

1 Going down the rabbit hole: a review on methods characterizing selection and
2 demography in natural populations

3 **Running title:** Linking evolution and genome-wide data

4 Yann X.C. Bourgeois¹, Khaled M. Hazzouri¹, Ben H. Warren²

5 ¹New York University Abu Dhabi, PO Box 129188, Saadiyat Island, Abu Dhabi, United
6 Arab Emirates

7 ²Department of Systematic and Evolutionary Botany, University of Zurich, Zollikerstrasse 107, 8008
8 Zurich, Switzerland

9

10 **Abstract**

- 11 1. Characterizing species history and identifying loci underlying local adaptation is
12 crucial in functional ecology, evolutionary biology, conservation and agronomy. The
13 ongoing and constant improvement of next-generation sequencing (NGS) techniques
14 has facilitated the production of an ever-increasing number of genetic markers across
15 genomes of non-model species.
- 16 2. The study of variation in these markers across natural populations has deepened the
17 understanding of how population history and selection act on genomes. Population
18 genomics now provides tools to better integrate selection into a historical framework,
19 and take into account selection when reconstructing demographic history. However,
20 this improvement has come with a burst of analytical tools that can confuse users.
- 21 3. Such confusion can limit the amount of information effectively retrieved from
22 complex genomic datasets. In addition, the lack of a unified analytical pipeline impairs
23 the diffusion of the most recent analytical tools into fields like conservation biology.
- 24 4. To address this need, we describe possible analytical protocols and link these with
25 more than 70 methods dealing with genome-scale datasets. We summarise the

26 strategies they use to infer demographic history and selection, and discuss some of
27 their limitations. A website listing these methods is available at
28 www.methodspopgen.com.

29 **Keywords**

30 **Coalescent, Software, Molecular evolution**

31 **Introduction**

32 Multiple historical and selective factors shape the genetic makeup of populations. The advent
33 of Next-Generation Sequencing (NGS) in the last 10 years has enhanced our understanding on
34 how intermingled these factors are, and how they can impact genomic variation. Important
35 results have been gathered on model species, or species of economic interest. Such results
36 include, among other examples, an improved understanding of the history of human
37 migrations, admixture and adaptation (e.g. Sabeti *et al.*, 2002; Abi-Rached *et al.*, 2011; Li and
38 Durbin, 2011), the origin of domesticated species (e.g. Axelsson *et al.*, 2013; Schubert *et al.*,
39 2014), and the genetic basis of local adaptation in both model and non-model species (e.g.
40 Legrand *et al.*, 2009; Kolaczkowski *et al.*, 2011; Roux *et al.*, 2013; Kubota *et al.*, 2015). The
41 amount of population genomic data that is aimed at elucidating the history of natural
42 populations has increased enormously in the last five years, even for non-model species.
43 Studying genetic variation at the genome level allows the demographic factors shaping
44 species history to be characterised. Further, understanding demographic history is important
45 in correctly identifying loci under selection. Such data can even help in conservation efforts
46 by identifying locally adapted genes that can be used to define relevant conservation units
47 (Fraser and Bernatchez, 2001).

48

49 In the last 10 years, developments in NGS have continually improved the throughput of data,
50 while reducing time and cost of their production. These methods have become more
51 affordable for teams studying evolutionary processes in biology, and many new methods to
52 infer demography and selection have been developed. However, these methodological
53 advances have brought increased analytical complexity to the field, and an inflation in the
54 number of methods covering any one topic. As a consequence, it has become increasingly
55 difficult for all potential users to follow developments and be sure of selecting the most
56 appropriate method for the question and data in hand.

57

58 An overarching theme that concerns new users in a wide range of contexts is understanding
59 patterns of heterogenous diversity along the genome. Patterns of nucleotide variation in
60 genomes are shaped by both intrinsic and extrinsic factors. Even within a single isolated
61 panmictic population, interaction between recombination, selection and historical variation in
62 population size will lead to heterogeneous diversity along the genome. At the scale of several
63 connected populations or even between emerging species, these processes will affect the rate
64 at which migration homogenizes the genome (Wolf and Ellegren, 2016).

65

66 A prime example is the situation of a researcher primarily interested in identifying signatures
67 of recent positive selection in a species of interest. Since a new mutation will see its frequency
68 increase in a population where it provides a selective advantage (*i.e.* hard selective sweep), a
69 large region around it can remain uniform, especially if selection is strong (Sabeti *et al.*, 2002;
70 McVean, 2007; Vitti *et al.*, 2013). This can lead to an increase in linkage disequilibrium (LD)
71 between variants associated to the advantageous mutation, as well as a decrease in the age of
72 the positively selected alleles and their nucleotide diversity. If positive selection occurs only

73 in some populations, it may be possible to observe an increase in differentiation at this locus
74 (Charlesworth *et al.*, 1997). To detect this signature of selection, some methods can track
75 particularly long haplotypes and linkage disequilibrium along the genome. Others will rather
76 focus on allele frequency spectrum and nucleotide diversity. Association methods will take
77 advantage of preliminary knowledge of a phenotype or environment to identify loci displaying
78 correlated allele frequencies. A few methods aim at inferring the whole history of coalescence
79 and recombination along genomes, but still make simplifying assumptions and often require
80 whole-genome resequencing data, which remain unaffordable for many teams.

81

82 Therefore, the choice of methods of any such researcher will depend on the available data and
83 specifics of the question being addressed. One key aspect is that all these methods and
84 questions do not have the same requirements in terms of reference genomes and marker
85 density. For example, recent discussion of RAD-markers has been interesting from this
86 perspective (Lowry *et al.*, 2016; Catchen *et al.*, 2017). The density of markers obtained along
87 a genome depends on the choice of the restriction enzyme, and this choice must take into
88 account the average extent of LD. Genome scans of selection will lose power if this density is
89 not enough to cover mutations in strong linkage with variants under selection.

90

91 In the absence of any unified framework, combining several tools is necessary to interpret
92 results. It must be borne in mind that recombination rates vary along the genome, which can
93 possibly bias tests based on LD. It can therefore be important to characterize the
94 recombination landscape in natural populations, requiring the use of another method (e.g.
95 LDHat, Table 1). Background selection can lead to signatures of high differentiation that
96 mimic disruptive selection (Charlesworth *et al.*, 1997). An assessment of genetic diversity

97 within populations, haplotype frequencies and possibly association with phenotype in each
98 population would therefore be needed to explore this possibility (Charlesworth *et al.*, 1997).
99 Demographic history impacts patterns of LD, allele age and frequencies at the genome scale,
100 and affects the efficiency of selection at specific genes. This calls for at least basic checking
101 of demographic structure and history and ideally building neutral demographic models to
102 estimate the expected frequency of outliers without involving selection. In addition, most
103 methods estimating selection coefficients require estimating effective population sizes.
104 Finally, including markers under selection can bias demographic inference by skewing allele
105 frequency spectra and LD, which requires careful data filtering and removal of outliers.

106

107 In this simplified example, we see that a reciprocal feedback between different aspects of
108 evolutionary genomics is needed (Figure 1). Combining approaches is one of the current
109 grand challenges in evolutionary biology (Cushman, 2014). While large-scale collaborations
110 and sharing of skills between researchers allow for detailed analyses, a regularly updated list
111 of methods would be valuable for smaller research teams to quickly start new projects and
112 evaluate their experimental design.

113

114 In addition to methodological and technical challenges, the widespread use of sophisticated
115 analytical tools is made difficult by the lack of communication between fields (Shafer *et al.*,
116 2015), little user-friendliness of software, inflation of data formats (Lischer and Excoffier,
117 2012) and the ever-increasing number of methods made available. Fields like landscape
118 genetics and phylogeography have largely focussed on identifying general patterns in
119 populations history and species diversification. Other researchers are more interested in
120 identifying specific genes that are involved in adaptation in natural populations. All these

121 views contribute to our understanding of causation in biology, an effort that has included
122 genetics, developmental science and ecology (Laland *et al.*, 2011). A global summary of
123 methods used in these different fields would therefore facilitate communication between
124 disciplines.

125

126 The last extensive review of methods in population genetics was performed 10 years ago
127 (Excoffier and Heckel, 2006). Since then there has been increasing drive to translate these
128 methods into approaches applicable to genomic data and non-model species. This drive has
129 confirmed the value of population genomics on non-model species in understanding
130 biological diversity at various scales (Mandoli and Olmstead, 2000; Jenner and Wills, 2007;
131 Abzhanov *et al.*, 2008; White *et al.*, 2010; Ellegren *et al.*, 2012; Weber *et al.*, 2013; Poelstra
132 *et al.*, 2014). Such advances are needed to broaden our view about the evolutionary process
133 and improve sampling of distant clades. Ultimately, this process should provide a more
134 balanced picture than the one brought by the study of a few model species (Abzhanov *et al.*,
135 2008). Genomic approaches also have the potential to improve conservation genetic inference
136 by scaling up the amount of data available (Shafer *et al.*, 2015). Much effort has recently been
137 made in facilitating the diffusion of sometimes complex, state-of-the-art methods. Their
138 application to species with little background data has become more accessible, bringing the
139 potential to add much valuable information.

140

141 In this paper, we propose possible pipelines (Figures 1, 2 and 3) to help choose appropriate
142 methods dealing with current questions in population genomics and genetics of adaptation in
143 natural populations. We begin with a succinct review of methods available to obtain genome-
144 wide polymorphism data (Box 1) before focusing on i) methods devoted to the study of

145 population structure and quantitative characterization of population history (Table 1 and 2)
146 and ii) methods aimed at identifying selected loci (Table 3). We end this review by detailing
147 how these analyses can be combined, and present future directions that may be taken by the
148 field of population genomics. The tables and a summary of the methods discussed in this
149 paper will be kept updated to follow improvements, and are available at
150 www.methodspopgen.com.

151

152 Box 1. Common sequencing methods

153 *RAD-seq*: Reduced representation allows broad sampling of variants across the genome by
154 sequencing DNA fragments flanking restriction sites. Such sampling is not specific to any
155 particular kind of region (e.g. coding or non-coding). Some of the best-known reduced
156 representation techniques include RAD-sequencing (Baird *et al.*, 2008) and Genotyping by
157 Sequencing (GBS; Elshire *et al.*, 2011). Their main interest is their low cost and that they do
158 not require any reference genome (see Davey *et al.*, 2011 for details), although a reference can
159 be useful to identify outlier genomic regions and retrieve linkage disequilibrium information
160 between markers. Use of a reference genome also limits the bias due to paralogy and mapping
161 errors (Hand *et al.*, 2015). Reduced representation allows many individuals to be genotyped at
162 once, and so is widely used for the study of population structure, demography and selection. It
163 does not cover all mutations in the genome and the choice of the restriction enzyme is crucial
164 to control for the density of markers. This choice further controls the mean sequencing depth,
165 the number of mutations close to genes under selection, and the accurate calling of genotypes.
166 The number of SNPs ranges from thousands to millions, which is usually enough to retrieve
167 substantial information about demography and sometimes selection (see Puritz *et al.*, 2014 for
168 a detailed summary of reduced-representation techniques). As a general word of caution, note
169 that RAD-sequencing and related methods display specific properties that can bias genome-

170 wide estimates of diversity, e.g. allelic dropout (Arnold *et al.*, 2013, Puritz *et al.* 2014).
171 However, this type of marker remains valuable for phylogenetic estimation, even for distantly
172 related species (Cariou *et al.*, 2013), and allelic dropout can be compensated for by focusing
173 only on markers sequenced in all individuals. Variations on the original RADseq protocol
174 have been developed to overcome some of these caveats (ddRAD, Peterson *et al.*, 2012;
175 ezRAD, Toonen *et al.*, 2013; 2b-RAD, Wang *et al.*, 2012). Many pipelines have been
176 specifically designed to account for RAD-seq specificities, including Stacks (Catchen *et al.*,
177 2011), TASSEL-UNEAK (Lu *et al.*, 2013) or TASSEL-GBS for GBS data (Glaubitz *et al.*,
178 2014).

179

180 *Targeted sequencing:* This class of methods allows sequencing and genotyping the same set
181 of genomic fragments or single nucleotide polymorphisms (SNP arrays) across individuals,
182 and has been recently promoted to study non-model species (Jones and Good, 2016). Since
183 the specificity of the probe does not have to be very high, the same probe can be used among
184 closely related species (Nicholls *et al.*, 2015). Conservation of the target genomic region
185 under study is important. High conservation may lead to higher efficiency of capture but can
186 artificially reduce representation of polymorphic regions. Different technologies allow for
187 targeted sequence capture that can be classified by enrichment methods (hybridization-based;
188 PCR-based; molecular inversion probe-based; see Mamanova *et al.*, 2010). Commercial
189 products, such as Agilent's SureSelect, MYcroarray's MYbaits or Roche NimbleGen's
190 SeqCap offer these methods or a derivation (Grover *et al.*, 2012).

191 Targeted sequencing reduces the genomic representation compared to whole genome
192 sequencing and it allows for multiple individuals to be multiplexed, lowering the cost of
193 sequencing per sample. In addition, the complexity of analysis is reduced compared to whole
194 genome sequencing (WGS), since only a subset of genomic regions is sequenced. By allowing

195 an improvement in spatial and temporal sampling, targeted sequencing can reconstruct
196 dispersal routes and migration between varieties and subspecies (Nadeau *et al.*, 2012; da
197 Fonseca *et al.*, 2016). Another commonly used technique includes Single nucleotide
198 polymorphism (SNP) genotyping arrays have frequently been used in studies aimed at
199 detecting phenotype/genotype associations or to study population structure (Gautier *et al.*,
200 2010; Johnston *et al.*, 2011). However, re-genotyping of ascertained SNPs in a new population
201 can lead to bias which can be problematic for demographic inference (Albrechtsen *et al.*,
202 2010; Lachance and Tishkoff, 2013).

203

204 *RNAseq*: RNAseq can be used with and without a reference genome. In the latter case, like
205 any other reduced representation method, it does not provide information of linkage among
206 genes. It has applications on many different evolutionary time scales. Since it mostly
207 sequences coding regions, a deep phylogeny can be constructed with conserved orthologs.
208 Depth of coverage is gene expression dependent, so calling genotypes varies across genes and
209 which must be taken into consideration (Gayral *et al.*, 2013). If a reference genome is
210 available, it is possible to call variants (Piskol *et al.*, 2013). This method is cost-effective and
211 an alternative to whole genome sequencing. However, common variant callers do not behave
212 well with RNAseq due to reads encompassing intronic regions as well as bias introduced
213 during the sequencing library preparation. One of the common variant calling pipelines
214 available is GATK which suggests best practices for calling variants on RNAseq
215 (<https://software.broadinstitute.org/gatk/best-practices/>). Another variant calling protocol
216 specifically designed for RNAseq is Opossum (Oikkonen and Lise, 2017), which can be used
217 with haplotype-based callers such as Platypus and GATK haplotypeCaller. This software
218 maintains precision and improves the sensitivity of SNP calling compared to the GATK best
219 practice pipeline. RVboost (Wang, Davila, *et al.*, 2014) was developed using the method of

220 variant prioritization, using a so-called boosting method that uses a set of high-confidence
221 variants to set a model of good quality variants. All RNA variants are then prioritized and
222 called based on this model. It outperforms Variant Quality Score Recalibration (VQSR) from
223 the Genome Analysis Tool Kit (GATK) and the RNA-Seq variant calling pipeline SNPiR
224 (Piskol *et al.*, 2013). RVboost can identify false variants introduced by random hexamer
225 priming during library preparation.

226

227 *Whole genome resequencing:* Whole-genome resequencing requires a well assembled
228 reference and is more expensive than RAD-seq or targeted sequencing, especially for species
229 with long and complex genomes. Some methods do not actually require any reference
230 sequence to call SNPs from raw reads, like kSNP2 (Gardner and Hall, 2013) or DiscoSNP
231 (Uricaru *et al.*, 2015). However, this limits the main interest of this approach, since mapping
232 back on a reference has the potential to provide a complete overview of structural and coding
233 variation. It also allows the use of powerful methods to track signatures of selection (see
234 below). Pooled sequencing (Futschik and Schlötterer, 2010) can be an option to reduce costs,
235 but generally restricts analyses to methods focusing on allele frequencies. Since individual
236 information is not available, variation in Linkage Disequilibrium across individuals (LD)
237 cannot be exploited. Shallow sequencing (1-5X per individual) may be a way to partly
238 overpass this last issue for a similar cost (Buerkle and Gompert, 2013), but should not be used
239 for methods requiring phasing and unbiased individual genotypes.

240 Shallow shotgun sequencing also allows retrieving complete plastomes, due to the
241 representation bias of mitochondrial or chloroplast sequences. Plastome sequences can
242 provide insightful information into the evolutionary history of populations or species, and
243 recent work has successfully used shallow sequencing to reconstruct mitochondrial or
244 chloroplast sequences in plants (Malé *et al.*, 2014), animals (Hahn *et al.*, 2013) or old and

245 altered museum samples (Besnard *et al.*, 2016). Methods such as MITObim (Hahn *et al.*,
246 2013) provide an automated and relatively user-friendly way to reconstitute plastome
247 sequences, which can then be analyzed as a single non-recombining marker for phylogeny or
248 population genetics.

