1   # Going down the rabbit hole: a review on how to link genome-wide data with

2   ecology and evolution in natural populations

3   Yann X.C. Bourgeois

4

5   ## Abstract

6   Characterizing species history and assessing the nature and extent of local adaptation is

7   crucial in conservation, agronomy, functional ecology and evolutionary biology. The ongoing

8   and constant improvement of next-generation sequencing (NGS) techniques has facilitated the

9   production of an increasingly growing amount of genetic markers across genomes of non-

10  model species. The study of variation at these markers across natural populations has

11  deepened the understanding of how population history and selection act on genomes.

12  However, this improvement has come with a burst of analytical tools that can confuse naïve

13  users. This confusion can limit the amount of information effectively retrieved from complex

14  genomic datasets. In addition, the lack of a unified analytical pipeline impairs the diffusion of

15  the most recent analytical tools into fields like conservation biology. This requires efforts be

16  made in providing introduction to these methods. In this paper I describe possible analytical

17  protocols and recent methods dealing with analysis of genome-scale datasets, clarify the

18  strategy they use to infer demographic history and selection, and discuss some of their

19  limitations.

20

21

22

23

24

25

26

27

## 28    Introduction

29    Genetic makeup of populations is shaped by multiple historical and selective factors. The

30    advent of Next-Generation Sequencing (NGS) in the last 20 years has enhanced our

31    understanding on how intermingled these factors are, and how they can impact genomic

32    variation. Important results have been gathered on model species, or species with an

33    economical interest. Such results include, among other examples, an improved perspective on

34    human history of migrations, admixture and adaptation (e.g. Sabeti *et al.*, 2002; Abi-Rached

35    *et al.*, 2011; Li and Durbin, 2011), elucidating the origin of domesticated species (e.g.

36    Axelsson *et al.*, 2013; Schubert *et al.*, 2014), or characterizing the genetic bases of local

37    adaptation in model or near-model species (e.g. Legrand *et al.*, 2009; Kolaczkowski *et al.*,

38    2011; Roux *et al.*, 2013; Kubota *et al.*, 2015). These studies have brought insights at an

39    unprecedented scale on the links between genotype, phenotype and environment. Most of

40    these studies relied on a precise knowledge of both population history and patterns of

41    selection, together with functional validation of variants associated to selected phenotypes.

42    Translation of these methods into non-model species is part of a shift in evolutionary sciences

43    that aims at better understanding biological diversity at various scales (Mandoli and

44    Olmstead, 2000; Jenner and Wills, 2007; Abzhanov *et al.*, 2008). Recent breakthroughs

45    brought by the study of initially non-model species (e.g. White *et al.*, 2010; Ellegren *et al.*,

46    2012; Weber *et al.*, 2013; Poelstra *et al.*, 2014) have confirmed the value of population

47    genomics from this perspective. These advances are needed to broaden our view about the

48    evolutionary process and improve sampling of distant clades.  Ultimately, this process should

49    provide a more balanced picture than the one brought by the study of a few model species

50    (Abzhanov *et al.*, 2008). Genomic approaches also have the potential to improve conservation

51    genetic inference by scaling up the amount of data available (Shafer *et al.*, 2015).

52    However, the widespread use of sophisticated analytical tools remains challenged by the lack

53    of communication between fields (Shafer *et al.*, 2015), little user-friendliness of software and

54    the ever-increasing amount of tools made available. Much effort has been put recently in

55    addressing these issues, but a lack of clarity subsists and many uncertainties remain. The

56    application of sometimes complex methods to species with little background has nonetheless

57    become more accessible, and has the potential to bring valuable information.

58    In this paper, I propose various methods and suggestions to deal with usual questions in

59    population genomics and genetics of adaptation in natural populations. I begin with a succinct

60  review of methods available to obtain genome-wide polymorphism data before focusing on i)

61  methods devoted to the study of population demographic history (Figure 1) and ii) methods

62  aiming at detecting signatures of selection (Figure 2).

63

## 64  Glossary

65  SNP: single nucleotide polymorphism.

66  Variant calling: identifying confidently genomic variants from alignment data (in SAM/BAM

67  format, see Li *et al.*, 2009). Classical SNP callers include the Genome Analysis Toolkit or

68  GATK (McKenna *et al.*, 2010), freebayes (Garrison and Marth, 2012), samtools (Li *et al.*,

69  2009) or Platypus (Rimmer *et al.*, 2014). Other tools call large-scale variants such as

70  inversions, translocations or copy-number variation (see main text).

71  Phasing: a process which identifies the alleles that are co-located on the same chromosome

72  copy.

73  Pooled sequencing: a protocol where tens or hundreds samples are pooled in a single library

74  prior sequencing (Futschik and Schlötterer, 2010). This prevents any individual identification

75  of each sample.

76

## 77  Obtaining genetic markers and linking them to a genome

### 78  **Common sequencing methods**

79  I consider here two main ways of dealing with genomics in non-model species: reduced

80  representation (Davey *et al.*, 2011) and whole-genome resequencing. Reduced representation

81  allows sampling homogeneously variants across the genome by sequencing DNA fragments

82  flanking restriction sites. Some of the best-known reduced representation techniques include

83  RAD-sequencing ( Baird *et al.*, 2008) and Genotyping by Sequencing or GBS (Elshire *et al.*,

84  2011). Their main interest is their relatively low cost and that they do not require any

85  reference genome (see Davey *et al.*, 2011 for details). The amount of SNPs ranges from

86  thousands to millions, which is most of the time enough to retrieve substantial information

87  about demography and sometimes selection (see Puritz *et al.*, 2014 for a detailed summary of

88  reduced-representation techniques).

