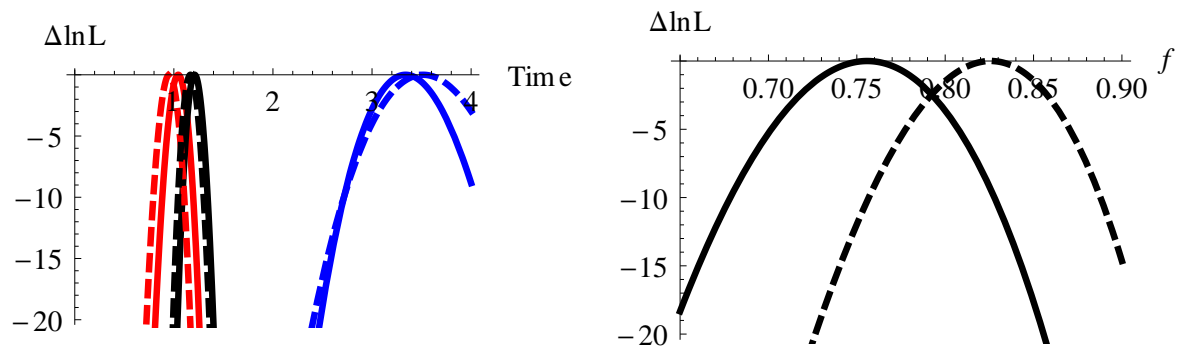


Figure 4: A) $\Delta \ln L$ plots for the times of divergence (T_1 (black) and T_2 (blue)) and admixture T_{gf} (red). B) $\Delta \ln L$ for the admixture proportion f . Estimates from the *a* data are shown as solid lines, those from the replicate data set *b* as dashed lines.

Table 1: Summary statistics for meta-assemblies and filtered datasets

Species	N50	# of contigs	Total bases	GC %
<i>Biorhiza pallida</i>	1,075	1,163,314	805,102,378	32.9
<i>Belizinaella gibbera</i>	643	817,710	443,963,639	36.1
<i>Biorhiza pallida</i> , post Filter	734	232,097	113,583,710	36.7
<i>Belizinaella gibbera</i> , post Filter	508	290,379	111,785,775	35.9
<i>a</i> data >300 bp	803	84,822	61,012,720	38.6
<i>b</i> data >300 bp	768	77,752	54,117,641	39.0
<i>a</i> data 2 kb	2,000	2,640	5,280,000	40.1
<i>b</i> data 2 kb	2,000	2,231	4,462,000	40.2

Summaries are shown for in- and outgroup meta-assemblies (first four rows) and the *a* and *b* triplet data before and after length filtering. N50 is defined as the length N for which 50% of all sequenced bases are assembled in a contig of length $< N$.

Table 2: Genetic diversity and relative frequencies of mutational types in *B. pallida* alignments.

Dataset	θ_w	<i>W</i>	<i>C</i>	<i>E</i>	<i>W/C</i>	<i>W/E</i>	<i>C/E</i>
<i>a</i> data > 300 bp	0.00188	0.325	0.214	0.263	0.040	0.058	0.100
<i>b</i> data > 300 bp	0.00147	0.269	0.244	0.283	0.044	0.060	0.100
<i>a</i> data 2 kb	0.00089	0.338	0.220	0.267	0.027	0.049	0.098
<i>b</i> data 2 kb	0.00079	0.276	0.250	0.287	0.035	0.054	0.099

Before (>300bp) and after (>2kb) length filtering.

Table 3: Support for alternative scenarios of divergence and admixture in the oak gall wasp *B. pallida* (WaCaE, 1kb data)

Model	k			
Panmixia	1			-589.3
Polytomy	2			-88.7
gene flow		$(W_1, (C_2, E_3))$	$(C_1, (E_2, W_3))$	$(E_1, (C_2, W_3))$
A)	2 → 1	5	-9.1, (T_1)	-18.8
B)	3 → 1	5	0	-88.7, (T_1, T_2)
C)	2/3 → 1	5	-4.8	-88.7, (T_1, f^*)
D)	1 → 2	5	-25.7, (f)	-88.7, (T_{gf}, f^*)
E)	1 → 3	5	-18.0	-18.2, (T_1)
F)	1 → 2/3	5	-25.7 (f^*)	-18.2, (f^*)
	2 pop.	2	-260.8	-88.7, (T_1, T_2)
	3 pop.	2	-25.7	-88.7, (T_2)
	2 pop. N_e	3	-48.5	-93.7
	3 pop. N_e	4	-20.8	-90.1
				-88.7, (T_2)

Support ($\Delta \ln L$) relative to the best model for alternative histories of refugial populations of *B. pallida* estimated from the *a* dataset (Model B in Fig. 2 has highest support and is shown in bold). The labelling of populations (1–3) and of models (A–F) corresponds to that in Fig. 2; all scenarios involving unidirectional admixture were assessed for each of the three possible orders of population divergence (columns 1–3). Models of strict divergence without admixture between two (2 pop., i.e. $T_1 = 0$) or three (3 pop.) populations were fitted assuming either a single or two different N_e for ancestral populations. Parameters for which the maximum likelihood estimate is 0 (i.e. the model reduces to a simpler nested model) are indicated in brackets (f^* refers to complete admixture, i.e. $f = 1$).

Table 4: Parameter estimates under the best supported model (see Table 3).

dataset	μ het.	$\ln L$	f	$\theta (N_e)$	$T_{gf} (t_{gf})$	$T_1 (t_1)$	$T_2 (t_2)$
<i>a</i> , 1kb	no	-9269.3	0.76 (0.72, 0.79)	0.69 (52,000)	1.04 (54KY) (51–58KY)	1.21 (63KY) (60 – 66 KY)	3.34 (173KY) (158 – 189KY)
<i>b</i> , 1kb	no	-8815.1	0.83 (0.80, 0.86)	0.64 (43,000)	0.95 (41KY) (38–44KY)	1.17 (50KY) (51 – 57 KY)	3.51 (151KY) (135 – 168KY)
<i>a</i> , 1kb	yes	-8769.7	0.76 (0.72, 0.79)	0.61 (45,900)	1.10 (50KY) (47–54KY)	1.26 (58KY) (55 – 60 KY)	3.45 (158KY) (143 – 172KY)
<i>b</i> , 1kb	yes	-8444.0	0.82 (0.79, 0.85)	0.58 (39,100)	0.97 (38KY) (35–40KY)	1.17 (51KY) (49 – 54 KY)	3.47 (136KY) (121 – 151KY)

Maximum likelihood estimates are given for different triplet combinations and analyses with and without mutational heterogeneity (μ het.; see Methods). Both effective population size and divergence time parameters are scaled relative to the rate of coalescence, i.e. in $2N_e$ generations. Absolute values are given in brackets, calibrated using a direct, genome-wide mutation rate for *Drosophila* (Keightley *et al.*, 2009) and assuming two generations per year. 95 % confidence intervals of scaled parameter values are given in brackets below the point estimate. f is the admixture proportion, θ is the scaled mutation rate, N_e is the effective population size, T_{gf} is the time ago of admixture (t_{gf} is the calibrated estimate), T_1 and T_2 are the younger and older splitting times in the population topology, with t_1 and t_2 the absolute ages, respectively.