249

250 Population structure and data description

251 **Population structure and diversity**

252 Description of the data is essential to assess the proportion of loci displaying a consistent
253 pattern, and characterize how genetic diversity is partitioned within species. Genetic diversity
254 and its genome-wide variance are directly impacted by variation in many factors including
255 effective population sizes, population structure, inbreeding, migration, and recombination
256 rates. Their characterization must be performed prior to any analysis to get insights into the
257 forces and constraints acting on populations.

258

259 A key aspect when describing a new dataset is the assessment of relatedness between
260 individuals or localities. Neglecting population structure can dramatically bias demographic
261 inference, especially when gene flow is not accounted for or panmixia is assumed (Chikhi *et*
262 *al.*, 2010; Heller *et al.*, 2013). It also biases the detection of loci under selection (e.g. Nielsen
263 *et al.*, 2007). Cryptic population structure is typically a confounding effect in studies of
264 phenotype-genotype association studies, when a given feature or trait is disproportionately
265 found in a population or a set of related individuals (Balding, 2006). Fortunately, the
266 abundance of SNP data produced by typical genomic studies is often enough to thoroughly
267 assess relatedness between individuals.

268

269 Many tools currently exist to infer population structure (Table 1, Figure 2). An elegant and
270 efficient class of methods relies on using multivariate approaches such as principal component
271 analysis (PCA) to infer relatedness between individuals and populations without *a priori*
272 knowledge. Since these methods do not have underlying assumptions based on population
273 genetics, they are suitable for analyzing species displaying polyploidy or mixed-ploidy
274 (Dufresne *et al.*, 2014). A detailed review of these methods has been already performed
275 (Jombart *et al.*, 2009) and an exhaustive list of their applications is beyond the scope of this
276 review. These approaches have been especially useful to study the consistency between
277 geographical and genetic structure in human populations of Europe (Novembre *et al.*, 2008).
278 They were also recently applied to RAD-sequenced populations of a freshwater crustacean
279 (*Daphnia magna*). Procrustes rotation (Novembre *et al.*, 2008) was used to match
280 geographical coordinates with PCA axes, showing how isolation by distance has shaped
281 genetic structure (Fields *et al.* 2015).

282

283 Methods for estimating the relatedness of individuals are suited to studies relying on pedigree
284 information, or if there are reasons to suspect that familial relationships can play a major role
285 in shaping genetic structure of the population(s) considered. When each individual in a study
286 is sampled from a different location or environment, estimating relatedness also provides a
287 way to assess the genetic distance between individuals. Genetic distance can then be
288 compared with geographical or ecological distance. For example, in a recent study using more
289 than 1000 *Arabidopsis thaliana* genomes, estimates of relatedness have allowed the
290 identification of putatively relictual populations that may have persisted in Europe since the
291 last Ice Age (Alonso-Blanco *et al.*, 2016).

292

293 Approaches such as Structure (Pritchard *et al.*, 2000) and fastSTRUCTURE (Raj *et al.*, 2014)
294 have been widely used to determine hierarchical population structure and admixed
295 populations by grouping individuals in clusters. The optimal number of clusters (K) can then
296 be determined based on likelihood, although examining population structure for a range of K
297 can allow substructure to be better identified. The main interest of these approaches is that
298 they provide a measure of coancestry coefficients, which are the proportions of an individual
299 genome originating from multiple ancestral gene pools. Such information is more difficult to
300 retrieve with approaches such as PCA. There have been criticisms however about whether
301 ambiguous assignment could be actually interpreted as a signal of admixture, and detailed
302 inference requires thorough model testing and estimating the goodness of fit of a model with
303 admixture (see Falush *et al.*, 2016).

304

305 **Heterogenous patterns of divergence between species along their genomes**

306 Advantageous alleles can migrate from one population to another, resist introgression from
307 other populations, reach fixation and erase diversity around them. This is one scenario leading
308 to heterogenous patterns of divergence along the genome, the so-called islands of divergence
309 (Wolf and Ellegren, 2016). Alternative scenarios leading to similar patterns were recently
310 highlighted (Cruickshank and Hahn, 2014). Understanding the origin of genomic regions
311 under selection highlights the evolutionary history of adaptive alleles (e.g. Abi-Rached *et al.*,
312 2011) and contributes to our understanding of the origin and maintenance of reproductive
313 isolation. Studies focusing on hybrid zones and introgression have provided inspiring
314 examples (Hedrick, 2013), as demonstrated by recent work focusing on patterns of
315 heterogenous gene flow in *Mytilus* mussels (Roux *et al.*, 2014), localized introgression and

316 inversions at a color locus in *Heliconius* butterflies (The Heliconius Genome Consortium *et*
317 *al.*, 2012) and adaptive introgression of anticoagulant resistance alleles in mice (Song *et al.*,
318 2011). Descriptive statistics computed along genomes provide valuable information in this
319 context. One may for example plot the distribution of a differentiation measure such as F_{ST}
320 (Weir and Cockerham, 1984) between populations, mean linkage disequilibrium or nucleotide
321 diversity. Such an approach has been used in *Ficedula* flycatchers, which uncovered clear
322 genomic islands of divergence and the higher differentiation on sexual chromosomes due to
323 ongoing reproductive isolation (Ellegren *et al.*, 2012). Other approaches, such as chromosome
324 painting (Table 1), extend PCA and Structure-like methods by incorporating information
325 about the relative order of markers in the genome, allowing identification of regions for which
326 ancestry differs from the rest of the genome.

327

328 **Heterogeneous structure in space: landscape genomics**

329 Landscape (as well as seascape and lakescape) genetics has widely contributed to our
330 understanding of how ecological and geographical variation affects species history and
331 adaptation (Manel and Holderegger, 2013). Of central importance in this field is the
332 identification of how populations are connected and how organisms move in the landscape
333 matrix. Environmental heterogeneity has a strong impact on how genetic diversity is shaped
334 by migration success between populations, for example after a range expansion (Wegmann *et*
335 *al.*, 2006). A spatially explicit perspective provides context to understand the evolution of
336 locally adapted genes. Moreover, identifying how and where populations (or closely related
337 species, see Roux *et al.* 2016) hybridize is crucial when it comes to characterizing
338 colonization trajectories, tension zones and secondary contacts (Gay *et al.*, 2008; Bierne *et al.*,
339 2011).

340

341 Some methods can explicitly use spatial information to inform clustering, allowing improved
342 consideration of the effect of landscape heterogeneity on selection against migrants and drift.
343 This spatial perspective can be useful to visualize the location and shape of hybrid zones
344 (Guedj and Guillot, 2011). Landscape genetics has valuable application in management and
345 conservation, where it is useful to identify the relevant evolutionary significant units
346 displaying spatial and ecological divergence. Furthermore, researchers are often interested in
347 testing the impact of ecological variation on genetic structure. Mantel tests have been popular
348 to investigate relationships between ecological variables and genetic differentiation while
349 accounting for geographical distances. However, these tests are biased by spatial
350 autocorrelation, assume linear dependence between variables, and do not allow testing the
351 relative contribution of each variable (Legendre and Fortin, 2010; Guillot and Rousset, 2013).
352 Methods such as BEDASSLE (Bradburd *et al.*, 2013) can be used to complement these
353 approaches, and identify which combination of geographical and ecological distance limits
354 dispersal. However, disentangling these effects has proved to be complex and a deeper
355 analysis of genes more strongly impacted by either geography or ecology may be more
356 informative when it comes to the proximate causes of reduced dispersion and differentiation,
357 such as biased dispersal (Edelaar and Bolnick, 2012; Bolnick and Otto, 2013) or selection
358 against migrants (Hendry, 2004). Landscape genomics now extends its focus to adaptive
359 genetic variation, and benefits from new methods targeting signatures of selection (Figure 2
360 and below).

361

362 Population history

363 **Phylogeny**

364 Phylogenetics has a long history that is linked to the broader topic of systematics (Moritz &
365 Hillis 1996; Baum & Smith 2013). Since their inception in the 1980s, molecular phylogenetic
366 methods have been used to address a wide range of problems at different taxonomic scales,
367 including intraspecific population history. Recent advances in molecular phylogenetic
368 methods, and the employment of different types of NGS data is well beyond the scope of this
369 review (see e.g. Moriarty Lemmon & Lemmon 2013; Cruaud et al. 2014; Wen et al. 2015).
370 Rather we focus on the use of phylogeny within the context of studies of intra-specific
371 population history and selection. In this respect, both Maximum Likelihood and Bayesian
372 approaches have become popular to investigate evolutionary relationships between individuals
373 from different populations, even when divergence is very recent (e.g. Wagner *et al.*, 2013).
374 These methods are implemented in softwares such as RAxML (Stamatakis, 2014) and
375 BEAST2 (Drummond and Rambaut, 2007). Ultimately, all molecular phylogenies reconstruct
376 the genealogy of the genes with which they have been constructed. Therefore, a basic
377 assumption when using them to infer lineage history at any taxonomic level (populations,
378 species, and higher taxonomic units) is that the gene tree is representative of lineage history.
379 This assumption is likely to be particularly weak at the population level, since the influences
380 of gene flow, selection, and incomplete lineage sorting are strong at this scale, and may cause
381 gene trees to deviate from population history. Nonetheless, such phylogenies can provide a
382 useful starting point for inferences that are complemented with other methods.

383

384 When using genome-wide data at the population level, methods specifically dedicated to
385 reconstructing multiple species coalescent models (MSC) such as *BEAST (STAR-BEAST)
386 should be preferred over concatenation (Edwards *et al.* 2016), since they allow discordance
387 between species trees and individual gene trees to be identified. Note that these methods can

388 be strongly biased when it comes to estimate divergence times and effective population sizes
389 (Leaché *et al.*, 2014). The impact of gene flow and recombination on phylogenetic methods is
390 however an alley of research that will allow better integration between phylogeny and
391 population genetics (Edwards *et al.*, 2016). Such integration is particularly needed for species
392 and populations that are in the “grey zone of speciation” (Roux *et al.*, 2016). Recent advances
393 in MSC methods handling extremely short, non-recombining fragments (see Chou *et al.*, 2015
394 for a comparison) are promising, especially for datasets such as those produced by GBS.

395 While useful to infer topologies, caution is advised when using branches lengths obtained
396 from SNP-only datasets, e.g. to calculate divergence times between different groups or species
397 (Leaché *et al.*, 2015). For this purpose, it might therefore be easier to extract from the data
398 both variant and invariant sites at several genes or RAD contigs, and analyze the whole
399 sequences in a software like BEAST2. Network methods implemented in Splitstree (Huson
400 and Bryant, 2006), make less assumptions and account for potentially conflicting signals due
401 to high gene flow. Unfortunately, such methods remain mostly descriptive.

402

403 **Approximate Bayesian Computation**

404 Phylogenetic methods tend to be slow for large datasets, and generally do not attempt to
405 account for many effects that are crucial in population genetic interpretation, such as gene
406 flow and recent demographic events within species. A more suitable framework for
407 microevolutionary studies relies on coalescence theory. Population geneticists first developed
408 coalescent theory as a way of modeling the genealogy of alleles from a sample of a large
409 population. Going backward in time, alleles merge (coalesce) in a stochastic way until
410 reaching their most recent common ancestor (Kingman, 1982). Obtaining demographic
411 estimates (e.g. time in years) for parameters usually requires that mutation rate and generation

412 time be known or at least reasonably well estimated, for example from closely-related species
413 with similar life history.

414

415 Computationally fast approaches include Approximate Bayesian Computation (ABC), which
416 compares the empirical data with a set of simulated data produced by coalescent simulations
417 under scenarios predefined by the user (Table 2). By measuring the distance between carefully
418 chosen summary statistics describing each simulation with those from the observed dataset, it
419 is possible to infer which scenario explains the data the best. More information on how to
420 perform ABC analyses are described by Csilléry *et al.* (2010). The main advantage of ABC is
421 that it allows handling any type of marker and arbitrarily complex models, contrary to
422 methods like IMA where the model is predefined. However, using summary statistics leads to
423 the loss of potentially useful information (Robert *et al.*, 2011).

424

425 **Likelihood methods based on the allele frequency spectrum (AFS)**

426 Recently, new likelihood methods based on the AFS emerged to facilitate and speed up the
427 analysis of large SNP datasets. Different patterns of gene flow and demographic events all
428 shape the AFS in specific ways (e.g. alleles are likely to occur at more similar frequencies if
429 divergence is recent or if populations are highly connected). These approaches quickly
430 estimate parameters using composite likelihoods, and do not explicitly take into account
431 correlations induced by LD between physically linked markers (but see ABLE, Table 2). This
432 might limit power to detect recent demographic events (e.g. migration, Jenkins *et al.*, 2012).
433 Including SNPs that are physically close together should not strongly bias parameter
434 estimation. However, such an approach prevents direct comparisons of likelihoods from
435 different models. Therefore, physically independent SNPs should be used to consider

436 composite likelihoods as quasi likelihoods for model comparison (Excoffier *et al.*, 2013).
437 Note that the AFS can also be used as a set of summary statistics for ABC inference. Using
438 allele frequencies estimated from pooled datasets is also feasible, as illustrated by a recent
439 study on hybridization in *Populus* species where AFS was estimated from pooled whole
440 genome resequencing data (Christe *et al.*, 2016).

441

442 The number of mutations found in a given length of DNA sequence directly depends on the
443 mutation rate. One drawback when using SNP data without considering monomorphic sites is
444 that the mutation rate per generation can not be used to convert parameters into demographic
445 estimates (Excoffier *et al.*, 2013). Another possibility consists of calibrating parameter
446 estimates by including a fixed parameter in the analysis, such as population size or divergence
447 time. An issue specific to SNP arrays is ascertainment bias, which is the systematic deviation
448 of allele frequencies from theoretical expectations due to the choice of individuals used at the
449 step of SNP discovery. For example, if SNPs found in one population are the only ones
450 genotyped in another population, a whole set of markers polymorphic in the second
451 population but not in the first will be missed, biasing the AFS (Lachance and Tishkoff, 2013).

452

453 Reaching a high level of precision when estimating demographic parameters can be
454 challenging when information is lacking about the evolutionary history of the species
455 considered. However, even when such information is lacking it is possible to compare the
456 likelihoods of different demographic scenarios, a procedure that has been successfully applied
457 to many species to shed light on the process of speciation (Roux *et al.*, 2016).

458

459 **Methods using whole-genome resequencing**

460 Recently, methods have been developed to infer variation in population sizes with time using
461 the whole genome of just one diploid individual. This began with the Pairwise Sequentially
462 Markovian Coalescent (PSMC, Li and Durbin, 2011), and extensions have been made to this
463 model to allow for several genomes. Such methods have the advantage of requiring only a few
464 individuals, and no *a priori* knowledge of population history. One general drawback,
465 however, is that they are limited to rather simple scenarios, and do not handle more than two
466 populations as yet (but see diCal2, Table 2). While powerful, they are sensitive to
467 confounding factors such as population structure (Orozco-terWengel, 2016) that lead to false
468 signatures of expansion or bottleneck. They also do not allow extremely recent demographic
469 events to be investigated, since the coalescence of two alleles from a single individual in the
470 recent past (a few tens to hundreds generations) is infrequent. Moreover, most of these
471 methods require the data to be phased (but see SMC++, Table 2), for example with fastPhase
472 (Scheet and Stephens, 2006) or BEAGLE (Browning and Browning, 2011). In addition,
473 phasing errors can lead to strong biases in parameters estimates for recent times (Terhorst *et*
474 *al.*, 2016). An extension of these methods takes into account population structure and aims to
475 identify the number of islands contributing to a single genome, assuming it is sampled from a
476 Wright n-island meta-population (Mazet *et al.*, 2015). Such developments should improve the
477 amount of information retrieved from only a few genomes. However, natural populations are
478 structured and connected in complex ways, which can bias demographic inferences, even for
479 popular markers such as mitochondrial sequences (Heller *et al.*, 2013).

480

481 Methods based on tracts of identity-by-descent (IBD, Palamara and Pe'er, 2013) constitute an
482 interesting alternative for more complex model testing when whole genome or densely
483 genotyped datasets are available in large number. Such methods allow recent demographic
484 events to be inferred with relative precision. They are used to predict the length of haplotypes

485 shared by two individuals that are inherited from a common ancestor without recombination.
486 However, IBD detection requires large cohorts and accurate phasing, and therefore application
487 of these methods has been largely restricted to human populations so far (Browning and
488 Browning, 2011; Palamara and Pe'er, 2013). Another approach has used tracts of identity-by-
489 state to perform demographic inference over a range of time-scales (IBS, Harris and Nielsen,
490 2013). IBS tracts are directly observable since they are simply the intervals between pairwise
491 differences in an alignment of sequences and do not require any assumption about coancestry
492 to be defined. The method predicts the length distribution of IBS tracts for pairs of haplotypes
493 under a range of demographic parameters. These predicted spectra are then compared to
494 empirical data under a likelihood framework, as with methods based on the AFS.

495

496 There is currently a tradeoff to be made between methods allowing for arbitrarily complex
497 models that are defined *a priori* by the user (e.g. ABC), and methods that allow population
498 history to be inferred agnostically (e.g. PSMC). While the first category of methods are
499 typically the highest performers at inferring complex population history from a moderate
500 number of markers, it is currently only the second category of methods that are able to make
501 use of the full information provided by whole genome data. Using both methods can therefore
502 help in accurately retrieving the evolutionary history of a given species. For example, a recent
503 study on maize demographic and selective history used both $\hat{\rho}a\hat{\rho}$ and Markovian
504 Coalescent methods to characterize the bottleneck and expansion associated with
505 domestication (Beissinger *et al.*, 2016).