3

89    Whole-genome resequencing requires a reference (at least at a draft stage) and is much more

90    expensive, especially for species with long and complex genomes. However, this approach

91    gives a complete overview of structural and coding variation, and allows some of the most

92    powerful methods currently available to track signatures of selection (see below). Pooled

93    sequencing (Futschik and Schlötterer, 2010) can be an option to reduce the costs, but restricts

94    the analysis to methods focusing on allele frequencies, losing most of the information

95    provided by variation in Linkage Disequilibrium (LD).

96    Shallow sequencing (1-5X per individual) may be a way to partly overpass this last issue for a

97    similar cost (Buerkle and Gompert, 2013), but should not be used for methods requiring

98    phasing and unbiased individual genotypes. Shallow shotgun sequencing also allows

99    retrieving complete plastomes, due to the representation bias of mitochondrial or chloroplast

100   sequences. Plastome sequences can provide insightful information about the evolutionary

101   history of populations or species. Recent work has successfully used shallow sequencing to

102   reconstruct mitochondrial or chloroplast sequences in plants (Malé *et al.*, 2014), animals

103   (Hahn *et al.*, 2013) or old and altered museum samples (Besnard *et al.*, 2016). Methods such

104   as MITObim (Hahn *et al.*, 2013) provide an automated and relatively user-friendly way to

105   reconstitute plastome sequences, which can then be analyzed as a single non-recombining

106   marker for phylogeny or population genetics.

**Obtain positional information for markers**

108   Whole-genome resequencing requires at least a draft genome, and reduced representations

109   methods can also benefit from a reference, either to order markers or retrieve information

110   about the nearest gene of a focal SNP. Methods inferring selection from haplotype extension

111   and patterns of LD (described further below) require that the relative order of markers on

112   genome sequence is known. A reference also allows analyzing separately sex chromosomes

113   (that can be haploid) and autosomes to correct for variation in ploidy between males and

114   females in gonochoric organisms. Obtaining a draft reference from deep Illumina sequencing

115   is now relatively common, but requires a good knowledge of assembly methods to choose the

116   tool adapted to the focal species. Initiatives such as Assemblathon (Bradnam *et al.*, 2013)

117   have provided valuable insights and advices from this regard. Once a draft is produced,

118   annotation of features is recommended since it allows linking variation at a locus to its

119   putative function. This requires either RNA-seq data to be mapped back on the reference or at

120   least that an annotation from a relatively close species is available.

4

121   It is possible to avoid these steps for species having a close relative already sequenced. Short-

122   reads alignment algorithms like BWA (Li and Durbin, 2009) generally assume relatively low

123   divergence between reads and reference. For species having less than 3% divergence, reads

124   may be directly mapped back onto the nearest genome. For more distantly related species, a

125   possible strategy would be using RAD-seq or GBS, build contigs for each locus with methods

126   like Stacks (Catchen *et al.*, 2011) or PyRAD (Eaton, 2014), and map those loci on the

127   reference with BLAT (Kent, 2002) or LASTZ (Schwartz *et al.*, 2003). Using a related

128   reference requires that synteny is conserved between species. While this assumption is

129   reasonable in, e.g., birds (Derjusheva *et al.*, 2004), it becomes more doubtful in other clades,

130   like in plants (Molinari *et al.*, 2008; Soltis *et al.*, 2015). Before conducting a NGS study, it is

131   therefore important to know how genomes vary in their structure across related species. Some

132   methods do not even require any reference sequence to call SNPs from raw reads, like kSNP2

133   (Gardner and Hall, 2013). It is however advised to cautiously filter reads prior calling, since

134   the method does not distinguish between sequencing errors and actual variants.

135   Checking for the presence of large structural variants can be informative when performing

136   whole-genome resequencing. Structural variants include duplications and copy number

137   variation (CNV), deletions, inversions or translocations. Neglecting this variation can lead to

138   call spurious SNPs, for example in regions which are single copy in the reference but display

139   CNV in some individual. This can distort estimates of nucleotide diversity or homozygosity,

140   biasing analyses based on LD or allele frequencies. These variations can be partly masked by

141   filtering SNPs on the basis of Hardy-Weinberg equilibrium or sequencing depth. However,

142   more quantitative methods are available that allow to precisely characterize the nature and the

143   position of this type of variation, like Delly (Rausch *et al.*, 2012) or Lumpy (Layer *et al.*,

144   2014).  Regions that display changes in genomic structure can then be excluded for analyses

145   requiring accurate estimates of diversity (e.g. Rasmussen *et al.*, 2014). On the other hand,

146   these variations can be used for studying association with traits of interest.

147

148   Assessing population history

149   **Exploring population structure**

150   Checking for population structure is an essential step when performing analyses on genome-

151   level datasets. Neglecting it can bias demographic inferences (Chikhi *et al.*, 2010; Heller *et*

5

152    *al.*, 2013) or the detection of loci under selection (e.g. Nielsen *et al.*, 2007); thus, checking for

153    outlier individuals and assessing the global structure is required prior any more sophisticated

154    analysis. A simple approach that does not assume any *a priori* grouping is the Principal

155    Component Analysis (PCA), based on analyzing variance-covariance structure among

156    genotypes, which can be performed on both individual and pooled data. Methods such as

157    SMARTPCA (Patterson *et al.*, 2006) or EIGENSTRAT (Patterson *et al.*, 2006) emerged from

158    this framework. There are many software solutions and packages allowing to perform this

159    type of analysis, such as SNPRelate (Zheng *et al.*, 2012), implemented in Bioconductor

160    (Huber *et al.*, 2015), PLINK (Purcell *et al.*, 2007) or GenAbel (Aulchenko *et al.*, 2007). For

161    large whole-genome data or high-density RAD-seq, reducing SNP redundancy by

162    subsampling unlinked markers (having low LD or large physical distance between them) is a

163    way to reduce computation time while keeping the relevant information.

164    Taking into account the relatedness of individuals is recommended, for example to evaluate

165    the amount of inbreeding within a population. When each individual in a study is sampled

166    from a different location or environment, estimating relatedness also provides a way to assess

167    the genetic distance between them, in relation with geographical or ecological distance (e.g.

168    Fields *et al.*, 2015). VCFTools (Danecek *et al.*, 2011) provides two ways calculating

169    relatedness; unadjusted Ajk (Yang *et al.*, 2010) and a kinship coefficient also implemented in

170    KING (Manichaikul *et al.*, 2010). It also allows calculating Hardy-Weinberg equilibrium.

171    Population stratification and relatedness can also be explored in PLINK based on pairwise

172    identity-by-state (IBS) distance or identity by descent (IBD).

173    Other approaches such as Structure (Pritchard *et al.*, 2000) and fastSTRUCTURE (Raj *et al.*,

174    2014) allow determining hierarchical population structure by grouping individuals in clusters

175    without any *a priori*. FastSTRUCTURE is computationally faster and more efficient with

176    large SNP datasets. These methods are also more efficient at detecting signatures of

177    admixture. Geneland (Guillot *et al.*, 2012), available as a R package, allows determining the

178    optimal number of population in a dataset by optimizing linkage and Hardy-Weinberg

179    equilibrium within clusters, and is also able to incorporate geographic coordinates in the

180    model to delineate their spatial organization. It can be useful to characterize the location and

181    shape of hybrid zones.

182    In order to properly test for the existence of hierarchical population structure, methods based

183    on differentiation measures (like $F_{st}$) can be used to build phylogenetic trees. POPTREE

184    (Takezaki *et al.*, 2010) allows to use various differentiation metrics to infer relationships

185   between populations. TreeMix (Pickrell and Pritchard, 2012) is a method building a

186   population tree based on the covariance matrix of population allele frequencies. It allows

187   tracking admixture events but requires the populations to be defined *a priori* (e.g. by a

188   Structure analysis). Other methods can use individual SNP data to reconstruct phylogenies,

189   like PhyML (Guindon *et al.*, 2010) or RAxML (Stamatakis, 2014). Splitstree (Huson and

190   Bryant, 2006) is a user-friendly software to compute phylogenies and networks on SNP

191   datasets and incorporate various methods for phylogeny reconstruction. Other pipelines, like

192   SNPhylo (Lee *et al.*, 2014), propose a complete framework from SNP filtering to tree

193   reconstruction that might help obtaining reliable topologies.

194   While useful to infer topologies, caution is advised when using branches lengths obtained

195   from SNP-only datasets, e.g. to calculate divergence times between different groups or

196   species (Leache *et al.*, 2015). For this purpose, it might therefore be easier to extract genes or

197   RAD contigs from the data and analyze them as DNA sequences in a software like BEAST2

198   (Drummond and Rambaut, 2007). In RAxML, a recent correction for bias on branch length

199   has been implemented that requires the number of monomorphic sites to be known (Leache *et*

200   *al.*, 2015) when providing only SNP alignment. Dating species or population divergence and

201   changes in population sizes using SNP data is also possible in SNAPP (Bryant *et al.*, 2012),

202   although the method requires long computing times when many markers are included. For

203   dating purpose and resolution of individual and population/species trees, BEAST2 and

204   BEAST* can also be used on sequence data for moderate-sized datasets (Drummond and

205   Rambaut, 2007).

206   As a general word of caution, it is important to remind that RAD-sequencing and related

207   methods display specific properties that can bias genome-wide estimates of diversity, like

208   allelic dropout (Arnold *et al.*, 2013). However, this type of markers remains valuable for

209   phylogenetic estimation, even for distantly related species (Cariou *et al.*, 2013).

210   To assess how diversity is partitioned across the different groups inferred by the methods

211   described previously, it is advisable to perform an Analysis of Molecular Variance

212   (AMOVA). Arlequin (Excoffier and Lischer, 2010) is particularly suited for this task. More

213   generally, investigating patterns of nucleotide diversity, inbreeding, $F_{st}$ or variation in LD

214   between populations and across the genome is useful to have a preliminary idea of the amount

215   of gene flow, admixture and variation in population sizes. These statistics can be easily

216   retrieved with VCFTools or PopGenome (Pfeifer *et al.*, 2014).