506

507 Screening for selection and association

508 **Selection and its impact on sequence variation**

509 While demographic forces such as drift and migration will affect the whole genome, selection
510 is expected to be specific to particular portions of the genome, and therefore yield
511 discrepancies with genome-wide polymorphism (Lewontin and Krakauer, 1973). Selection
512 affects allele frequencies and polymorphism in predictable ways at the scale of single
513 populations (Charlesworth, 2006; Charlesworth and Charlesworth, 2010). Several statistics
514 summarize them, such as π , the nucleotide diversity (Nei and Li, 1979), Tajima's D (Tajima,
515 1989), and Fay and Wu's H (Fay and Wu, 2000). Using a combination of these statistics
516 allows targets of selection to be identified with greater precision, and minimizes the
517 confounding effects of demography (Nielsen *et al.*, 2005). This approach has been used to
518 develop composite tests, such as the composite likelihood ratio (CLR) test (Nielsen *et al.*,
519 2005) that aim to detect recent selective sweeps.

520

521 **Methods based on population subdivision**

522 When an allele is under positive selection in a population, its frequency tends to rise to
523 fixation, unless gene flow from other populations or strong drift prevents this from happening
524 (Charlesworth *et al.*, 1997). It is therefore possible to contrast patterns of differentiation
525 between populations adapted to their local environment to detect loci under divergent
526 selection (e.g. displaying a high F_{st}). However, it is essential to control for population
527 structure, as it may strongly affect the distribution of differentiation measures and produce
528 high rates of false positives. First attempts to take into account population structure and
529 variation in gene flow included FDIST2 (Beaumont and Nichols, 1996). This method models
530 populations as islands and is aimed at detecting loci under selection by contrasting
531 heterozygosity to F_{st} between populations. More sophisticated methods are now available
532 (Table 3), dedicated to the detection of outliers in large genomic datasets. Most of them
533 correct for relatedness across samples, and can test association between allele frequencies and

534 environmental features (see the extensive review by François *et al.*, 2015). These methods are
535 particularly well suited for the study of RAD-sequencing data, for which allele frequencies are
536 often the only information available in the absence of any reference genome.

537 Detecting association between environment and allele frequencies does not necessarily imply
538 a role for local adaptation. For example, in the case of secondary contact, intrinsic genetic
539 incompatibilities can lead to the emergence of tension zones that may shift until they reach an
540 environmental barrier where they can be trapped (Bierne *et al.*, 2011). Characterizing
541 population history is required to draw conclusions about the possible involvement of a
542 genomic region in adaptation to environment. The sampling strategy must take into account
543 the particular historical and demographic features of the species investigated to gain power
544 (Nielsen *et al.*, 2007). The sequencing strategy must also be carefully considered to control for
545 spatial autocorrelation of genotypes due to isolation by distance and shared demographic
546 history.

547

548 **Genome-wide association**

549 The methods described above focus on allele frequencies at the population scale, but do not
550 test association with traits that vary between individuals within populations (e.g. resistance to
551 a pathogen, symbiotic association, individual size or flowering time). For this task, methods
552 performing Genome-wide association analysis (GWAS) are better suited. The recent
553 development of multivariate methods such as PCAdapt (Duforet-Frebourg *et al.*, 2016) also
554 allow loci putatively under selection to be identified in admixed or continuous populations
555 without requiring information about individual phenotype.

556

557 Uncovering the genetic basis of complex, polygenic traits remains challenging, even in model
558 species (Pritchard and Di Rienzo, 2010; Rockman, 2012). It may be unavoidable as a first step
559 to focus only on traits that are under relatively simple genetic determinism. This can,
560 however, lead to the overrepresentation of loci of major phenotypic effect, a fact that should
561 be acknowledged when discussing the impact of selection on genome variation. The fact that
562 loci of major effect are the easiest to target does not imply that they are necessarily the main
563 substrate of selection (Rockman, 2012). Association methods may help targeting variants
564 undergoing soft sweeps, weak selection or those involved in polygenic control of traits
565 (Pritchard *et al.*, 2010). In such cases, signatures of selection may be subtle and sometimes
566 difficult to retrieve from allele frequency data.

567

568 **Detecting selection with methods focusing on LD**

569 LD is increased and diversity is decreased near a selected allele, especially after recent
570 selection. A class of methods are aimed at targeting those regions that display an excess of
571 long homozygous haplotypes, such as the extended haplotype homozygosity (EHH) test
572 (Sabeti *et al.*, 2002). It is also possible to compare haplotype extension across populations,
573 with the Cross Population Extended Haplotype Homozygosity test (XP-EHH (McCarroll *et al.*,
574 2007)) or Rsb (the standardized ratio of EHH at a given SNP site (Tang *et al.*, 2007)).
575 Individuals included in the analysis should be as distantly related as possible to improve
576 precision and avoid an excess of false positives. These methods require data to be phased in
577 order to reconstruct haplotypes. Statistics dedicated to the detection of selection on standing
578 variation or on multiple alleles (so called soft sweeps) are also available, like the nSL
579 statistics (Ferrer-Admetlla *et al.*, 2014) in selscan or the H2/H1 statistics (Garud *et al.*, 2015),
580 although further studies are still needed to understand to what extent hard and soft sweeps can

581 actually be distinguished (Schridder *et al.*, 2015), as well as their relative importance (Messer
582 and Petrov, 2013; Jensen, 2014).

583

584 Even hard selective sweeps can be challenging to detect with LD-based statistics (Jensen,
585 2014). It is advisable to combine several approaches to improve confidence when pinpointing
586 candidate genes for selection. Methods based on LD alone can sometimes miss the actual
587 variants under selection due to the impact of recombination on local polymorphism that can
588 mimic soft or ongoing hard sweeps (Schridder *et al.*, 2015).

589

590 All LD-based approaches are more powerful with a relatively high density of markers, such as
591 the ones obtained from whole-genome sequencing, SNP-arrays or high-density RAD-seq, and
592 benefit from using statistics focusing on polymorphism and allele sharing. In a recent study of
593 local adaptation in sticklebacks (Roesti *et al.*, 2015), these statistics have been used on dense
594 RAD-sequencing data to look for recent selection at loci displaying high differentiation (F_{ST}).
595 This approach has allowed new candidate loci to be pinpointed, and has confirmed the
596 involvement of those implicated previously (e.g. the *Ectodysplasin* gene). In addition, the
597 identification of large regions displaying high divergence and LD has revealed the importance
598 of large-scale structural variation in shaping genome structure, such as inversions (Roesti *et*
599 *al.*, 2015).

600

601 **Detecting and characterizing selection with the coalescent**

602 If a candidate locus or genomic region has been identified, it is possible to use coalescent
603 simulations to evaluate the strength of selection and estimate the age of alleles. A software
604 such as msms (Ewing and Hermisson, 2010), which is also available in PopGenome, can then

605 be used. However, this requires that population history is known in order to control for other
606 phenomena such as population structure and gene flow. An advantage of full coalescent
607 methods is that they provide a relatively complete picture of the history of individual loci.
608 This can be achieved by modeling coalescence and recombination, and considering variation
609 in mutation rate. However, such methods have long been computationally intensive, and thus
610 difficult to apply to whole genomes. Fortunately, recent computational improvements make
611 their application to whole genomes feasible. A good example is ARGWeaver (Rasmussen *et*
612 *al.*, 2014), which has allowed candidate genes for long-term balancing selection to be
613 recovered from human data. This method uses ancestral recombination graphs to model the
614 genealogy of each non-recombining block in the genome. Ancestral recombination graphs
615 (ARG) are a generalization of the coalescent and describe the sequence of genealogies along a
616 sample of recombining sequence. Genealogies are estimated for each non-recombining block,
617 and recombination between adjacent blocks is described by breaking the branch leading to the
618 recombining haplotype and allowing it to re-coalesce to the rest of the tree. This succession of
619 local trees joined by recombination events provides a full description of the genealogical
620 history of the data and is therefore a promising approach to characterize positive, purifying or
621 balancing selection while taking into account variation in recombination and mutation rate.

622

623 **Identifying variants of functional interest**

624 Characterizing the number of synonymous versus non-synonymous mutations is another
625 approach to detect whether a specific gene is undergoing purifying or positive selection.
626 However, this approach requires an annotated genome. An excess of non-synonymous
627 mutations can signal positive or balancing selection, or a relaxation of selective constraints on
628 a given gene. Annotation of mutations can be done with SNPdat (Doran and Creevey, 2013),
629 or directly in PopGenome, which can also perform tests of selection such as the MK test at the

630 genome scale (McDonald and Kreitman, 1991). Another popular test of selection is the
631 comparison of non-synonymous and synonymous mutations between orthologs from different
632 species, and can be performed in packages such as PAML (Yang, 2007). To recover
633 information about the putative function of a gene or a genomic region, it may be useful to
634 perform a genome ontology (GO) enrichment analysis, using tools such as BLAST2GO
635 (Conesa *et al.*, 2005).

636 While suggestive, genome scans for selection and association in natural populations cannot be
637 considered as conclusive evidence for the function of a given gene, and need to be combined
638 with functional evidence (Vitti *et al.*, 2013). Such evidence can sometimes be provided by
639 variation in the expression of a candidate gene highlighted by RNA-sequencing data. More
640 often, developmental studies are required, a step that is not always possible for non-model
641 organisms. Pinpointing the exact genetic mutation leading to a change in phenotype is
642 challenging even when combining several tests for selection, and requires whole-genome
643 sequencing data to obtain a near-exhaustive list of mutations. It has been proposed to combine
644 QTL analyses with population genomics to facilitate identification of candidate loci
645 (Stinchcombe and Hoekstra, 2008). Essentially, controlled crosses allow genomic regions
646 associated with a selected phenotype to be identified, while the study of variation in natural
647 populations facilitates the fine-mapping of selected variants in natural populations. However,
648 this requires that the species of interest can be raised in a laboratory or greenhouse, which is
649 unpractical for many research teams. An alternative is the study of candidate genes, for which
650 an extensive description of functional variation is available. For example, in a recent study on
651 passerines (bananaquits), GBS data have been used to obtain a neutral distribution to which
652 patterns of substitution and differentiation were compared at candidate genes for color
653 variation (Uy *et al.*, 2016). Another study on color polymorphism in *Peromyscus* mice used a
654 combination of field experiments, targeted sequencing of candidate genes and neutral regions,

655 and genome-scans for selection. Tests for association between these data were able to show
656 how selection on many mutations at the same locus drive adaptive phenotypic divergence
657 (Linnen *et al.*, 2013).

658

659 The combination of tests aimed at different signatures of selection can allow the size of
660 candidate regions to be reduced. For example, combining results from environmental
661 association mapping and genomic scans for selection allows the identification of candidate
662 genes for which a function can be proposed (François *et al.*, 2015). Another common
663 approach relies on the combination of different tests targeting signatures of selection, typically
664 those using the allele frequency spectrum and those using haplotype length. A test of this type
665 has been proposed in human genetics (Grossman *et al.*, 2013), and is called the composite of
666 multiple signals (CMS) test. Nevertheless, signatures of selection can be elusive, and
667 obtaining an exhaustive list of genes under positive selection is unlikely. Further advances
668 will require that methods targeting selection be able to better take into account epistatic
669 interaction and weak selection.

670

671 **Suggestions and perspectives**

672 **Estimating selection and demography jointly along a heterogeneous genome**

673 As stated by Lewontin and Krakauer in 1973, "while natural selection will operate differently
674 for each locus and each allele at a locus, the effect of breeding structure is uniform over all
675 loci and all alleles". Since then, traditional studies on selection have mostly considered that
676 demographic processes act on all loci in the same way across a genome, and that positive
677 selection is mostly rare. This traditional approach has thus tended to disconnect the study of
678 selection from the study of demography (Li *et al.*, 2012).

679

680 However, this assumption may be incorrect, and a joint understanding of demography and
681 selection is crucial from this perspective (Figure 3). For example, the large effective
682 population sizes of *Drosophila* have been hypothesized to facilitate a widespread effect of
683 selection across the genome (Sattath *et al.*, 2011; discussion in Li *et al.*, 2012), making both
684 demographic inference and detection of outliers difficult. Other confounding factors include
685 variation in recombination and mutation rates, and background selection (Ewing and Jensen,
686 2016), which are difficult to assess with precision in non-model species. Moreover, it has
687 been shown in the last few years that loci involved in reproductive isolation are often also
688 involved in local adaptation. This, combined with variation in introgression rates along the
689 genome, can bias inference about selection and demography (Bierne *et al.*, 2011; Roux *et al.*,
690 2014). Genomic regions with low recombination rates can lead to reduced polymorphism, and
691 be mistaken for signatures of purifying selection.

692

693 These issues can only be addressed by going beyond categorization between methods
694 assigned to either the study of selection or demography, and using the results obtained by one
695 method to inform the other. Such an approach was taken by Tine *et al.* 2014 in investigating
696 the two different lineages of the European Sea Bass, using a RAD-sequencing approach. Tine
697 *et al.* took into account variation in recombination rate along the genome to interpret
698 signatures of reduced polymorphism as possibly being the result of selection, low
699 recombination, or a combination of the two (Tine *et al.*, 2014). Since differentiation along the
700 genome seemed to reveal islands resisting gene flow, they could fit a model incorporating
701 variation in introgression rates. This provided improved fit to the data and suggested that
702 islands of differentiation are most likely to be due to locally reduced gene flow after
703 secondary contact. This example illustrates how a combination of descriptive statistics and

704 coalescent analyses can be used to retrieve information from genomic data about both
705 selection and demography.

706

707 Most methods do not actually estimate demography and selection jointly, but rather rely on a
708 process where neutral expectations are first drawn from a set of SNPs presumed to be neutral
709 (e.g. intergenic SNPs), followed by a step where the likelihood of a marker being under
710 selection is evaluated. Methods such as BAYPASS or PCAdapt (Table 3) are convenient in
711 both describing population structure and providing preliminary insights into the proportion of
712 loci that do not follow neutral expectations. If this proportion is high, it would suggest recent
713 introgression or an excess of markers displaying high LD (e.g. due to large inversions).

714 However, when this proportion is not too high, outliers can be removed to avoid bias
715 (Schridder *et al.*, 2016) and the remaining loci used to compare neutral models and estimate
716 demographic parameters (e.g. using an ABC framework). These estimated parameters can
717 then be used to simulate sequences or independent SNPs and generate a neutral expectation.
718 Loci that are more likely to be neutral can be used to further calibrate tests for selection such
719 as FLK or BAYPASS (Lotterhos and Whitlock, 2014).

720

721 Some recent methods are especially relevant to study both demography and selection at once,
722 while taking into account variation in recombination and mutation rates. For whole-genome
723 data, methods reconstructing ancestral recombination graphs (such as ARGWeaver) have high
724 potential. They allow genealogies to be retrieved along the genome as well as the timing of
725 coalescence events. Such information is ultimately useful for making inferences regarding
726 selection and migration. Recently this method was used in human paleogenomics to
727 quantitatively characterize introgression between modern humans, Neandertals and

728 Denisovans using only a few whole genomes (Kuhlwilm *et al.*, 2016). However, the approach
729 has a high computing and sequencing cost, and is therefore not suitable for studies requiring
730 sampling of many individuals.

731

732 Caution must prevail when attempting to apply sophisticated methods to disentangle selection
733 and demography. In a recent review, Cruickshank & Hahn suggest that IMA2, which is
734 commonly used to estimate migration rates, is not able to reliably distinguish between loci
735 resisting gene flow, and those under selection in the absence of gene flow (Cruickshank and
736 Hahn, 2014). In the specific case they highlight (*Oryctolagus cuniculus* rabbits, Sousa *et al.*,
737 2013), a descriptive statistic that should have captured introgression signatures (d_{xy}) did not
738 reveal any evidence for differential gene flow between loci categorized by IMA2. This
739 controversy illustrates that basic description of the data is needed prior using more
740 sophisticated methods. Note however that Cruickshank & Hahn did not address the case of
741 secondary contact, and other methods such as ABC may better detect interruption in gene
742 flow (Sup. Text in Roux *et al.*, 2016).

743

744 To sum up, the field of population genomics is now moving towards both better integrating
745 the demographic framework in inferences of selection, and, conversely, taking into account
746 selection when reconstructing demographic history. The joint inference of loci under selection
747 and quantification of demographic dynamics is of crucial importance in fields such as
748 landscape genomics or the study of ongoing speciation. It should provide insights into the role
749 of selection, recombination and gene flow in promoting or impairing local adaptation to new
750 habitats. The growing availability of genome-wide data for non-model species is therefore
751 promising, but requires caution and high stringency in our interpretation of observed patterns.

752 With the decreasing cost of sequencing, it has been suggested that NGS will rapidly broaden
753 our perspective on complex evolutionary processes, from biogeography (Lexer *et al.*, 2013) to
754 the genetic basis of traits (Hohenlohe, 2014) and the maintenance of polymorphisms (Hedrick,
755 2006). While genome heterogeneity in migration, mutation and recombination rates do not
756 necessarily make impossible any conclusion about evolutionary dynamics, they have the
757 potential to blur inferences. The study of DNA sequence variation is already challenging in its
758 own right. Nonetheless, in order to be informative about processes such as selection and
759 demography it should ultimately be combined with other disciplines such as ecology and
760 functional analyses (Habel *et al.*, 2015). This can be done for example by assessing the
761 function of selected genes, the consistency of demographic history with information retrieved
762 from the fossil record or geological history, and the broader integration of population
763 genomics with other fields and methods whenever possible, such as niche modeling, common
764 garden experiments or the study of macro-evolutionary patterns of selection and
765 diversification.

766

767 **Beyond SNPs: studying structural variation, transposable elements and epigenetic**
768 **modifications**

769 Most genome-scale studies of selection and demography have so far focused on SNPs, since
770 they are relatively easy to detect with current technology and their mutation mechanism
771 produces mostly biallelic alleles, making them easier to use for statistical tests. However,
772 many other heritable genetic alterations can affect genomes, including insertions of
773 transposable elements (Villanueva-Cañas *et al.*, 2017), epigenetics modifications such as
774 methylation (Danchin *et al.*, 2011), duplications, inversions, deletions and translocations
775 (Iskow *et al.*, 2012). One of the main issues with this type of variation is that their diversity
776 and their impact on the genome can make them difficult to detect in a systematic way (Iskow

777 *et al.*, 2012), especially for species for which only a draft genome is available. It is however
778 possible to use variation in such genetic alternations to study selection, for example by using
779 differentiation statistics, association to environment or extension of haplotypes. Combining
780 information about variant position and SNP variation in flanking regions is also a powerful
781 way to detect variants under selection (Villanueva-Cañas *et al.*, 2017) as highlighted by a
782 recent study of transposable element insertions in *Drosophila* (Kofler *et al.*, 2012). Recent
783 work also shows that classical summary statistics such as Tajima's D can be adapted to non-
784 SNP datasets, such as methylations (Wang and Fan, 2014).

785 Sets of neutral SNPs can be used to control for demography and relatedness between samples
786 when inferring selection. For example, this type of approach has recently been adopted in
787 studies of selection on methylation patterns. In a recent *Molecular Ecology* issue (Verhoeven
788 *et al.*, 2016), a study using bisulfite precipitation in Valley Oak trees (Gugger *et al.* 2016) was
789 able to place methylated variants associated to climatic variables near to genes known to be
790 involved in response to environment. Another study could show a stronger pattern of Isolation
791 by Distance for methylation-sensitive AFLPs than for regular AFLPs and microsatellites,
792 suggesting a stronger impact of environment on methylation patterns than expected under
793 neutrality (Herrera *et al.*, 2016).

794

795 Another potential issue with this type of variation is that there is currently a lack of tools able
796 to simulate their models of mutation, complicating any comparison with neutral models built
797 from SNPs. This is particularly true for transposable elements, for which the assumption of
798 mutation-drift equilibrium is challenging, making comparisons of their allele frequency
799 spectrum with neutral SNPs potentially difficult. For example, a recent burst of transposition
800 can lead to an excess of low frequency elements and recent insertions compared to the
801 expectation under equilibrium, even if transposable elements are not under purifying selection

802 (Bergman and Bensasson, 2007; Blumenstiel *et al.*, 2014). More generally, neutral models
803 would benefit from new ways to model the appearance of genomic variation through time for
804 non-SNP data. This would provide even more conservative assessments of negative and
805 positive selection.

806

807 Acknowledgements

808 The University of Basel and New York University Abu Dhabi have supported YB's research
809 in this area. We want to thank two anonymous reviewers, Stephane Boissinot, Joris Bertrand
810 and Anne Roulin for their insightful comments on previous versions of the manuscript.

811

812 References

- 813 Abi-Rached L, Jobin M, Kulkarni S, McWhinnie A, Dalva K, Gragert L, *et al.* (2011). The
814 shaping of modern human immune systems by multiregional admixture with archaic
815 humans. *Science* **334**: 89–95.
- 816 Abzhanov A, Extavour CG, Groover A, Hodges SA, Hoekstra HE, Kramer EM, *et al.* (2008).
817 Are we there yet? Tracking the development of new model systems. *Trends Genet* **24**:
818 353–60.
- 819 Albrechtsen A, Nielsen FC, Nielsen R (2010). Ascertainment biases in SNP chips affect
820 measures of population divergence. *Mol Biol Evol* **27**: 2534–2547.
- 821 Alexander DH, Novembre J (2009). Fast Model-Based Estimation of Ancestry in Unrelated
822 Individuals. *Genome Res*: 1655–1664.
- 823 Alonso-Blanco C, Andrade J, Becker C, Bemm F, Bergelson J, Borgwardt KMM, *et al.*
824 (2016). 1,135 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis*
825 *thaliana*. *Cell* **166**: 481–491.
- 826 Arnold B, Corbett-Detig RB, Hartl D, Bomblies K (2013). RADseq underestimates diversity
827 and introduces genealogical biases due to nonrandom haplotype sampling. *Mol Ecol* **22**:

- 828 3179–90.
- 829 Aulchenko YS, Ripke S, Isaacs A, van Duijn CM (2007). GenABEL: An R library for
830 genome-wide association analysis. *Bioinformatics* **23**: 1294–1296.
- 831 Axelsson E, Ratnakumar A, Arendt M-L, Maqbool K, Webster MT, Perloski M, *et al.* (2013).
832 The genomic signature of dog domestication reveals adaptation to a starch-rich diet.
833 *Nature* **495**: 360–4.
- 834 Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, *et al.* (2008). Rapid SNP
835 discovery and genetic mapping using sequenced RAD markers. *PLoS One* **3**: e3376.
- 836 Balding DJ (2006). A tutorial on statistical methods for population association studies. *Nat*
837 *Rev Genet* **7**: 781–91.
- 838 Baran Y, Pasaniuc B, Sankararaman S, Torgerson DG, Gignoux C, Eng C, *et al.* (2012). Fast
839 and accurate inference of local ancestry in Latino populations. *Bioinformatics* **28**: 1359–
840 1367.
- 841 Beaumont MA, Balding DJ (2004). Identifying adaptive genetic divergence among
842 populations from genome scans. *Mol Ecol* **13**: 969–980.
- 843 Beaumont MA, Nichols RA (1996). Evaluating loci for use in the genetic analysis of
844 population structure. *Proc R Soc London Biol Sci*: 1619–1626.
- 845 Beeravolu CR, Hickerson MJ, Frantz LAF, Lohse K (2016). Approximate Likelihood
846 Inference of Complex Population Histories and Recombination from Multiple Genomes.
847 *bioarXiv*: 1–31.
- 848 Beerli P, Palczewski M (2010). Unified framework to evaluate panmixia and migration
849 direction among multiple sampling locations. *Genetics* **185**: 313–26.
- 850 Beissinger TM, Wang L, Crosby K, Durvasula A, Hufford MB, Ross-Ibarra J (2016). Recent
851 demography drives changes in linked selection across the maize genome. *Nat Plants* **2**:
852 16084.
- 853 Bergman CM, Bensasson D (2007). Recent LTR retrotransposon insertion contrasts with
854 waves of non-LTR insertion since speciation in *Drosophila melanogaster*. *Proc Natl*
855 *Acad Sci U S A* **104**: 11340–11345.
- 856 Besnard G, Bertrand JAM, Delahaie B, Bourgeois YXC, Lhuillier E, Thébaud C (2016).

- 857 Valuing museum specimens: high-throughput DNA sequencing on historical collections
858 of New Guinea crowned pigeons (Goura). *Biol J Linn Soc* **117**: 71–82.
- 859 Bierne N, Welch J, Loire E, Bonhomme F, David P (2011). The coupling hypothesis: why
860 genome scans may fail to map local adaptation genes. *Mol Ecol* **20**: 2044–72.
- 861 Blumenstiel JP, Chen X, He M, Bergman CM (2014). An age-of-allele test of neutrality for
862 transposable element insertions. *Genetics* **196**: 523–538.
- 863 Boistard S, Rodriguez W, Jay F, Mona S, Austerlitz F (2016). Inferring Population Size
864 History from Large Samples of Genome-Wide Molecular Data - An Approximate
865 Bayesian Computation Approach. *PLoS Genet*: 858–865.
- 866 Bolnick DI, Otto SP (2013). The magnitude of local adaptation under genotype-dependent
867 dispersal. *Ecol Evol* **3**: 4722–4735.
- 868 Bonhomme M, Chevalet C, Servin B, Boitard S, Abdallah JM, Blott S, *et al.* (2010).
869 Detecting Selection in Population Trees: The Lewontin and Krakauer Test Extended.
870 *Genetics*: 241–262.
- 871 Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu CH, Xie D, *et al.* (2014). BEAST 2: A
872 Software Platform for Bayesian Evolutionary Analysis. *PLoS Comput Biol* **10**: 1–6.
- 873 Bradburd GS, Ralph PL, Coop GM (2013). Disentangling the effects of geographic and
874 ecological isolation on genetic differentiation. *Evolution (N Y)* **67**: 3258–3273.
- 875 Brisbin A, Bryc K, Byrnes J, Zakharia F, Omberg L, Degenhardt J, *et al.* (2012). PCAdmix:
876 Principal Components-Based Assignment of Ancestry along Each Chromosome in
877 Individuals with Admixed Ancestry from Two or More Populations. *Hum Biol* **84**: 343–
878 364.
- 879 Browning BL, Browning SR (2011). A fast, powerful method for detecting identity by
880 descent. *Am J Hum Genet* **88**: 173–182.
- 881 Bryant D, Bouckaert R, Felsenstein J, Rosenberg NA, Roychoudhury A (2012). Inferring
882 species trees directly from biallelic genetic markers: Bypassing gene trees in a full
883 coalescent analysis. *Mol Biol Evol* **29**: 1917–1932.
- 884 Buerkle CA, Gompert Z (2013). Population genomics based on low coverage sequencing:
885 how low should we go? *Mol Ecol* **22**: 3028–35.

- 886 Cadzow M, Boocock J, Nguyen HT, Wilcox P, Merriman TR, Black MA (2014). A
887 bioinformatics workflow for detecting signatures of selection in genomic data. *Front*
888 *Genet* **5**: 1–8.
- 889 Cariou M, Duret L, Charlat S (2013). Is RAD-seq suitable for phylogenetic inference? An in
890 silico assessment and optimization. *Ecol Evol* **3**: 846–852.
- 891 Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011). Stacks: building
892 and genotyping Loci de novo from short-read sequences. *G3 (Bethesda)* **1**: 171–82.
- 893 Catchen JM, Hohenlohe PA, Bernatchez L, Funk WC, Andrews KR, Allendorf FW (2017).
894 Unbroken: RADseq remains a powerful tool for understanding the genetics of adaptation
895 in natural populations. *Mol Ecol Resour*: 362–365.
- 896 Charlesworth D (2006). Balancing selection and its effects on sequences in nearby genome
897 regions. *PLoS Genet* **2**: e64.
- 898 Charlesworth B, Charlesworth D (2010). *Elements of evolutionary genetics*.
- 899 Charlesworth B, Nordborg M, Charlesworth D (1997). The effects of local selection, balanced
900 polymorphism and background selection on equilibrium patterns of genetic diversity in
901 subdivided populations. *Genet Res, Camb* **70**: 155–174.
- 902 Chifman J, Kubatko L (2014). Quartet inference from SNP data under the coalescent model.
903 *Bioinformatics* **30**: 3317–3324.
- 904 Chikhi L, Sousa VC, Luisi P, Goossens B, Beaumont MA (2010). The confounding effects of
905 population structure, genetic diversity and the sampling scheme on the detection and
906 quantification of population size changes. *Genetics* **186**: 983–995.
- 907 Chou J, Gupta A, Yaduvanshi S, Davidson R, Nute M, Mirarab S, *et al.* (2015). A
908 comparative study of SVDquartets and other coalescent-based species tree estimation
909 methods. *BMC Genomics* **16**: S2.
- 910 Christe C, Stolting KN, Paris M, Frayisse C, Bierne N, Lexer C (2016). Adaptive evolution
911 and segregating load contribute to the genomic landscape of divergence in two tree
912 species connected by episodic gene flow. *Mol Ecol*.
- 913 Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M (2005). Blast2GO: A
914 universal tool for annotation, visualization and analysis in functional genomics research.

- 915 *Bioinformatics* **21**: 3674–3676.
- 916 Cornuet J-M, Santos F, Beaumont MA, Robert CP, Marin J-M, Balding DJ, *et al.* (2008).
917 Inferring population history with DIY ABC: a user-friendly approach to approximate
918 Bayesian computation. *Bioinformatics* **24**: 2713–9.
- 919 Cruickshank TE, Hahn MW (2014). Reanalysis suggests that genomic islands of speciation
920 are due to reduced diversity, not reduced gene flow. *Mol Ecol* **23**: 3133–3157.
- 921 Csilléry K, Blum MGB, Gaggiotti OE, François O (2010). Approximate Bayesian
922 Computation (ABC) in practice. *Trends Ecol Evol* **25**: 410–8.
- 923 Csilléry K, François O, Blum MGB (2012). abc: an R package for approximate Bayesian
924 computation (ABC). *Methods Ecol Evol* **3**: 475–479.
- 925 Cushman SA (2014). Grand challenges in evolutionary and population genetics: The
926 importance of integrating epigenetics, genomics, modeling, and experimentation. *Front*
927 *Genet* **5**: 1–5.
- 928 Danchin É, Charmantier A, Champagne F a, Mesoudi A, Pujol B, Blanchet S (2011). Beyond
929 DNA: integrating inclusive inheritance into an extended theory of evolution. *Nat Rev*
930 *Genet* **12**: 475–86.
- 931 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, *et al.* (2011). The
932 variant call format and VCFtools. *Bioinformatics* **27**: 2156–2158.
- 933 Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML (2011). Genome-
934 wide genetic marker discovery and genotyping using next-generation sequencing. *Nat*
935 *Rev Genet* **12**: 499–510.
- 936 Degiorgio M, Huber CD, Hubisz MJ, Hellmann I, Nielsen R (2016). Genetics and population
937 analysis SWEEPfinder 2: Increased sensitivity, robustness, and flexibility.
938 *Bioinformatics*.
- 939 DeGiorgio M, Lohmueller KE, Nielsen R (2014). A model-based approach for identifying
940 signatures of ancient balancing selection in genetic data. *PLoS Genet* **10**: e1004561.
- 941 Doran AG, Creevey CJ (2013). Snpdat: easy and rapid annotation of results from de novo snp
942 discovery projects for model and non-model organisms. *BMC Bioinformatics* **14**: 45.
- 943 Drummond AJ, Rambaut A (2007). BEAST: Bayesian evolutionary analysis by sampling

- 944 trees. *BMC Evol Biol* **7**: 214.
- 945 Duforet-Frebourg N, Luu K, Laval G, Bazin E, Blum MGB (2016). Detecting genomic
946 signatures of natural selection with principal component analysis: Application to the
947 1000 genomes data. *Mol Biol Evol* **33**: 1082–1093.
- 948 Dufresne F, Stift M, Vergilino R, Mable BK (2014). Recent progress and challenges in
949 population genetics of polyploid organisms: An overview of current state-of-the-art
950 molecular and statistical tools. *Mol Ecol* **23**: 40–69.
- 951 Edelaar P, Bolnick DI (2012). Non-random gene flow: an underappreciated force in evolution
952 and ecology. *Trends Ecol Evol* **27**: 659–65.
- 953 Edwards S V., Xi Z, Janke A, Faircloth BC, McCormack JE, Glenn TC, *et al.* (2016).
954 Implementing and testing the multispecies coalescent model: A valuable paradigm for
955 phylogenomics. *Mol Phylogenet Evol* **94**: 447–462.
- 956 Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, *et al.* (2012). The
957 genomic landscape of species divergence in *Ficedula* flycatchers. *Nature* **491**: 756–60.
- 958 Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, *et al.* (2011). A
959 Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species.
960 *PLoS One* **6**: e19379.
- 961 Ewing G, Hermisson J (2010). MSMS: A coalescent simulation program including
962 recombination, demographic structure and selection at a single locus. *Bioinformatics* **26**:
963 2064–2065.
- 964 Ewing GB, Jensen JD (2016). The consequences of not accounting for background selection
965 in demographic inference. *Mol Ecol* **25**: 135–141.
- 966 Excoffier L, Dupanloup I, Huerta-Sanchez E, Sousa VC, Foll M (2013). Robust Demographic
967 Inference from Genomic and SNP Data. *PLoS Genet* **9**.
- 968 Excoffier L, Foll M (2011). Fastsimcoal: a Continuous-Time Coalescent Simulator of
969 Genomic Diversity Under Arbitrarily Complex Evolutionary Scenarios. *Bioinformatics*
970 **27**: 1332–4.
- 971 Excoffier L, Heckel G (2006). Computer programs for population genetics data analysis: a
972 survival guide. *Nat Rev Genet* **7**: 745–58.