217 **Investigating population history with coalescent methods**

218 The coalescent has first emerged to provide population geneticists a way of modeling alleles
219 genealogy from a sample taken from a large population. Going backward in time, alleles
220 merge (coalesce) in a stochastic way until reaching their most recent common ancestor
221 (Kingman, 1982). A variety of methods used and enriched this theoretical framework to
222 resolve complex population histories and their associated demographic parameters, such as
223 divergence times, effective population sizes or gene flow. These parameters are usually scaled
224 by mutation rate per generation. Converting those parameters into demographic estimates
225 (e.g. time in years) requires that mutation rate and generation time be known or at least
226 reasonably well estimated, for example from other close species with similar life history.
227 Most well-known coalescent-based tools dedicated to population genetics include IMa (Hey
228 and Nielsen, 2007), Migrate-n (Beerli and Palczewski, 2010) or Lamarc (Kuhner, 2009).
229 Lamarc is the only one taking into account recombination in the model, the other ones
230 requiring non-recombining blocks of sequence or markers to be used. Although they are
231 powerful, these methods tend to be computationally slow (Excoffier *et al.*, 2013), since they
232 require a full evaluation of the likelihood function associated to the model, a procedure that
233 can be complex with hundreds or thousands of markers.

234 A way to bypass this issue has been the use of Approximate Bayesian Computation (ABC)
235 methods, which compare to the actual data a set of simulated data produced by coalescent
236 simulations under predefined scenarios. By measuring the distance between carefully chosen
237 summary statistics describing each simulation with those from the observed dataset, it is
238 possible to infer which scenario explains the data the best. DIYABC (Cornuet *et al.*, 2008) is
239 a popular and user-friendly software allowing to perform a full ABC analysis (from
240 simulations to model comparison), although it does not allow yet to model continuous gene
241 flow between populations. Another approach, which provides more control to the user,
242 consists in using coalescent simulators such as ms (Hudson, 2002) or fastsimcoal2 (Excoffier
243 and Foll, 2011). A pipeline allowing to perform all these steps is also available in
244 ABCToolbox (Wegmann *et al.*, 2010). Fastsimcoal is a bit slower to simulate data, but is
245 more user-friendly than ms, and more effective when simulating recombination for sequence
246 data. Once simulations are done, one can compute summary statistics for the simulated
247 datasets (e.g. with Arlequin when using fastsimcoal2), then use packages like abc in R
248 (Csilléry *et al.*, 2012) to perform model choice, cross-validation, estimate model
249 misclassification and demographic parameters. More information on how to perform a proper

250    ABC analysis can be found in the work by Csilléry *et al.* (2010). The main advantage of ABC

251    is that it allows handling arbitrarily complex models, unlike methods like IMa where the

252    model is predefined. However, using summary statistics leads to the loss of potentially useful

253    information.

254    More recently, new methods based on the allele frequency spectrum (AFS) emerged to

255    facilitate and speed up the analysis of large SNP datasets. Different patterns of gene flow and

256    demographic events all shape the AFS in specific ways (e.g. more alleles are likely to be

257    found at similar frequencies in two recently diverged or highly connected populations). ∂a∂i

258    (Gutenkunst *et al.*, 2009) does not rely on computationally intensive coalescent simulations

259    but rather on a diffusion approximation of alleles, and computes likelihoods for the alternative

260    models provided by the user. However, its current implementation does not handle more than

261    three populations. More recently, another likelihood-based approach has been implemented in

262    fastsimcoal2 (Excoffier *et al.*, 2013), that uses coalescent simulations and handles arbitrarily

263    complex scenarios while not being limited by the number of populations included. These two

264    methods assume that SNPs are under linkage equilibrium. Including SNPs in strong LD

265    should not particularly bias model comparison, but can be an issue when estimating

266    parameters (see fastsimcoal manual for more details). Note that the AFS can also be used as a

267    set of summary statistics for ABC inference. Using allele frequencies estimated from pooled

268    datasets should be feasible, although no study explored this possibility to my knowledge.

269    One drawback when using SNP data without considering monomorphic sites is that the

270    mutation rate per generation is not directly taken into account. For example, in DIYABC, it

271    does not matter when a mutation appears in the simulated genealogy, as long as it happens

272    only once before coalescence, a reasonable assumption for SNP markers. However, this

273    prevents any conversion of parameters into demographic estimates by using mutation rate.

274    Again, it is also possible to extract the complete DNA sequence for a set of randomly selected

275    markers and perform analyses on this dataset including monomorphic sites. Another

276    possibility consists in a calibration of parameter estimates by including in the analysis a fixed

277    parameter, such as population size or divergence time. This approach is also feasible when

278    estimating parameters from the allele frequency spectrum, like in ∂a∂i or fastsimcoal2.

279    When whole genome data are available, it is then possible to use methods such as those based

280    on Pairwise Sequentially Markovian Coalescent (PSMC), that require only a single diploid

281    genome (Li and Durbin, 2011). This method allows tracking changes in population size across

282    discrete time intervals. While powerful, PSMC is sensitive to confounding factors such as

283     population structure (Orozco-terWengel, 2016) that leads to false signatures of expansion or

284     bottleneck. It also does not allow studying recent demographic events. This is due to the fact

285     that coalescence events for only two alleles from a single individual in the recent past are

286     infrequent. However, extensions of the model allowing for several genomes have been

287     developed to precise population history in the recent past, like MSMC (Schiffels and Durbin,

288     2014) or diCal (Sheehan *et al.*, 2013). As these methods require that heterozygous positions

289     be properly called, it is required to correct for low depth of coverage (less than 8-10X) if

290     needed. Recently an ABC framework, implemented in PopSizeABC, has been proposed to

291     infer demographic variation from single genomes (Boistard *et al.*, 2016). The summary

292     statistics used describe variation in LD and the AFS, while being robust to sequencing errors.