- 973 Excoffier L, Lischer HEL (2010). Arlequin suite ver 3.5: a new series of programs to perform
974 population genetics analyses under Linux and Windows. *Mol Ecol Resour* **10**: 564–7.
- 975 Falush D, Dorp L van, Lawson D (2016). A tutorial on how (not) to over-interpret
976 STRUCTURE/ADMIXTURE bar plots. *bioRxiv*: 66431.
- 977 Fay JC, Wu CI (2000). Hitchhiking under positive Darwinian selection. *Genetics* **155**: 1405–
978 13.
- 979 Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R (2014). On detecting incomplete soft
980 or hard selective sweeps using haplotype structure. *Mol Biol Evol* **31**: 1275–1291.
- 981 Ferretti L, Ramos-Onsins SE, Pérez-Enciso M (2013). Population genomics from pool
982 sequencing. *Mol Ecol* **22**: 5561–5576.
- 983 Foll M, Gaggiotti O (2008). A genome-scan method to identify selected loci appropriate for
984 both dominant and codominant markers: a Bayesian perspective. *Genetics* **180**: 977–93.
- 985 da Fonseca RR, Albrechtsen A, Themudo GE, Ramos-Madrigal J, Sibbesen JA, Maretty L, *et*
986 *al.* (2016). Next-generation biology: Sequencing and data analysis approaches for non-
987 model organisms. *Mar Genomics* **30**: 1–11.
- 988 François O, Martins H, Caye K, Schoville SD (2015). Controlling False Discoveries in
989 Genome Scans for Selection. *Mol Ecol* **55**: in press.
- 990 Fraser DJ, Bernatchez L (2001). Adaptive evolutionary conservation: Towards a unified
991 concept for defining conservation units. *Mol Ecol* **10**: 2741–2752.
- 992 Frichot E, Mathieu F, Trouillon T, Bouchard G, François O (2014). Fast and efficient
993 estimation of individual ancestry coefficients. *Genetics* **196**: 973–983.
- 994 Frichot E, Schoville SD, Bouchard G, François O (2013). Testing for associations between
995 loci and environmental gradients using latent factor mixed models. *Mol Biol Evol* **30**:
996 1687–1699.
- 997 Futschik A, Schlötterer C (2010). The next generation of molecular markers from massively
998 parallel sequencing of pooled DNA samples. *Genetics* **186**: 207–18.
- 999 Gardner SN, Hall BG (2013). When whole-genome alignments just won't work: KSNP v2
1000 software for alignment-free SNP discovery and phylogenetics of hundreds of microbial
1001 genomes. *PLoS One* **8**.

- 1002 Garrigan D (2013). POPBAM: Tools for evolutionary analysis of short read sequence
1003 alignments. *Evol Bioinforma* **2013**: 343–353.
- 1004 Garud NR, Messer PW, Buzbas EO, Petrov DA (2015). Recent Selective Sweeps in North
1005 American *Drosophila melanogaster* Show Signatures of Soft Sweeps. *PLoS Genet* **11**: 1–
1006 32.
- 1007 Gautier M (2015). Genome-Wide Scan for Adaptive Divergence and Association with
1008 Population-Specific Covariates. *Genetics* **201**: 1555–1579.
- 1009 Gautier M, Laloë D, Moazami-Goudarzi K (2010). Insights into the genetic history of French
1010 cattle from dense SNP data on 47 worldwide breeds. *PLoS One* **5**: 1–11.
- 1011 Gautier M, Vitalis R (2012). Rehh An R package to detect footprints of selection in genome-
1012 wide SNP data from haplotype structure. *Bioinformatics* **28**: 1176–1177.
- 1013 Gautier M, Vitalis R (2013). Inferring population histories using genome-wide allele
1014 frequency data. *Mol Biol Evol* **30**: 654–668.
- 1015 Gay L, Crochet P-A, Bell D a, Lenormand T (2008). Comparing clines on molecular and
1016 phenotypic traits in hybrid zones: a window on tension zone models. *Evolution* **62**:
1017 2789–806.
- 1018 Gayral P, Melo-Ferreira J, Glémin S, Bierne N, Carneiro M, Nabholz B, *et al.* (2013).
1019 Reference-Free Population Genomics from Next-Generation Transcriptome Data and the
1020 Vertebrate-Invertebrate Gap. *PLoS Genet* **9**.
- 1021 Glaubitz JC, Casstevens TM, Lu F, Harriman J, Elshire RJ, Sun Q, *et al.* (2014). TASSEL-
1022 GBS: A high capacity genotyping by sequencing analysis pipeline. *PLoS One* **9**.
- 1023 Gronau I, Hubisz MJ, Gulko B, Danko CG, Siepel A (2011). Bayesian inference of ancient
1024 human demography from individual genome sequences. *Nat Genet* **43**: 1031–1034.
- 1025 Grossman SR, Andersen KG, Shlyakhter I, Tabrizi S, Winnicki S, Yen A, *et al.* (2013).
1026 Identifying recent adaptations in large-scale genomic data. *Cell* **152**: 703–713.
- 1027 Grover CE, Salmon A, Wendel JF (2012). Targeted sequence capture as a powerful tool for
1028 evolutionary analysis. *Am J Bot* **99**: 312–9.
- 1029 Guedj B, Guillot G (2011). Estimating the location and shape of hybrid zones. *Mol Ecol*
1030 *Resour* **11**: 1119–1123.

- 1031 Guillot G, Renaud S, Ledevin R, Michaux J, Claude J (2012). A unifying model for the
1032 analysis of phenotypic, genetic, and geographic data. *Syst Biol* **61**: 897–911.
- 1033 Guillot G, Rousset F (2013). Dismantling the Mantel tests. *Methods Ecol Evol* **4**: 336–344.
- 1034 Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O (2010). New
1035 algorithms and methods to estimate maximum-likelihood phylogenies: Assessing the
1036 performance of PhyML 3.0. *Syst Biol* **59**: 307–321.
- 1037 Günther T, Coop G (2013). Robust identification of local adaptation from allele frequencies.
1038 *Genetics* **195**: 205–220.
- 1039 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD (2009). Inferring the joint
1040 demographic history of multiple populations from multidimensional SNP frequency data.
1041 *PLoS Genet* **5**.
- 1042 Habel J, Zachos F, Dapporto L, Rödder D, Radespiel U, Tellier A, *et al.* (2015). Population
1043 genetics revisited – towards a multidisciplinary research field. *Biol J Linn Soc* **115**: 1–12.
- 1044 Hahn C, Bachmann L, Chevreux B (2013). Reconstructing mitochondrial genomes directly
1045 from genomic next-generation sequencing reads - A baiting and iterative mapping
1046 approach. *Nucleic Acids Res* **41**.
- 1047 Hand BK, Hether TD, Kovach RP, Muhlfeld CC, Amish SJ, Boyer MC, *et al.* (2015).
1048 Genomics and introgression: Discovery and mapping of thousands of species-diagnostic
1049 SNPs using RAD sequencing. *Curr Zool* **61**: 146–154.
- 1050 Harris K, Nielsen R (2013). Inferring Demographic History from a Spectrum of Shared
1051 Haplotype Lengths. *PLoS Genet* **9**.
- 1052 Hedrick PW (2006). Genetic Polymorphism in Heterogeneous Environments: The Age of
1053 Genomics. *Annu Rev Ecol Syst* **37**: 67–93.
- 1054 Hedrick PW (2013). Adaptive introgression in animals: Examples and comparison to new
1055 mutation and standing variation as sources of adaptive variation. *Mol Ecol* **22**: 4606–
1056 4618.
- 1057 Heled J, Drummond AJ (2010). Bayesian Inference of Species Trees from Multilocus Data.
1058 *Mol Biol Evol* **27**: 570–580.
- 1059 Hellenthal G, Busby GBJ, Band G, Wilson JF, Capelli C, Falush D, *et al.* (2014). A Genetic

- 1060 Atlas of Human Admixture History. *Science* **343**: 747–751.
- 1061 Heller R, Chikhi L, Siegmund HR (2013). The Confounding Effect of Population Structure
1062 on Bayesian Skyline Plot Inferences of Demographic History. *PLoS One* **8**.
- 1063 Hendry AP (2004). Selection against migrants contributes to the rapid evolution of
1064 ecologically dependent reproductive isolation. *Evol Ecol Res* **6**: 1219–1236.
- 1065 Herrera CM, Medrano M, Bazaga P (2016). Comparative spatial genetics and epigenetics of
1066 plant populations: Heuristic value and a proof of concept. *Mol Ecol* **25**: 1653–1664.
- 1067 Hey J, Nielsen R (2007). Integration within the Felsenstein equation for improved Markov
1068 chain Monte Carlo methods in population genetics. *Proc Natl Acad Sci U S A* **104**: 2785–
1069 90.
- 1070 Hohenlohe PA (2014). Ecological genomics in full colour. *Mol Ecol* **23**: 5129–31.
- 1071 Hudson RR (2002). Generating samples under a Wright–Fisher neutral model of genetic
1072 variation. *Bioinformatics* **18**: 337–338.
- 1073 Huson DH, Bryant D (2006). Application of phylogenetic networks in evolutionary studies.
1074 *Mol Biol Evol* **23**: 254–267.
- 1075 Iskow RC, Gokcumen O, Lee C (2012). Exploring the role of copy number variants in human
1076 adaptation. *Trends Genet* **28**: 245–257.
- 1077 Jenkins PA, Song YS, Brem RB (2012). Genealogy-Based Methods for Inference of
1078 Historical Recombination and Gene Flow and Their Application in *Saccharomyces*
1079 *cerevisiae*. *PLoS One* **7**.
- 1080 Jenner RA, Wills MA (2007). The choice of model organisms in evo-devo. *Nat Rev Genet* **8**:
1081 311–319.
- 1082 Jensen JD (2014). On the unfounded enthusiasm for soft selective sweeps. *Nat Commun* **5**:
1083 5281.
- 1084 Johnston SE, McEwan JC, Pickering NK, Kijas JW, Beraldi D, Pilkington JG, *et al.* (2011).
1085 Genome-wide association mapping identifies the genetic basis of discrete and
1086 quantitative variation in sexual weaponry in a wild sheep population. *Mol Ecol* **20**: 2555–
1087 2566.
- 1088 Jombart T, Devillard S, Balloux F, Falush D, Stephens M, Pritchard J, *et al.* (2010).

- 1089 Discriminant analysis of principal components: a new method for the analysis of
1090 genetically structured populations. *BMC Genet* **11**: 94.
- 1091 Jombart T, Devillard S, Dufour a-B, Pontier D (2008). Revealing cryptic spatial patterns in
1092 genetic variability by a new multivariate method. *Heredity (Edinb)* **101**: 92–103.
- 1093 Jombart T, Pontier D, Dufour A-B (2009). Genetic markers in the playground of multivariate
1094 analysis. *Heredity (Edinb)* **102**: 330–41.
- 1095 Jones MR, Good JM (2016). Targeted capture in evolutionary and ecological genomics. *Mol*
1096 *Ecol* **25**: 185–202.
- 1097 Jostins L, McVean G (2016). Trinculo: Bayesian and frequentist multinomial logistic
1098 regression for genome-wide association studies of multi-category phenotypes.
1099 *Bioinformatics* **32**: 1898–1900.
- 1100 Kempainen P, Knight CG, Sarma DK, Hlaing T, Prakash A, Maung Maung YN, *et al.*
1101 (2015). Linkage disequilibrium network analysis (LDna) gives a global view of
1102 chromosomal inversions, local adaptation and geographic structure. *Mol Ecol Resour*:
1103 1031–1045.
- 1104 Kern AD, Schrider DR (2016). Discoal: flexible coalescent simulations with selection.
1105 *Bioinformatics* **32**: 3839–3841.
- 1106 Kingman JFC (1982). The coalescent. *Stoch Process their Appl* **13**: 235–248.
- 1107 Kofler R, Betancourt AJ, Schlötterer C (2012). Sequencing of pooled DNA samples (Pool-
1108 Seq) uncovers complex dynamics of transposable element insertions in *Drosophila*
1109 *melanogaster*. *PLoS Genet* **8**.
- 1110 Kofler R, Orozco-terWengel P, De Maio N, Pandey RV, Nolte V, Futschik A, *et al.* (2011).
1111 PoPoolation: a toolbox for population genetic analysis of next generation sequencing
1112 data from pooled individuals. *PLoS One* **6**: e15925.
- 1113 Kofler R, Pandey RV, Schlötterer C (2011). PoPoolation2: identifying differentiation between
1114 populations using sequencing of pooled DNA samples (Pool-Seq). *Bioinformatics* **27**:
1115 3435–6.
- 1116 Kolaczkowski B, Kern AD, Holloway AK, Begun DJ (2011). Genomic differentiation
1117 between temperate and tropical Australian populations of *Drosophila melanogaster*.

- 1118 *Genetics* **187**: 245–60.
- 1119 Korneliussen TS, Albrechtsen A, Nielsen R (2014). ANGSD: Analysis of Next Generation
1120 Sequencing Data. *BMC Bioinformatics* **15**: 356.
- 1121 Kubota S, Iwasaki T, Hanada K, Nagano AJ, Fujiyama A, Toyoda A, *et al.* (2015). A Genome
1122 Scan for Genes Underlying Microgeographic-Scale Local Adaptation in a Wild
1123 Arabidopsis Species. *PLoS Genet* **11**: 1–26.
- 1124 Kuhlwilm M, Gronau I, Hubisz MJ, de Filippo C, Prado-Martinez J, Kircher M, *et al.* (2016).
1125 Ancient gene flow from early modern humans into Eastern Neanderthals. *Nature* **530**:
1126 429–433.
- 1127 Lachance J, Tishkoff SA (2013). SNP ascertainment bias in population genetic analyses: Why
1128 it is important, and how to correct it. *Bioessays* **35**: 780–786.
- 1129 Laland KN, Sterelny K, Odling-Smee J, Hoppitt W, Uller T (2011). Cause and effect in
1130 biology revisited: is Mayr’s proximate-ultimate dichotomy still useful? *Science* **334**:
1131 1512–6.
- 1132 Lee T-H, Guo H, Wang X, Kim C, Paterson AH (2014). SNPhylo: a pipeline to construct a
1133 phylogenetic tree from huge SNP data. *BMC Genomics* **15**: 162.
- 1134 Legendre P, Fortin MJ (2010). Comparison of the Mantel test and alternative approaches for
1135 detecting complex multivariate relationships in the spatial analysis of genetic data. *Mol*
1136 *Ecol Resour* **10**: 831–844.
- 1137 Legrand D, Tenaillon MI, Matyot P, Gerlach J, Lachaise D, Cariou M-L (2009). Species-wide
1138 genetic variation and demographic history of *Drosophila sechellia*, a species lacking
1139 population structure. *Genetics* **182**: 1197–206.
- 1140 Lewontin RC, Krakauer J (1973). Distribution of gene frequency as a test of the theory of the
1141 selective neutrality of polymorphisms. *Genetics* **74**: 175–195.
- 1142 Lexer C, Mangili S, Bossolini E, Forest F, Stölting KN, Pearman PB, *et al.* (2013). ‘Next
1143 generation’ biogeography: towards understanding the drivers of species diversification
1144 and persistence (M Carine, Ed.). *J Biogeogr* **40**: 1013–1022.
- 1145 Li H, Durbin R (2011). Inference of human population history from individual whole-genome
1146 sequences. *Nature* **475**: 493–496.

- 1147 Li J, Li H, Jakobsson M, Li S, Sjödin P, Lascoux M (2012). Joint analysis of demography
1148 and selection in population genetics: Where do we stand and where could we go? *Mol*
1149 *Ecol* **21**: 28–44.
- 1150 Linnen CR, Poh Y-P, Peterson BK, Barrett RDH, Larson JG, Jensen JD, *et al.* (2013).
1151 Adaptive evolution of multiple traits through multiple mutations at a single gene. *Science*
1152 **339**: 1312–1316.
- 1153 Lischer HEL, Excoffier L (2012). PGDSpider: An automated data conversion tool for
1154 connecting population genetics and genomics programs. *Bioinformatics* **28**: 298–299.
- 1155 Liu L, Yu L (2011). Estimating species trees from unrooted gene trees. *Syst Biol* **60**: 661–667.
- 1156 Lotterhos KE, Whitlock MC (2014). Evaluation of demographic history and neutral
1157 parameterization on the performance of FST outlier tests. *Mol Ecol* **23**: 2178–2192.
- 1158 Lowry DB, Hoban S, Kelley JL, Lotterhos KE, Reed LK, Antolin MF, *et al.* (2016). Breaking
1159 RAD: An evaluation of the utility of restriction site associated DNA sequencing for
1160 genome scans of adaptation. *Mol Ecol Resour.*
- 1161 Lu F, Lipka AE, Glaubitz J, Elshire R, Cherney JH, Casler MD, *et al.* (2013). Switchgrass
1162 Genomic Diversity, Ploidy, and Evolution: Novel Insights from a Network-Based SNP
1163 Discovery Protocol. *PLoS Genet* **9**.
- 1164 Malé PJG, Bardon L, Besnard G, Coissac E, Delsuc F, Engel J, *et al.* (2014). Genome
1165 skimming by shotgun sequencing helps resolve the phylogeny of a pantropical tree
1166 family. *Mol Ecol Resour* **14**: 966–975.
- 1167 Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A, *et al.* (2010). Target-
1168 enrichment strategies for next-generation sequencing. *Nat Methods* **7**: 111–118.
- 1169 Mandoli DF, Olmstead R (2000). The importance of emerging model systems in plant
1170 biology. *J Plant Growth Regul* **19**: 249–252.
- 1171 Manel S, Holderegger R (2013). Ten years of landscape genetics. *Trends Ecol Evol* **28**: 614–
1172 621.
- 1173 Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen W-M (2010). Robust
1174 relationship inference in genome-wide association studies. *Bioinformatics* **26**: 2867–
1175 2873.