293     This last method does not require phasing, which should limit the impact of phasing errors.

294     A recent extension of these methods takes into account population structure and aims at

295     identifying the number of islands contributing to a single genome, assuming it is sampled

296     from a Wright n-island meta-population (Mazet *et al.*, 2015). Such developments should help

297     increasing the amount of information retrieved from only a few genomes. However, it is

298     essential to keep in mind that natural populations are structured and connected in complex

299     ways, which can bias demographic inferences, even for popular markers such as

300     mitochondrial sequences (Heller *et al.*, 2013).

301     Reaching a high level of precision in demographic parameters estimation can be challenging

302     when perspective is lacking about the evolutionary history of the species considered. At larger

303     time-scales, the lack of fossil record can make difficult the calibration of molecular clocks.

304     Thus, for some species, only qualitative interpretation will be possible.

305

## 306     Screening for selection and association

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

308     The impact of selection on genetic variation has been extensively studied, but still remains a

309     central topic in evolutionary biology. Here I describe some features that are associated to

310     different types of selection.

311     Selection acts both on correlations i) between alleles and environment at selected loci and ii)

312     between alleles from different loci, either directly under selection or not. This is reflected

313    respectively by i) variation in polymorphism within and between populations and ii) linkage

314    disequilibrium between loci (Figure 2). A new mutation will see its frequency increase in a

315    population where it provides a selective advantage (hard sweep). When such an allele arose

316    recently, a large region around it can remain uniform, especially if selection is strong. As the

317    allele rises quickly in frequency, it has too little time to recombine with other ancestral

318    variants. This leads to an increase of linkage disequilibrium between variants associated to the

319    advantageous mutation, as well as a decrease in nucleotide diversity around the selected locus.

320    If selection occurs in one population but not others, it may be possible to observe a local

321    increase in differentiation, like higher $F_{st}$ values. If selection acts on standing variation or

322    recurrent mutation, signature of selection can be less clear as several haplotypes surround the

323    mutation under positive selection (see however Messer and Petrov, 2013; Jensen, 2014 for a

324    discussion about the relative importance of soft selective sweeps).

325    Another type of selection is balancing selection, an umbrella term grouping all selective

326    processes that lead to the maintenance of genetic polymorphism at a locus and to an excess of

327    common alleles. Such processes include divergent selection (the same allele is under positive

328    selection in one population and selected against in another one), negative frequency-

329    dependent selection (a rare allele is preferably selected) or heterozygote advantage. In the

330    case of recent balancing selection, the signature of selection is similar to a partial selective

331    sweep, with the recently selected allele displaying reduced diversity and higher LD than the

332    ancestral one. In the case of long-term balancing selection, there is an accumulation of genetic

333    polymorphism around the selected loci, leading to the maintenance of haplotypes older and

334    more diverse than in the rest of the genome. This increase in diversity can be associated to

335    higher local estimates of effective population sizes and effective recombination rates. As

336    alleles are older, coalescence times tend to be higher and can sometimes predate speciation,

337    leading to trans-species polymorphism (see Charlesworth, 2006 for a detailed review). In

338    some cases, an allele under balancing selection is stabilized at a single equilibrium frequency

339    across populations, which can lead to a signature of lower differentiation compared to

340    genomic background. There is still a lack of methods aiming at detecting specifically

341    balancing selection compared to positive selection and recent hard sweeps (but see Fijarczyk

342    and Babik, 2015).

343    In the following parts I present tools that can be used to detect signatures of selection. The

344    methods that these tools implement fall into three main categories (partly reviewed in Vitti *et*

345    *al.*, 2013), corresponding to the signature they try to target: i) study of variation in allele

346     frequencies and polymorphism, ii) study of variation in linkage disequilibrium and iii)

347     reconstruction of allele genealogies using the coalescent. Most of these methods assume that

348     markers are ordered along a genome; although they can also be used to extract individual

349     markers under selection that can be then be aligned (except for most LD-based methods).

350     **Methods focusing on polymorphism**

351     While demographic forces such as drift and migration will affect the whole genome in a

352     similar way, local effects of selection should produce discrepancies with genome-wide

353     polymorphism (Lewontin and Krakauer, 1973). Selection affects allele frequencies and

354     polymorphism in predictable ways at the scale of single populations. Several statistics

355     summarize them, like $\pi$, the nucleotide diversity (Nei and Li, 1979), Tajima's D (Tajima,

356     1989) or Fay and Wu's H (Fay and Wu, 2000). They are sensitive to population demographic

357     history, that they allow characterizing as summary statistics in, e.g., ABC analyses. They have

358     nonetheless the potential to highlight genomic regions displaying clear signatures of selection,

359     or to confirm selection at candidate genes. For example, balancing selection should lead to an

360     excess of common polymorphisms, similar to a recent bottleneck, leading to high Tajima's D

361     and $\pi$ values. Purifying selection leads to an opposite pattern, similar to a recent population

362     expansion, with an excess of rare variants and low diversity. More sophisticated methods

363     using allele frequency spectrum have been developed to detect positive selection, such as the

364     recent improvement of the composite likelihood ratio (CLR) test (Nielsen *et al.*, 2005)

365     performed in SweepFinder2 (Degiorgio *et al.*, 2016).

366     PopGenome (Pfeifer *et al.*, 2014) is a powerful R package that allows calculating AFS

367     statistics (including the CLR test) across many genomes, as well as a variety of statistics on

368     linkage disequilibrium and diversity. It also allows performing coalescent simulations to

369     contrast observed polymorphism to neutral expectations. It is probably one of the most

370     comprehensive tools to perform genome-wide analyses. Other possibilities include VCFTools,

371     POPBAM (Garrigan, 2013) or Biopython libraries. For pooled data, Popoolation (Kofler,

372     Orozco-terWengel, *et al.*, 2011; Kofler, Pandey, *et al.*, 2011) provides ways to calculate

373     Tajima's D and nucleotide diversity, as well as measures of differentiation between

374     populations.

375     Understanding the origin of genomic regions under selection highlights the evolutionary

376     history of adaptive alleles (e.g. Abi-Rached *et al.*, 2011). Advantageous alleles can migrate

377     from one population to another, or resist introgression from other populations (genomic

378   islands of speciation/adaptation). The relative importance of these islands resisting gene flow

379   after secondary contact has been recently discussed (Cruickshank and Hahn, 2014). Methods

380   aiming at characterizing heterogeneity in introgression rates are in this context useful and can

381   also refine the demographic history. A recent ABC framework has been developed to

382   characterize this heterogeneity  (Roux *et al.*, 2014). Methods such as PCAdmix (Brisbin *et al.*,

383   2012) can also be used to estimate the relative contributions of putative sources to a given

384   sink population across the genome. A common test for introgression, available in

385   PopGenome, is the ABBA-BABA test, summarized by Patterson's D (Durand *et al.*, 2011).

386   Another possibility lies in the comparison of absolute and relative measures of divergence

387   (Cruickshank and Hahn, 2014), such as $d_{xy}$ and $F_{st}$, which can be calculated in PopGenome.

388   Absolute measures of divergence are correlated to the time since coalescence. In the case of

389   local introgression, both statistics should be reduced. For balancing selection, the decline in

390   $F_{st}$ is due to an excess of shared ancestral alleles, which should not impact $d_{xy}$, or should even

391   make it higher than genomic background. However, these methods do not prevent false

392   positives and results should (as usual) be interpreted with caution (Martin *et al.*, 2015).

393   When an allele is under positive selection in a population, its frequency tends to rise until

394   fixation, unless gene flow from other populations or strong drift prevents it. It is therefore

395   possible to contrast patterns of differentiation between populations adapted to their local

396   environment to detect loci under divergent selection (e.g. displaying a high $F_{st}$). However, it is

397   essential to control for population structure, as it may strongly affect the distribution of

398   differentiation measures and produce high rates of false positive. First attempts to take into

399   account population structure and variation in gene flow included FDIST2 (Beaumont and

400   Nichols, 1996), which modeled populations as islands and aimed at detecting loci under

401   selection by contrasting heterozygosity to $F_{st}$ between populations. An extension of this

402   model, able to take into account predefined hierarchical population structure, is implemented

403   in Arlequin. More sophisticated methods are now available, dedicated to the detection of

404   outliers in large genomic datasets. Most of them correct for relatedness across samples, and

405   are reviewed extensively in the work by Francois *et al.*(2015). Some methods, like LFMM

406   (Frichot *et al.*, 2013), aim at detecting variants correlated to environmental factors.

407   Association methods may help targeting variants undergoing soft sweeps, weak selection or

408   involved in polygenic control of traits (Pritchard *et al.*, 2010), for which signatures of

409   selection are subtle and sometimes difficult to retrieve from allele frequencies data.

410  Other methods perform a "naïve scan" for outliers on the basis of differentiation, like

411  BAYESCAN (Foll and Gaggiotti, 2008) which considers all populations to drift at different

412  rates from a single ancestral pool. Most recent methods, like BAYENV (Günther and Coop,

413  2013) and its recent improvement, BAYPASS (Gautier, 2015), model demographic history by

414  computing a kinship matrix between populations. Contrasting allele frequencies for each

415  locus to the ones expected given this matrix allows testing deviation from neutrality. Those

416  two last methods also include Bayesian tools to test for association with environmental

417  features, facilitating further interpretation. BAYENV and BAYPASS also allow using pooled-

418  sequencing data, making these methods polyvalent and possibly useful to many research

419  teams. However, detecting association between environment and allele frequencies does not

420  necessarily imply a role for local adaptation. For example, in the case of secondary contact,

421  intrinsic genetic incompatibilities can lead to the formation of tension zones that may shift

422  until they reach an environmental barrier where they can be trapped (Bierne *et al.*, 2011).

423  Again, characterizing population history is required to conclude about the possible

424  involvement of a genomic region in adaptation to environment.  Sampling strategy must take

425  into account the particular historical and demographic features of the species investigated to

426  gain power (Nielsen *et al.*, 2007). The sequencing strategy has also to be carefully picked.

427  Reduced representation methods do not cover all mutations in the genome and are thus more

428  likely to miss those actually under selection. Special care in the choice of the restriction

429  enzyme and determining the expected density of markers is needed to retrieve enough

430  mutations close to genes under selection.

431  The methods described above focus on allele frequencies at the population scale, but do not

432  allow characterizing properly association with a trait varying between individuals within

433  populations (e.g. resistance to a pathogen, symbiotic association, individual size or flowering

434  time). For this task, methods performing Genome-wide association analysis (GWAS) are

435  better suited. Methods such as GenAbel in R (Aulchenko *et al.*, 2007) or PLINK (Purcell *et*

436  *al.*, 2007) are powerful tool. Taking into account relatedness between samples and population

437  history (e.g. using EIGENSTRAT or PC-adjustment corrections in GenAbel or stratified

438  analyses in PLINK) is required to correct for false positives. This is especially recommended

439  for species that undergo episodes of selfing or strong bottlenecks, for which sampling

440  unrelated individuals may be unfeasible.