- 1176 Martin SH, Van Belleghem SM (2016). Exploring evolutionary relationships across the
1177 genome using topology weighting. *bioRxiv*: 69112.
- 1178 Mazet O, Rodriguez W, Chikhi L (2015). Demographic inference using genetic data from a
1179 single individual: Separating population size variation from population structure. *Theor*
1180 *Popul Biol* **104**: 46–58.
- 1181 McCarroll SA, Sabeti PC, Frazer KA, Varilly P, Fry B, Ballinger DG, *et al.* (2007). Genome-
1182 wide detection and characterization of positive selection in human populations. *Nature*
1183 **449**: 913–8.
- 1184 McDonald JH, Kreitman M (1991). Adaptive protein evolution at the Adh locus in
1185 *Drosophila*. *Nature* **351**: 652–4.
- 1186 McVean G (2007). The structure of linkage disequilibrium around a selective sweep. *Genetics*
1187 **175**: 1395–406.
- 1188 McVean G, Awadalla P, Fearnhead P (2002). A coalescent-based method for detecting and
1189 estimating recombination from gene sequences. *Genetics* **160**: 1231–1241.
- 1190 Messer PW, Petrov DA (2013). Population genomics of rapid adaptation by soft selective
1191 sweeps. *Trends Ecol Evol* **28**: 659–669.
- 1192 Mirarab S, Warnow T (2015). ASTRAL-II: Coalescent-based species tree estimation with
1193 many hundreds of taxa and thousands of genes. *Bioinformatics* **31**: i44–i52.
- 1194 Myers S (2005). A Fine-Scale Map of Recombination Rates and Hotspots Across the Human
1195 Genome. *Science* **310**: 321–324.
- 1196 Nadeau NJ, Whibley A, Jones RT, Davey JW, Dasmahapatra KK, Baxter SW, *et al.* (2012).
1197 Genomic islands of divergence in hybridizing *Heliconius* butterflies identified by large-
1198 scale targeted sequencing. *Philos Trans R Soc Lond B Biol Sci* **367**: 343–53.
- 1199 Nei M, Li WH (1979). Mathematical model for studying genetic variation in terms of
1200 restriction endonucleases. *Proc Natl Acad Sci U S A* **76**: 5269–73.
- 1201 Nicholls JA, Pennington RT, Koenen EJM, Hughes CE, Hearn J, Bunnefeld L, *et al.* (2015).
1202 Using targeted enrichment of nuclear genes to increase phylogenetic resolution in the
1203 neotropical rain forest genus *Inga* (Leguminosae: Mimosoideae). *Front Plant Sci* **6**: 710.
- 1204 Nielsen R, Hellmann I, Hubisz M, Bustamante C, Clark AG (2007). Recent and ongoing

- 1205 selection in the human genome. *Nat Rev Genet* **8**: 857–868.
- 1206 Nielsen R, Williamson S, Kim Y, Nielsen R, Williamson S, Kim Y, *et al.* (2005). Genomic
1207 scans for selective sweeps using SNP data Genomic scans for selective sweeps using
1208 SNP data. *Genome Res*: 1566–1575.
- 1209 Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, *et al.* (2008). Genes mirror
1210 geography within Europe. *Nature* **456**: 98–101.
- 1211 Oikkonen L, Lise S (2017). Making the most of RNA-seq: Pre-processing sequencing data
1212 with Opossum for reliable SNP variant detection. *Wellcome Open Res* **2**: 6.
- 1213 Orozco-terWengel P (2016). The devil is in the details: the effect of population structure on
1214 demographic inference. *Heredity (Edinb)* **116**: 349–350.
- 1215 Palamara PF, Pe'er I (2013). Inference of historical migration rates via haplotype sharing.
1216 *Bioinformatics* **29**: 180–188.
- 1217 Pavlidis P, Laurent S, Stephan W (2010). MsABC: A modification of Hudson's ms to
1218 facilitate multi-locus ABC analysis. *Mol Ecol Resour* **10**: 723–727.
- 1219 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012). Double digest RADseq: an
1220 inexpensive method for de novo SNP discovery and genotyping in model and non-model
1221 species. *PLoS One* **7**: e37135.
- 1222 Pfeifer B, Wittelsburger U, Ramos-Onsins SE, Lercher MJ (2014). PopGenome: An efficient
1223 swiss army knife for population genomic analyses in R. *Mol Biol Evol* **31**: 1929–1936.
- 1224 Pickrell JK, Pritchard JK (2012). Inference of population splits and mixtures from genome-
1225 wide allele frequency data. *PLoS Genet* **8**: e1002967.
- 1226 Piskol R, Ramaswami G, Li JB (2013). Reliable identification of genomic variants from
1227 RNA-seq data. *Am J Hum Genet* **93**: 641–651.
- 1228 Poelstra JW, Vijay N, Bossu CM, Lantz H, Ryll B, Baglione V, *et al.* (2014). The genomic
1229 landscape underlying phenotypic integrity in the face of gene flow in crows. *Science* **344**:
1230 1410–1414.
- 1231 Price A, Patterson NJ, Plenge RM, Weinblatt ME, Shadick N a, Reich D (2006). Principal
1232 components analysis corrects for stratification in genome-wide association studies. *Nat*
1233 *Genet* **38**: 904–9.

- 1234 Pritchard JK, Pickrell JK, Coop G (2010). The Genetics of Human Adaptation: Hard Sweeps,
1235 Soft Sweeps, and Polygenic Adaptation. *Curr Biol* **20**: R208–R215.
- 1236 Pritchard JK, Di Rienzo A (2010). Adaptation – not by sweeps alone. *Nat Rev Genet* **11**: 665–
1237 667.
- 1238 Pritchard JK, Stephens M, Donnelly P (2000). Inference of population structure using
1239 multilocus genotype data. *Genetics* **155**: 945–959.
- 1240 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, *et al.* (2007).
1241 PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage
1242 Analyses. *Am J Hum Genet* **81**: 559–575.
- 1243 Puritz JB, Matz M V., Toonen RJ, Weber JN, Bolnick DI, Bird CE (2014). Demystifying the
1244 RAD fad. *Mol Ecol* **23**: 5937–5942.
- 1245 Raj A, Stephens M, Pritchard JK (2014). FastSTRUCTURE: Variational inference of
1246 population structure in large SNP data sets. *Genetics* **197**: 573–589.
- 1247 Rasmussen MD, Hubisz MJ, Gronau I, Siepel A (2014). Genome-Wide Inference of Ancestral
1248 Recombination Graphs. *PLoS Genet* **10**.
- 1249 Robert CP, Cornuet J-M, Marin J-M, Pillai NS (2011). Lack of confidence in approximate
1250 Bayesian computation model choice. *Proc Natl Acad Sci U S A* **108**: 15112–7.
- 1251 Rockman M V (2012). The QTN program and the alleles that matter for evolution: all that’s
1252 gold does not glitter. *Evolution (N Y)* **66**: 1–17.
- 1253 Roesti M, Kueng B, Moser D, Berner D (2015). The genomics of ecological vicariance in
1254 threespine stickleback fish. *Nat Commun* **6**: 8767.
- 1255 Roux C, Fraisse C, Castric V, Vekemans X, Pogson GH, Bierne N (2014). Can we continue to
1256 neglect genomic variation in introgression rates when inferring the history of speciation?
1257 A case study in a *Mytilus* hybrid zone. *J Evol Biol* **27**: 1662–1675.
- 1258 Roux C, Fraïsse C, Romiguier J, Anciaux Y, Galtier N, Bierne N (2016). Shedding Light on
1259 the Grey Zone of Speciation along a Continuum of Genomic Divergence. *PLOS Biol*.
- 1260 Roux C, Pauwels M, Ruggiero M-V, Charlesworth D, Castric V, Vekemans X (2013). Recent
1261 and ancient signature of balancing selection around the S-locus in *Arabidopsis halleri*
1262 and *A. lyrata*. *Mol Biol Evol* **30**: 435–47.

- 1263 Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, Schaffner SF, *et al.* (2002).
1264 Detecting recent positive selection in the human genome from haplotype structure. **419**.
- 1265 Sattath S, Elyashiv E, Kolodny O, Rinott Y, Sella G (2011). Pervasive adaptive protein
1266 evolution apparent in diversity patterns around amino acid substitutions in *Drosophila*
1267 *simulans*. *PLoS Genet* **7**.
- 1268 Scheet P, Stephens M (2006). A fast and flexible statistical model for large-scale population
1269 genotype data: applications to inferring missing genotypes and haplotypic phase. *Am J*
1270 *Hum Genet* **78**: 629–44.
- 1271 Schiffels S, Durbin R (2014). Inferring human population size and separation history from
1272 multiple genome sequences. *Nat Genet* **46**: 919–25.
- 1273 Schrider DR, Mendes FK, Hahn MW, Kern AD (2015). Soft shoulders ahead: Spurious
1274 signatures of soft and partial selective sweeps result from linked hard sweeps. *Genetics*
1275 **200**: 267–284.
- 1276 Schrider DR, Shanku AG, Kern AD (2016). Effects of linked selective sweeps on
1277 demographic inference and model selection. *Genetics* **204**: 1207–1223.
- 1278 Schubert M, Jónsson H, Chang D, Der Sarkissian C, Ermini L, Ginolhac A, *et al.* (2014).
1279 Prehistoric genomes reveal the genetic foundation and cost of horse domestication. *Proc*
1280 *Natl Acad Sci* **111**: 201416991.
- 1281 Shafer AB a., Wolf JBW, Alves PC, Bergström L, Bruford MW, Brännström I, *et al.* (2015).
1282 Genomics and the challenging translation into conservation practice. *Trends Ecol Evol*
1283 **30**: 78–87.
- 1284 Sheehan S, Harris K, Song YS (2013). Estimating Variable Effective Population Sizes from
1285 Multiple Genomes: A Sequentially Markov Conditional Sampling Distribution
1286 Approach. *Genetics* **194**: 647–662.
- 1287 Song Y, Endepols S, Klemann N, Richter D, Matuschka FR, Shih CH, *et al.* (2011). Adaptive
1288 introgression of anticoagulant rodent poison resistance by hybridization between old
1289 world mice. *Curr Biol* **21**: 1296–1301.
- 1290 Sousa VC, Carneiro M, Ferrand N, Hey J (2013). Identifying loci under selection against gene
1291 flow in isolation-with-migration models. *Genetics* **194**: 211–233.

- 1292 Staab PR, Metzler D (2016). Coala: An R framework for coalescent simulation.
1293 *Bioinformatics* **32**: 1903–1904.
- 1294 Staab PR, Zhu S, Metzler D, Lunter G (2015). Scrm: Efficiently simulating long sequences
1295 using the approximated coalescent with recombination. *Bioinformatics* **31**: 1680–1682.
- 1296 Stamatakis A (2014). RAxML version 8: A tool for phylogenetic analysis and post-analysis of
1297 large phylogenies. *Bioinformatics* **30**: 1312–1313.
- 1298 Stinchcombe JR, Hoekstra HE (2008). Combining population genomics and quantitative
1299 genetics: finding the genes underlying ecologically important traits. *Heredity (Edinb)*
1300 **100**: 158–170.
- 1301 Stucki S, Orozco-Terwengel P, Bruford MW, Colli L, Masembe C, Negrini R, *et al.* (2016).
1302 High performance computation of landscape genomic models integrating local indices of
1303 spatial association. *Mol Ecol Resour*: 1–15.
- 1304 Szpiech ZA, Hernandez RD (2014). selscan: an efficient multithreaded program to perform
1305 EHH-based scans for positive selection. *Mol Biol Evol* **31**: 2824–2827.
- 1306 Tajima F (1989). Statistical method for testing the neutral mutation hypothesis by DNA
1307 polymorphism. *Genetics* **123**: 585–95.
- 1308 Takezaki N, Nei M, Tamura K (2010). POPTREE2: Software for constructing population
1309 trees from allele frequency data and computing other population statistics with windows
1310 interface. *Mol Biol Evol* **27**: 747–752.
- 1311 Tang K, Thornton KR, Stoneking M (2007). A new approach for using genome scans to
1312 detect recent positive selection in the human genome. *PLoS Biol* **5**: 1587–1602.
- 1313 Terhorst J, Kamm JA, Song YS (2016). Robust and scalable inference of population history
1314 from hundreds of unphased whole genomes. *Nat Genet* **49**: 303–309.
- 1315 The Heliconius Genome Consortium, Dasmahapatra KK, Walters JR, Briscoe AD, Davey JW,
1316 Whibley A, *et al.* (2012). Butterfly genome reveals promiscuous exchange of mimicry
1317 adaptations among species. *Nature* **487**: 94–98.
- 1318 Tine M, Kuhl H, Gagnaire P-A, Louro B, Desmarais E, Martins RST, *et al.* (2014). European
1319 sea bass genome and its variation provide insights into adaptation to euryhalinity and
1320 speciation. *Nat Commun* **5**: 5770.

- 1321 Toonen RJ, Puritz JB, Forsman ZH, Whitney JL, Fernandez-Silva I, Andrews KR, *et al.*
1322 (2013). ezRAD: a simplified method for genomic genotyping in non-model organisms.
1323 *PeerJ* **1**: e203.
- 1324 Uricaru R, Rizk G, Lacroix V, Quillery E, Plantard O, Chikhi R, *et al.* (2015). Reference-free
1325 detection of isolated SNPs. *Nucleic Acids Res* **43**: e11.
- 1326 Uy JAC, Cooper EA, Cutie S, Concannon MR, Poelstra JW, Moyle RG, *et al.* (2016).
1327 Mutations in different pigmentation genes are associated with parallel melanism in island
1328 flycatchers. *Proc R Soc B* **283**: 2115–2118.
- 1329 Verhoeven KJF, VonHoldt BM, Sork VL (2016). Epigenetics in ecology and evolution: What
1330 we know and what we need to know. *Mol Ecol* **25**: 1631–1638.
- 1331 Villanueva-Cañas JL, Rech GE, de Cara MAR, González J (2017). Beyond SNPs: how to
1332 detect selection on transposable element insertions. *Methods Ecol Evol* **8**: 728–737.
- 1333 Vitalis R, Gautier M, Dawson KJ, Beaumont MA (2014). Detecting and measuring selection
1334 from gene frequency data. *Genetics* **196**: 799–817.
- 1335 Vitti JJ, Grossman SR, Sabeti PC (2013). Detecting Natural Selection in Genomic Data. *Annu*
1336 *Rev Genet* **47**: 97–120.
- 1337 Wagner CE, Keller I, Wittwer S, Selz OM, Mwaiko S, Greuter L, *et al.* (2013). Genome-wide
1338 RAD sequence data provide unprecedented resolution of species boundaries and
1339 relationships in the Lake Victoria cichlid adaptive radiation. *Mol Ecol* **22**: 787–98.
- 1340 Wang C, Davila JI, Baheti S, Bhagwate A V., Wang X, Kocher JPA, *et al.* (2014). RVboost:
1341 RNA-seq variants prioritization using a boosting method. *Bioinformatics* **30**: 3414–3416.
- 1342 Wang J, Fan C (2014). A neutrality test for detecting selection on DNA methylation using
1343 single methylation polymorphism frequency spectrum. *Genome Biol Evol* **7**: 154–171.
- 1344 Wang M, Huang X, Li R, Xu H, Jin L, He Y (2014). Detecting recent positive selection with
1345 high accuracy and reliability by conditional coalescent tree. *Mol Biol Evol* **31**: 3068–
1346 3080.
- 1347 Wang S, Meyer E, McKay JK, Matz M V (2012). 2b-RAD: a simple and flexible method for
1348 genome-wide genotyping. *Nat Methods* **9**: 808–810.
- 1349 Weber JN, Peterson BK, Hoekstra HE (2013). Discrete genetic modules are responsible for

- 1350 complex burrow evolution in *Peromyscus* mice. *Nature* **493**: 402–5.
- 1351 Wegmann D, Currat M, Excoffier L (2006). Molecular diversity after a range expansion in
1352 heterogeneous environments. *Genetics* **174**: 2009–20.
- 1353 Wegmann D, Leuenberger C, Neuenschwander S, Excoffier L (2010). ABCtoolbox: a
1354 versatile toolkit for approximate Bayesian computations. *BMC Bioinformatics* **11**: 116.
- 1355 Weir BS, Cockerham CC (1984). Estimating F-Statistics for the Analysis of Population
1356 Structure. *Evolution (N Y)* **38**: 1358–1370.
- 1357 White BJ, Cheng C, Simard F, Costantini C, Besansky NJ (2010). Genetic association of
1358 physically unlinked islands of genomic divergence in incipient species of *Anopheles*
1359 *gambiae*. *Mol Ecol* **19**: 925–939.
- 1360 Wolf JBW, Ellegren H (2016). Making sense of genomic islands of differentiation in light of
1361 speciation. *Nat Rev Genet* **18**: 87–100.
- 1362 Yang Z (2007). PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**:
1363 1586–1591.
- 1364 Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS (2012). A high-performance
1365 computing toolset for relatedness and principal component analysis of SNP data.
1366 *Bioinformatics* **28**: 3326–3328.
- 1367 Zhou X, Stephens M (2012). Genome-wide efficient mixed model analysis for association
1368 studies. *Nat Genet* **44**: 821–824.
- 1369
- 1370
- 1371
- 1372
- 1373
- 1374
- 1375
- 1376

1377

Software	Class of method	Purpose	Species	Issues and warnings	Link	Reference
Arlequin	AMOVA (Analysis of Molecular Variance)	Characterizing hierarchical population structure	Arlequin allows for a variety of analyses of diversity (see below)	Requires a priori assignment of individuals to populations, data formatting is required prior analysis	http://cmpg.unibe.ch/software/arlequin35/Arlequin35Downloads.html	(Excoffier and Lischer, 2010)
ADMIXTURE	Clustering and characterizing admixture	Grouping individuals in clusters maximizing HW equilibrium and	Maximum Likelihood, claimed to be faster than Structure	Often slower than its counterparts	https://www.genetics.ucla.edu/software/admixture/index.html	(Alexander and Novembre, 2009)

Tables

Table 1. Summary of methods dedicated to data description and assessing population structure. Methods highlighted in bold can be combined in a pipeline within the R software.

			LD between loci				
FastSTRUCTURE	Clustering and characterizing admixture	Grouping individuals in clusters maximizing HW equilibrium and LD between loci	~100X faster than Structure	Approximate inference of the original Structure model	http://rajanil.github.io/fastStructure/	(Raj <i>et al.</i> , 2014)	
FineStructure/Globetrotter	Clustering and characterizing admixture	Chromosome painting, admixture and clustering	Estimates time since admixture, fast, specific tools for RAD-seq, set of scripts to facilitate analysis	Relies on Structure and fastStructure assumptions. Requires phased data.	http://paintmychromosomes.com/	(Hellenthal <i>et al.</i> , 2014)	
GENELAND	Clustering and characterizing admixture	Grouping individuals in spatially consistent clusters maximizing HW equilibrium	Takes into account spatial variation, supposed to detect weak structure, framed in R	Immigrant alleles are assumed to be found only in new immigrants	https://cran.r-project.org/web/packages/Geneland/	(Guillot <i>et al.</i> , 2012)	
PCAdmix	Clustering and characterizing admixture	Chromosome painting	Fast, uses HMM to smooth out windows and limit noise due to low confidence ancestry	Requires a priori definition of ancestral populations and phased haplotypes	https://sites.google.com/site/pcadmix/	(Brisbin <i>et al.</i> , 2012)	
sNMF	Clustering and characterizing admixture	Grouping individuals in clusters maximizing HW equilibrium and LD between loci	Fast (30X than ADMIXTURE)	Still slow computation time for large datasets	http://membres-timc.imag.fr/Olivier.Francois/snmf/index.htm	(Frichot <i>et al.</i> , 2014)	
STRUCTURE	Clustering and characterizing admixture	Grouping individuals in clusters maximizing HW equilibrium and LD between loci	User friendly interface. Bayesian inference.	Slow for large datasets. Requires specific input format	http://pritchardlab.stanford.edu/structure.html	(Pritchard <i>et al.</i> , 2000)	
TREEMIX	Clustering and characterizing admixture	Admixture graph, infers most likely admixture events in a tree	Based on allele frequencies and can be used for pooled data.	Requires multiple runs to properly assess the likelihood of each model	https://bitbucket.org/nygcresearch/treemix/src	(Pickrell and Pritchard, 2012)	
BEDASSLE	Differentiation and MCMC model testing	Identifies contribution of environment and geographical distance to populations differentiation	Less biased than Mantel tests, provides tools for model testing	Uses population-level data.	https://cran.r-project.org/web/packages/BEDASSLE/index.html	(Bradburd <i>et al.</i> , 2013)	
npstat	Differentiation/Diversity	Extracting summary statistics from pooled data	Explicitly corrects for sampling bias in pooled data. Allows computing tests using an outgroup (MK)	Mostly limited to summary statistics, but more complete than Popoolation.	https://github.com/lucaferretti/npstat	(Ferretti <i>et al.</i> , 2013)	

Stacks	Differentiation/Diversity/Phylogeny	Processing RAD-seq data and facilitate their analysis	test, Fay and Wu's H) and characterizing coding mutations. Designed for RAD-seq data, variety of output formats for downstream analyses. Allows to retrieve DNA sequences for each locus	NA	http://catchenlab.life.illinois.edu/stacks/	(Catchen <i>et al.</i> , 2011)
ANGSD	Differentiation/Diversity/Recombination	Computing summary statistics based on AFS and LD along genomes	Able to process BAM files, built-in procedures for data filtering, admixture analysis	Mostly limited to summary statistics	https://github.com/ANGSD/angsd	(Kornelussen <i>et al.</i> , 2014)
Arlequin	Differentiation/Diversity/Recombination	Computing summary statistics based on AFS and LD along genomes	Can output AFS for further analysis in fastsimcoal2	Slower than PopGenome, requires a private format	http://cmpg.unibe.ch/software/arlequin35/Ar135Downloads.html	(Excoffier and Lischer, 2010)
POPGenome	Differentiation/Diversity/Recombination	Computing summary statistics based on AFS and LD along genomes	Accepts VCF and GFF/GFT files, efficient and fast. Tests for admixture available (ABBA BABA test). Includes basic coalescence simulations (ms and msms)	Mostly limited to summary statistics (but coalescent simulations are possible). No built-in SNP calling module	http://catchenlab.life.illinois.edu/stacks/	(Pfeifer <i>et al.</i> , 2014)
Popoolation2/Popoolation TE	Differentiation/Diversity/Recombination	Extracting summary statistics from pooled data	Explicitly corrects for sampling bias in pooled data	Mostly limited to a few summary statistics. A pipeline dedicated to TE detection is also available	https://sourceforge.net/p/popoolation/wiki/Main/	(Kofler, Orozco-terWengel, <i>et al.</i> , 2011; Kofler, Pandey, <i>et al.</i> , 2011)
VCFTOOLS	Differentiation/Diversity/Recombination	Computing summary statistics based on AFS and LD along genomes	Fast. VCFTOOLS can also be used for SNP filtering	Less summary statistics than POPGenome	https://vcftools.github.io/man_latest.html	(Danecek <i>et al.</i> , 2011)
POPTREE2	Genetic differentiation	Visualizing a matrix of pairwise differentiation statistics as a tree	Can be used for pooled datasets, several statistics can be used	Differentiation measures alone do not necessarily retrieve the actual history of populations	http://www.med.kagawa-u.ac.jp/~genomelb/takezaki/poptree2/index.html	(Takezaki <i>et al.</i> , 2010)
Kimtree	Genetic distance	Estimating divergence time between populations	The method is conditional on a prior topology provided by	Times are given in diffusion time scale, and can be converted	http://www1.montpellier.inra.fr/CBGP/software/kimtree/index	(Gautier and Vitalis, 2013)

		and testing for topologies	the user. It computes DIC for a given topology, allowing to test for the best one.	in demographic times using independent estimates of N_e .	html	
Eigenstrat/smartpca	Multivariate analysis	Summarizing variance across loci and visualizing inter-individual genetic distance	Fast. Can use VCF files as an input	Requires careful interpretation (Jombard et al. 2009)	https://github.com/DReichLab/ElG/tree/master/EIGENSTRAT	(Price <i>et al.</i> , 2006)
SPRelate	Multivariate analysis	Summarizing variance across loci and visualizing inter-individual genetic distance	Fast. Can use VCF files as an input	Requires careful interpretation (Jombard et al. 2009)	https://bioconductor.org/packages/release/bioc/html/SNPRelate.html	(Zheng <i>et al.</i> , 2012)
DAPC (adegenet)	Multivariate analysis/Clustering	Maximizes divergence between groups identified by PCA	Fast. Less sensitive to HWE assumptions. Claims to be more efficient than Structure	Requires careful interpretation (Jombard et al. 2009)	http://adegenet.r-forge.r-project.org/	(Jombard <i>et al.</i> , 2010)
sPCA (adegenet)	Multivariate analysis/Clustering	Spatially explicit model to assess population structure	Spatially explicit and able to detect cryptic structure. Fast. Mendelian error checking, testing family structure, highly accurate kinship coefficient, association analysis, population structure inference	Does not take into account HW equilibrium or LD	http://adegenet.r-forge.r-project.org/	(Jombard <i>et al.</i> , 2008)
KING	Pedigree, Identity by descent/state	Estimating inbreeding and relatedness, multivariate analysis	LAMP also allows for association and pedigree analyses	Kinship coefficient also computed in VCFTOOLS	http://people.virginia.edu/~wc9c/KING/Download.htm	(Manichaikul <i>et al.</i> , 2010)
LAMP	Pedigree, Identity by descent/state	Chromosome painting, relatedness	Allows studying identity by descent and by state. PLINK is a multi-purpose tool, facilitating data analysis within the same software	Identifies local ancestry in windows (source of noise), requires phased data	http://lamp.icsi.berkeley.edu/lamp/	(Baran <i>et al.</i> , 2012)
PLINK	Pedigree, Identity by descent/state	Estimating inbreeding and relatedness	Computes unadjusted Ajk and kinship coefficient	NA	http://pngu.mgh.harvard.edu/~purcell/plink/	(Purcell <i>et al.</i> , 2007)
VCFTOOLS	Pedigree, Identity by descent/state	Estimating inbreeding and relatedness	Coalescence-based. Suitable for short loci (e.g. RAD-seq and	NA	https://vcftools.github.io/man_latest.html	(Danecek <i>et al.</i> , 2011)
ASTRAL-2	Phylogeny	Builds species trees using short non-recombining	Coalescence-based. Suitable for short loci (e.g. RAD-seq and	More reliable under high incomplete lineage sorting that	https://github.com/mirarab/ASTRAL	(Mirarab and Warnow, 2015)

		sequences	GBS)	SVDQuartets and NJst (Chou <i>et al.</i> 2015)		
BEAST2	Phylogeny	Network reconstruction and phylogenetic relationships	User friendly. Can be used to track changes in effective population sizes (Bayesian Skyline Plots). Possible to estimate divergence times	Slow for large datasets. Requires sequence data that can be produced by , e.g., Stacks for RAD-seq data	http://beast2.org/	(Drummond and Rambaut, 2007; Bouckaert <i>et al.</i> , 2014)
NJst (in phybase)	Phylogeny	Builds species trees using short non-recombining sequences	Coalescence-based. Suitable for short loci (e.g. RAD-seq and GBS)	See ASTRAL-2 and Chou <i>et al.</i> 2015	https://code.google.com/archive/p/phybase/downloads	(Liu and Yu, 2011)
PhyML	Phylogeny	Phylogenetic relationships	Maximum Likelihood inference of phylogenetic relationships. An online version is available	Should be used on complex of species or divergent populations with little migration	http://www.atgc-montpellier.fr/phyml/binaries.php	(Guindon <i>et al.</i> , 2010)
RxML	Phylogeny	Network reconstruction and phylogenetic relationships	Maximum Likelihood inference of phylogenetic relationships	Should be used on complex of species or divergent populations with little migration	http://sco.hits.org/exelixis/web/software/raxml/index.html	(Stamatakis, 2014)
SNAPP	Phylogeny	Phylogenetic relationships	Handles SNP data	Remains slow for medium to large datasets (>1,000SNPs)	http://beast2.org/snapp/	(Bryant <i>et al.</i> , 2012)
SNPhylo	Phylogeny	Network reconstruction and phylogenetic relationships	Complete pipeline from SNP filtering to tree reconstruction	Should be used on complex of species or divergent populations with little migration	http://chibba.pgml.uga.edu/snphylo/	(Lee <i>et al.</i> , 2014)
SVDQuartets	Phylogeny	Builds species trees using short non-recombining sequences	Coalescence-based. Suitable for short loci (e.g. RAD-seq and GBS)	See ASTRAL-2 and Chou <i>et al.</i> 2015	http://www.stat.osu.edu/~kubatko/software/SVDquartets/	(Chifman and Kubatko, 2014)
*BEAST	Phylogeny and species tree inference	Divergence time estimation and phylogenetic relationships	Outputs a species tree instead of concatenated gene tree. Allows for testing consistency between phylogenetic signals at different loci	Slow for large datasets. Requires sequence data. Not suited for situations where gene flow/admixture occurs	http://beast2.org/	(Heled and Drummond, 2010)

Splitstree	Phylogeny/Network	Network reconstruction and phylogenetic relationships	User friendly interface, proposes a variety of methods for networks reconstruction	Mostly descriptive	http://www.splitstree.org/	(Huson and Bryant, 2006)
LDHat	Recombination	Estimating variation in recombination rates along a genome	Handles unphased and missing data, underlying model can be used for organisms such as viruses or bacteria	Limited to 300 sequences, private format, model for recombination hotspots based on human data	http://ldhat.sourceforge.net/	(McVean <i>et al.</i> , 2002)
LDHot	Recombination	Identifying recombination hotspots	Specifically designed for detecting recombination hotspots	Requires data to be phased, working with LDHat	https://github.com/auton1/LDhot	(Myers, 2005)
TWISST	Topology weighting	Chromosome painting, clustering and branching between populations	Retrieves the most likely coalescence pattern between several taxa along the genome. Can be seen as an extension of the ABBA/BABA test	Needs a priori grouping of individuals into taxa. Requires at least 4 taxa. Impractical for more than 6 taxa. Windows size must include enough SNPs to retrieve the correct topology but at the risk that regions with different histories are included	https://github.com/simonmartin/twisst	(Martin and Van Belleghem, 2016)
BAYPASS/Bayenv	Variance/covariance matrix	Building a population covariance matrix across population allele frequencies, similar to TREEMIX	Can handle pooled data	Matrices are mostly designed to provide a neutral model for assessing selection, but can be used to infer population structure	http://www1.montpellier.inra.fr/CBGP/software/baypass/ ; https://bitbucket.org/tguenther/bayenv2/public/src	(Günther and Coop, 2013; Gautier, 2015)

Software	Class of method	Purpose	Specifities	Issues and warnings	Link	Reference
abc	ABC	Performs all steps for model-checking and parameters estimation for ABC analyses	Informative vignette, allows graphical representation, complete and robust	Does not perform coalescent simulations (but can be used in combination with coala)	https://cran.r-project.org/web/packages/abc/index.html	(Csilléry <i>et al.</i> , 2012)

Table 2. Summary of methods for demographic inference, simulations and scenarios comparisons. Methods available in R are highlighted in bold.

ABCToolbox	ABC	Complete ABC analysis, from simulations to model checking and parameters estimation	Modular, facilitates the computation of summary statistics	Current version is Beta (15/01/2016)	https://bitbucket.org/phaentu/abctoolbox-public/	(Wegmann <i>et al.</i> , 2010)
DIYABC	ABC	Complete ABC analysis, from simulations to model checking and parameters estimation	User-friendly	Does not allow to model continuous gene flow	http://www1.montpellier.inra.fr/CBGP/diyabc/	(Cornuet <i>et al.</i> , 2008)
PopSizeABC	ABC	Inferring change in N_e using whole-genome data	Supposed to better assess recent events. Uses a set of summary statistics for the AFS and LD between markers. Handles multiple individuals	Approximate bayesian approaches do not retrieve the whole information	https://forge-dga.jouy.inra.fr/projects/popsizabc/	(Boistard <i>et al.</i> , 2016)
coala	ABC/coalescent simulations	Combining coalescent simulators within a single framework	Facilitates the building of scenarios and computes summary statistics for simulations Performs coalescent simulations, parameter estimation and model testing using a fast likelihood method. Can handle arbitrarily complex scenarios for any type of marker Provides quantitative estimates for TMRCA and topologies at each locus.	Includes so far ms, msms and scrm	https://cran.r-project.org/web/packages/coala/index.html	(Staab and Metzler, 2016)
fastsimcoal2	ABC/Likelihood	Model comparison and parameters estimation	Model comparison and parameters estimation Retracing the whole process of recombination and coalescence along a genome	Summary statistics need to be calculated through Arlequin, slowing down their computation	http://cmpg.unibe.ch/software/fastsimcoal2/	(Excoffier <i>et al.</i> , 2013)
ARGWeaver	Ancestral Recombination Graphs/coalescence	Retracing the whole process of recombination and coalescence along a genome	Estimates effective population size. Provides tools to extract summary statistics for the topologies retrieved.	High computing cost. Requires phased whole-genome data.	https://github.com/mdrasmus/argweaver	(Rasmussen <i>et al.</i> , 2014)
G-PhoCS	Bayesian	Estimating population divergence and migration parameters using a coalescent framework	Bayesian + MCMC, handles ancient samples	Parameters scaled by mutation rate, no admixture	http://compgen.cshl.edu/GPhoCS/	(Gronau <i>et al.</i> , 2011)
IMa2	Bayesian	Inferring parameters from an isolation with migration model	Fully bayesian approach, can perform joint estimates of parameters in L-mode and test for nested models	IM model is the only one available. Discrete admixture cannot be tested. Long computation times. Recent splits lead to overestimate migration rates	https://bio.cst.temple.edu/~hey/software/software.htm#IMa2	(Hey and Nielsen, 2007)