441  It is important to keep in mind that uncovering the genetic bases of complex, polygenic traits

442  remains challenging, even in model species (Pritchard and Di Rienzo, 2010; Rockman, 2012).

443    It may be unavoidable in a first step to focus only on traits that are under a relatively simple

444    genetic determinism. This can however lead to an overrepresentation of loci of major

445    phenotypic effect, a fact that should be acknowledged when discussing the impact of selection

446    on genome variation. The fact that loci of major effect are easier to target does not imply that

447    they are the main substrate of selection (Rockman, 2012).

**Detecting selection with methods focusing on LD**

449    LD is increased and diversity is decreased in the vicinity of a selected allele, especially after

450    recent selection. A class of methods aims at targeting those regions that display an excess of

451    long homozygous haplotypes, such as the extended haplotype homozygosity (EHH) test

452    (Sabeti *et al.*, 2002). It is also possible to compare haplotype extension across populations,

453    with the XP-EHH test (McCarroll *et al.*, 2007) or Rsb (Tang *et al.*, 2007). Individuals

454    included in the analysis should be as distantly related as possible to improve precision and

455    avoid an excess of false positives. These approaches are more powerful with a relatively high

456    density of markers, such as the ones obtained from whole-genome sequencing or high-density

457    RAD-seq. They also require data to be phased in order to reconstruct haplotypes. This

458    procedure can be performed with fastPhase (Scheet and Stephens, 2006), BEAGLE

459    (Browning and Browning, 2011) or SHAPEIT2 (O'Connell *et al.*, 2014). The R package rehh

460    (Gautier and Vitalis, 2012) allows  calculating these statistics, as well as Sweep

461    (http://www.broadinstitute.org/mpg/sweep/index.html). Statistics dedicated to the detection of

462    soft sweeps are also available, like the H2/H1 statistics (Garud *et al.*, 2015), although further

463    studies are still needed to understand to what extent hard and soft sweeps can actually be

464    distinguished (Schrider *et al.*, 2015). This last statistics does not require data to be phased.

465    When the relative order of markers is not known, as it can be the case in RAD-seq studies

466    without a reference genome, LDna (Kemppainen *et al.*, 2015) can be used to target sets of

467    markers displaying strong linkage disequilibrium. This approach can be useful not only to

468    detect selection but also structural variation such as large inversions.

469    Even hard selective sweeps can be challenging to detect with LD-based statistics (Jensen,

470    2014). It is advisable to combine several approaches to reach a better confidence when

471    pinpointing candidate genes for selection. Methods based on LD alone can sometimes miss

472    the actual variants under selection due to the impact of recombination on local polymorphism

473    that can mimic soft or ongoing hard sweeps (Schrider *et al.*, 2015).

**Detecting and characterizing selection with the coalescent**

475    When a candidate locus has been identified, it is possible to use coalescent simulations to

476    evaluate the strength of selection and estimate the age of alleles. A software such as msms

477    (Ewing and Hermisson, 2010), which is also available in PopGenome, can then be used. This

478    requires that the neutral history of population be known in order to properly control for, e.g.,

479    population structure and gene flow.

480    An advantage of full coalescent methods is that they provide a relatively complete picture of

481    individual loci history, by modeling coalescence, recombination and taking into account

482    variation in mutation rate. They are however computationally intensive, and thus difficult to

483    apply to whole genomes. However, recent computational improvements make this procedure

484    feasible, as illustrated by ARGWeaver (Rasmussen *et al.*, 2014). This method uses ancestral

485    recombination graphs to model the genealogy of each non-recombining block in the genome.

486    It allows extracting genealogies for these blocks and provides estimates for local

487    recombination rate, coalescence time and local effective population size for each block. This

488    approach is promising to characterize positive, purifying or balancing selection while taking

489    into account variation in recombination and mutation rate. However, the high stochasticity in

490    parameters estimation can limit resolution when targeting single genes.

491    Other methods use the theoretical framework of the coalescent to target sites under positive

492    selection. A recent method (SCCT) using conditional coalescent trees (Wang *et al.*, 2014)

493    claims to be faster and more precise in targeting selective sweeps. BALLET (DeGiorgio *et al.*,

494    2014) is a promising method to characterize ancient balancing selection. Most of those

495    methods are designed for medium-to-high depth whole-genome resequencing, and require that

496    individual genotypes be phased and well characterized.

497    **Variants annotation**

498    Characterizing the amount of synonymous or non-synonymous mutations is another way to

499    detect whether a specific gene undergoes purifying or positive selection. An excess of non-

500    synonymous mutations can signal positive or balancing selection, or a relaxation of selective

501    constraints on a given gene. This requires that an annotated genome is available. Annotation

502    of mutations can be done with SNPdat (Doran and Creevey, 2013), or directly in PopGenome,

503    which can also perform at the genome scale tests of selection such as the MK test (McDonald

504    and Kreitman, 1991). The MK test compares the amount of fixed and polymorphic mutations

505    relative to an outgroup, according to their synonymous/non-synonymous state. Another

506    popular test of selection is the comparison of non-synonymous and synonymous mutations

16

507     between orthologs from different species and can be performed in packages such as PAML

508     (Yang, 2007).