Migrate-n	Bayesian	Inferring migration rates	Both ML and bayesian methods can be used to estimate parameters	Only estimates population sizes and migration rates. Not suited for large datasets. Private input format	http://popgen.sc.fsu.edu/Migrate/Migrate-n.html	(Beerli and Palczewski, 2010)
ABLE	Coalescence/Composite Likelihood	Model comparison and parameters estimation	Uses both allele frequency spectrum and linkage disequilibrium within blocks of a pre-specified size. Handles whole-genome data and RAD-seq.	Relies on ms syntax. Determining the most informative size for blocks requires performing pilot runs.	https://github.com/champost/ABLE	(Beeravolu <i>et al.</i> , 2016))
fastsimcoal2	coalescent simulations	Building any arbitrary scenario using a coalescent framework	Any arbitrary scenario can be implemented. Handles SNP, microsatellites and sequence data.	Does not handle selection. Slower than ms with no recombination, much faster with recombination (see manual)	http://cmpg.unibe.ch/software/fastsimcoal2/	(Excoffier and Foll, 2011)
ms, msms, msABC	coalescent simulations	Building any arbitrary scenario using a coalescent framework	Any arbitrary scenario can be implemented. Handles SNP, microsatellites and sequence data. msms can include selection in the model.	Can be difficult to handle for the naive user (but see coala)	http://www.bio.lmu.de/~pavlidis/home/?Software:msABC	(Hudson, 2002; Ewing and Hermisson, 2010; Pavlidis <i>et al.</i> , 2010)
scrm	coalescent simulations	Fast simulation of chromosome-scale sequences	Syntax similar to ms, handles any arbitrary scenario	Does not handle gene conversion and fixed number of segregating sites (unlike ms)	https://scrm.github.io/	(Staab <i>et al.</i> , 2015)
$\partial a \partial i$	Diffusion approximation of the AFS	Model comparison and parameters estimation	Run time does not depend on the number of SNPs included, does not require coalescent simulations, handles arbitrarily complex scenarios	Requires some knowledge of Python. Limited to 3 populations	https://bitbucket.org/guttenkunstlab/dadi	(Gutenkunst <i>et al.</i> , 2009)
DoRIS	IBD tract	Testing various demographic scenario	Uses variation in IBD tracts length to test for various demographic models.	IBD must be inferred first with, e.g., BEAGLE. Handles a limited set of demographic scenarios. Modification in the code is required for more complex scenarios	https://github.com/pierpal/DoRIS	(Palamara and Pe'er, 2013)

Unnamed	Identity by state tract	Predict observed patterns of Identity by state along a genome by fitting an appropriate, arbitrary complex demographic model	Allows bootstrapping and estimating confidence over parameter estimates with ms	Specific input format (similar to MSMC or ARGWeaver)	https://github.com/kellyharris/Inferring-demography-from-IBS	(Harris and Nielsen, 2013)
diCal2	Sequentially Markovian coalescent	Testing any arbitrary demographic scenario	Works with smaller, more fragmented datasets than PSMC. Handles more complex demographic models than MSMC (including admixture). Allows to track population size changes in time without a priori.	Requires phased whole genome data and a model to be defined	https://sourceforge.net/projects/dical2/	(Sheehan <i>et al.</i> , 2013)
MSMC	Sequentially Markovian coalescent	Inferring change in N_e and migration rates with time between two populations	Allows estimating variation in cross-coalescence rate between two populations	Limited to the study of 8 diploid individuals from 2 populations at once. Requires whole genome phased data and masking regions with insufficient sequencing depth	https://github.com/stschiff/msmc	(Schiffels and Durbin, 2014)
PSMC	Sequentially Markovian coalescent	Inferring change in N_e with time using a single diploid genome	Allows to track population size changes in time without a priori.	Limited to one population and one diploid individual. Better used within MSMC. Requires phased whole genome data and masking regions with insufficient sequencing depth	https://github.com/lh3/psmc	(Li and Durbin, 2011)
SMC++	Sequentially Markovian coalescent	Inferring change in N_e with time and splitting time between two populations	Can analyze hundreds of individuals at a time and does not require phasing	The ancestral allele is assumed to be the reference allele by default. Assumes a clean split for populations divergence	https://github.com/ppgenmethods/smcpp	(Terhorst <i>et al.</i> , 2016)

Table 3. Summary of common methods for identifying loci under selection. Methods available in R are highlighted in bold.

Software	Class of method	Purpose	Specifics	Issues and warnings	Link	Reference
----------	-----------------	---------	-----------	---------------------	------	-----------

ARGWeaver	Ancestral recombination graphs	Detecting selection by screening for variation in topology and age of alleles	Provides quantitative estimates for TMRCA and topologies at each locus. Can be used to infer demographic history. Especially useful to identify signature of long-term balancing selection (older coalescence times)	High computing cost. Requires phased whole-genome data.	https://github.com/mdrasmus/argweaver	(Rasmussen <i>et al.</i> , 2014)
GEMMA	Association	Detecting association with environmental/phenotypic features	Computationally efficient for large scale datasets	Imports data from PLINK format	http://www.xzlab.org/software.html	(Zhou and Stephens, 2012)
GENABEL	Association	Detecting association with environmental/phenotypic features	Modularity, facilitates correction for population structure/relatedness.	Imports data from PLINK format	http://www.genabel.org/	(Aulchenko <i>et al.</i> , 2007)
PLINK	Association	Detecting association with environmental/phenotypic features	Handles a variety of tests for population structure and relatedness	Population structure/kinship need to be assessed prior association analysis	http://pngu.mgh.harvard.edu/~purcell/plink/	(Purcell <i>et al.</i> , 2007)
Trinuclo	Association	Detecting association with environmental/phenotypic features	Specifically designed to handle categorical variables with more than 2 categories. Performs multinomial logistic regression and provides frequentist and bayesian frameworks.	Requires lapack library in Unix. Allows fine-mapping by testing for correlations between adjacent markers.	https://sourceforge.net/projects/trinuclo/	(Jostins and McVean, 2016)
SAMBADA	Association/Environmental association	Detecting association with environmental/phenotypic features	Designed to be fast, underlying models have been kept simple. Allows conversion from PLINK format. Takes into account spatial autocorrelation of individual genotypes. Allows correction for population structure	Does not work with pooled data. Possibly high levels of false positives. Relatedness between samples should be assessed independently. Should be used in combination with LFMM or BayPass.	http://asig.epfl.ch/sambada	(Stucki <i>et al.</i> , 2016)
discoal	Coalescence	Simulate selective sweeps under arbitrary demographic scenarios	More specifically designed for studying soft and hard sweeps	Redundant with msms	https://github.com/kernlab/discoal	Publication embargoed (Kern and Schrider, 2016)
msms	Coalescence	Simulate demographic scenarios including selection	Flexible, syntax similar to ms, handles arbitrarily complex models. Can be used in an ABC	Syntax can be difficult to handle for the naive user (but see coala)	http://www.mabs.at/ewing/msms/index.shtml	(Ewing and Hermisson, 2010)

			framework to include selection as a parameter to be estimated			
diCal-IBD	Coalescent with recombination/IBD	Predicting IBD tracts from demographic models	High IBD sharing suggests recent positive selection.	Uses diCal output to obtain expectations based on demographic scenarios	https://sourceforge.net/projects/dical-ibd/	https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4296155/
SweeD	Composite Likelihood test	Designed for whole genome data (or large continuous regions)	Supports Fasta and VCF formats. Estimates for selection coefficients. Designed for detecting recent positive selection.	Better suited for whole genome datasets	http://pop-gen.eu/wordpress/software/sweed	(Degiorgio <i>et al.</i> , 2016)
SCCT	Conditional coalescent tree	Detecting positive selection	Claims to be more precise at identifying selected sites	Requires whole-genome data. The ancestral state of alleles must be obtained through an outgroup	https://github.com/wavefancy/scct	(Wang, Huang, <i>et al.</i> , 2014)
LFMM	Environmental association	Detecting adaptation to environmental features	Corrects for population structure using latent factors, faster than BAYENV for large datasets	Only performs association with environment	http://membres-timc.imag.fr/Olivier.Francois/lfmm/software.htm	(Frichot <i>et al.</i> , 2013)
H12 test	LD	Detecting selection using signatures of high LD	Does not require phased data. Designed for detecting soft sweeps	Coalescent simulations are recommended to evaluate the likelihood of selection	https://github.com/ngarud/SelectionHapStats/	(Garud <i>et al.</i> , 2015)
LDna	LD	Detecting selection using signatures of high LD	Can be used to address population structure or detect large inversions or indel polymorphism through LD	The user needs to play with parameters to ensure robustness of SNPs significantly linked	https://github.com/petrikemppainen/LDna	(Kemppainen <i>et al.</i> , 2015)
rehh	LD	Detecting selection using signatures of high LD	Can compute both XP-EHH and Rsb. Handles several input formats	Requires phased data and high density of markers	https://cran.r-project.org/web/packages/rehh/index.html	(Gautier and Vitalis, 2012)
Selscan	LD	Detecting selection using signatures of high LD	Includes the nSL statistics dedicated to soft sweep detection	Does not include utilities to specify the ancestral state of alleles. Requires phased data and high density of markers	https://github.com/szpiech/selscan	(Szpiech and Hernandez, 2014)
BALLET	Likelihood test for balancing selection	Detecting balancing selection	Designed for detecting ancient balancing selection. Does not require phasing	Requires whole-genome data and recombination map. The ancestral state of alleles must be obtained through an outgroup	http://www.personal.psu.edu/mxd60/ballet.html	(DeGiorgio <i>et al.</i> , 2014)

Bayescan	Population differentiation	Detecting positive selection and local adaptation	Incorporates uncertainty on allele frequencies due to low sample sizes	Sensitive to priors on the ratio of selected/neutral sites. False positive rates can be high under scenarios of demographic expansion, admixture and isolation by distance	http://cmpg.unibe.ch/software/BayScan/	(Foll and Gaggiotti, 2008)
FDIST2	Population differentiation	Detecting positive selection and local adaptation	Allows to control for hierarchical population structure	False positive rate is high when an island model cannot be assumed	http://datadryad.org/resource/doi:10.5061/dryad.v8d05	(Beaumont and Balding, 2004)
PCAdapt	Population differentiation	Detecting positive selection and local adaptation	Does not require to define populations. Handles admixed populations and pooled datasets Can estimate the coefficients of selection. Calibration using a pseudo-observed dataset to obtain (can be used in combination with the R function <code>simulate.baypass()</code> in BayPass).	False positive rate can be high	http://membres-timc.imag.fr/Michael.Blum/PCAdapt.html	(Duforet-Frebourg <i>et al.</i> , 2016)
SelEstim	Population differentiation	Detecting positive selection and local adaptation	pseudo-observed dataset to obtain (can be used in combination with the R function <code>simulate.baypass()</code> in BayPass).	Assumes an island model.	http://www1.montpellier.inra.fr/CBGP/software/selection/elestim/	(Vitalis <i>et al.</i> , 2014)
Bayenv, BayPass	Population differentiation/Association	Detecting positive selection and adaptation to environmental features	Less sensitive to population demographic history than previous methods. Handle pooled datasets	Significance thresholds need to be determined from pseudo-observed datasets. Calibration with neutral SNPs is recommended. BayPass better estimates the kinship matrix	http://www1.montpellier.inra.fr/CBGP/software/baypass/ ; https://bitbucket.org/tgumenther/bayenv2_public/src	(Günther and Coop, 2013; Gautier, 2015)
FLK	Population differentiation/Association	Detecting positive selection and local adaptation	Less sensitive to population demographic history than previous methods	Requires an outgroup population	https://qgsp.jouy.inra.fr/index.php?option=com_content&view=article&id=50&Itemid=55	(Bonhomme <i>et al.</i> , 2010)
POP BAM	Summary statistics	Detecting selection using AFS, differentiation	Extracts summary statistics directly from BAM files	Does not allow for sophisticated filtering and SNP calling	http://popbam.sourceforge.net/	(Garrigan, 2013)
POPGenome	Summary statistics	Detecting selection using AFS, differentiation	Fast, embedded in R, allows using annotation files (GFF/GTF format).	Does not perform association, but can be used in combination with GENABEL within R	https://cran.r-project.org/web/packages/PopGenome/index.html	(Pfeifer <i>et al.</i> , 2014)

VCFTOOLS	Summary statistics	Detecting selection using AFS, differentiation	Extracts summary statistics from VCF files. Also allows VCF filtering and conversion	Set of summary statistics not as extensive as PopGenome	http://vcftools.sourceforge.net/	(Danecek <i>et al.</i> , 2011)
ANGSD	Summary statistics/Association	Detecting selection using AFS, differentiation, association with functional traits	Allows for association using generalized linear models	Descriptive statistics. P-values need to be evaluated through coalescent simulations.	http://www.popgen.dk/angsd/index.php/ANGSD	(Korneliussen <i>et al.</i> , 2014)
TASSEL	Summary statistics/Association	Detecting association with phenotype	User friendly (Java interface), corrects for relatedness, allows computing summary statistics (LD, diversity)	Requires relatedness to be assessed externally (with e.g. STRUCTURE)	http://www.maizegenetics.net/tassel	(Korneliussen <i>et al.</i> , 2014)
selectionTools	Summary statistics/LD	Detecting selection using AFS, differentiation and LD statistics	Allows combining several tools in a single pipeline. Includes phasing tools.	Set of available summary statistics remains limited (same as VCFtools + Fay and Wu's H)	https://github.com/MerrimanLab/selectionTools	(Cadzow <i>et al.</i> , 2014)

Figures

Figure 1. A possible general pipeline for analysing population genomics data using methods described in this paper. In red are indicated options that are generally not suited for pool-seq data. In green are indicated steps that require genome-wide datasets. ARG: Ancestral Recombination Graph (see main text).

Figure 2. Set of questions and relevant methods to characterize population structure and local adaptation. Proposed methods mostly use common data formats for input files, facilitating their integration in a single pipeline. PGDSpider (Lischer and Excoffier, 2012) can be used to automate file conversion for methods requiring private input format. The proposed methods are not exhaustive, see tables for a more detailed list.

Figure 3: Set of questions and relevant methods to characterize demography and selection. *: requires reference genome; ** requires reference genome and whole genome resequencing.

Sequence data: **whole genome**, SNP-array, RNAsequencing, target enrichment, RAD-seq

Is there a closely related reference genome? No → **Assembly by locus then BLAST on the closest reference (optional)**

Yes
Align on reference

Many tools assume diploid species.

Extract **individual genotypes** or allele frequencies (Pool-sequencing)

Checking population structure and admixture

Selection, recombination, admixture, demography all impact polymorphism

Checking variation in polymorphism genome-wide (e.g. Fst, LD, Tajima's D, ABBA/BABA tests)

Is there more than a single admixed population?

Yes

No

Detect selection using variation in genotypes and allele frequencies across populations and individuals (Fst outliers methods, association)

Detect selection using variation in genotypes across individuals (association)

Mutation rate, recombination rate, generation time are needed to recover demographic estimates for parameters. Most methods for whole genomes are limited by computation time and do not always allow for complex models

Demographic inference (Sequential Markov Coalescent)
Quantify local admixture, introgression, selection: ARG

Background selection can produce higher Fst, like positive selection. Check diversity in each population and their connectivity.

Is there a high density of markers replaced along the genome?

Yes

Perform LD-based tests (soft and hard sweeps)

Cross the results of all analyses
Which regions are under balancing selection?
Is there any sign of adaptive introgression?
Which set of SNPs is found in all or most analyses?
How does demographic history affect selection?

Identify a set of neutral markers and use it for model testing and estimate demographic history. Consider allowing for various introgression rates between populations.

Question	Data format	Software
<p>Is the dataset structured ? Is there inbreeding ? Characterizing hybridization</p>	FASTQ (RAD-sequencing)	STACKS. Outputs F-statistics and estimators of effective population size (π).
	VCF	VCFTOOLS/POPGENOME: relatedness between individuals, FST between populations, Hardy-Weinberg equilibrium Nucleotide diversity and estimates of effective population sizes. Signatures of population size change. PCA methods (SNPRelate in R)
	PLINK PED/BED file	fastSTRUCTURE, sNMF, ADMIXTURE. Provide coefficients of coancestry for each individual. Familial relationships: KING, PLINK.
	Private format (convertible from VCF with PGDspider)	Arlequin/Genepop: testing hierarchical structure of populations (AMOVA), FIS, FST.
<p>How does environment impact this structure and historical dispersal ?</p>	VCF, PLINK PED/BED format	adegenet and LEA packages in R. Highlight barriers to gene flow in the landscape.
	Private format	BEDASSLE. Identifies environmental features limiting gene flow. GENELAND. Highlight barriers to gene flow in the landscape.
<p>Is there any association of specific loci with environment / a relevant phenotype ?</p>	VCF, PLINK PED/BED format	GENABEL, TRINCULO (individual phenotypes) LFMM (LEA package)
	PLINK PED/BED format	SAMBADA
	Private format (convertible from VCF with PGDspider)	BAYENV

informs



Question	Data format	Software
<p>How does selection shape genome variation?</p>	VCF	VCFTOOLS, POPGENOME. Output diversity and LD statistics
	Modified IMPUTE format (can be obtained from VCFTOOLS)	R package rehh. LD-based tests of haplotype extension *
	Private format (PGDSpider from VCF)	BAYENV. FST-outlier method
	VCF and PLINK PED/BED files	PCAdapt. List loci atypically related to population structure
	Private format	ARGWeaver. Returns coalescence times and other statistics for non-recombining blocks along the genome. **
<p>Does population history shape potential for adaptation (e.g. admixture bringing new alleles, bottleneck reducing genetic diversity) ?</p>	PLINK PED/BED file	TreeMix. Identifies admixture events, their magnitude and direction between populations
	Private format (Arlequin, PGDSpider from VCF)	Fastsimcoal and R package abc. ABC and Likelihood methods for comparing arbitrarily complex demographic models
	VCF	SMC methods. SMC++ (no phasing). ** Variation in effective population sizes and divergence times between populations
<p>Origin of genomic islands of differentiation. Characterizing adaptive introgression.</p>	VCF	VCFTOOLS, POPGENOME. Output divergence statistics.
	PLINK	FineStructure. Identifies introgressed blocks along the genome and estimates times since admixture *
	BEAGLE (after phasing from VCF)	PCAdmix. Identifies introgressed blocks along the genome *
	Private format	ARGWeaver. Returns coalescence times and other statistics for non-recombining blocks along the genome. **

informs