509     To recover information about the putative function of a gene or a genomic region, it may be

510     useful to perform a genome ontology (GO) enrichment analysis. BLAST2GO (Conesa *et al.*,

511     2005) allows annotating genes by using a database of related species. It also allows

512     performing GO enrichment analysis. These analyses must be carefully interpreted, depending

513     on the level of divergence from the closest annotated species. It is important not to jump to

514     the conclusion that orthologous genes must share the same function. When interpreting the

515     link between selection and genetic variation, a careful review of literature can fruitfully

516     complete the conclusions made using GO enrichment analyses.

517

## Concluding remarks

519     In this contribution I highlighted different methods currently available to investigate how

520     history and selection shape diversity in natural populations. It is important to understand that

521     this dichotomy between selection and demography, while practical, remains artificial, and that

522     the study of one benefits from studying the other. With the decreasing cost of sequencing it

523     has been suggested that NGS should broaden quickly our perspective on complex

524     evolutionary processes, from biogeography (Lexer *et al.*, 2013) to genetic bases of traits

525     (Hohenlohe, 2014) or the maintenance of polymorphism (Hedrick, 2006). The study of DNA

526     sequence variation, while already challenging by itself, needs to be combined with other

527     disciplines such as ecology to be informative (Habel *et al.*, 2015). Although genome-scale

528     analyses can be insightful to this regard, it is necessary to be conscious of their limits and to

529     keep a biological perspective when interpreting their results. To do so, every analysis should

530     always begin with a proper understanding of the methods used, to avoid using them as black

531     boxes.

532     Before launching a project targeting thousands of markers in a species of interest, possibilities

533     and limits of the chosen protocol must be evaluated. Focusing on species history will not

534     necessarily require the same sampling strategy than focusing on local adaptation. While a

535     small number of markers and populations may be enough to recover global structure and infer

536     robustly demographic parameters, it will not provide enough resolution to target genes

537     involved in local adaptation. In many cases, a preliminary study focusing on a few markers

538     may already inform about the global species history and help to define an adapted design for

539     NGS.

540     There is a need for a more collaborative and open culture in biology, allowing the free access

541     to data and favoring good practices to allow repeatability of analyses (Nekrutenko and Taylor,

542     2012), although this cultural shift remains challenging (e.g. Mills *et al.*, 2015; Whitlock *et al.*,

543     2015). However, current challenges are not limited to data sharing, but also include dealing

544     with the inflation of bioinformatics tools that sometimes overlap. Instead of working

545     independently, researchers designing those tools could collaborate to propose free, robust and

546     unified pipelines (Prins *et al.*, 2015). Such initiatives, like Galaxy (Goecks *et al.*, 2010) or

547     Bioconductor (Huber *et al.*, 2015) are nonetheless emerging ; this should facilitate the

548     emergence of a unified framework to limit the time dedicated to data analysis and focus on

549     biological questions.

550

## 551     References

552

553     Abi-Rached L, Jobin M, Kulkarni S, McWhinnie A, Dalva K, Gragert L, *et al.* (2011). The

554        shaping of modern human immune systems by multiregional admixture with archaic

555        humans. *Science (80- )* **334**: 89–95.

556     Abzhanov A, Extavour CG, Groover A, Hodges SA, Hoekstra HE, Kramer EM, *et al.* (2008).

557        Are we there yet? Tracking the development of new model systems. *Trends Genet* **24**:

558        353–60.

559     Arnold B, Corbett-Detig RB, Hartl D, Bomblies K (2013). RADseq underestimates diversity

560        and introduces genealogical biases due to nonrandom haplotype sampling. *Mol Ecol* **22**:

561        3179–90.

562     Aulchenko YS, Ripke S, Isaacs A, van Duijn CM (2007). GenABEL: An R library for

563        genome-wide association analysis. *Bioinformatics* **23**: 1294–1296.

564     Axelsson E, Ratnakumar A, Arendt M-L, Maqbool K, Webster MT, Perloski M, *et al.* (2013).

565        The genomic signature of dog domestication reveals adaptation to a starch-rich diet.

566        *Nature* **495**: 360–4.

18

567    Baird N a, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis Z a, *et al.* (2008). Rapid SNP
568        discovery and genetic mapping using sequenced RAD markers. *PLoS One* **3**: e3376.

569    Beaumont MA, Nichols RA (1996). Evaluating loci for use in the genetic analysis of
570        population structure. *Proc R Soc London Biol Sci*: 1619–1626.

571    Beerli P, Palczewski M (2010). Unified framework to evaluate panmixia and migration
572        direction among multiple sampling locations. *Genetics* **185**: 313–26.

573    Besnard G, Bertrand JAM, Delahaie B, Bourgeois YXC, Lhuillier E, Thébaud C (2016).
574        Valuing museum specimens: high-throughput DNA sequencing on historical collections
575        of New Guinea crowned pigeons (Goura). *Biol J Linn Soc* **117**: 71–82.

576    Bierne N, Welch J, Loire E, Bonhomme F, David P (2011). The coupling hypothesis: why
577        genome scans may fail to map local adaptation genes. *Mol Ecol* **20**: 2044–72.

578    Boistard S, Rodriguez W, Jay F, Mona S, Austerlitz F (2016). Inferring Population Size
579        History from Large Samples of Genome-Wide Molecular Data - An Approximate
580        Bayesian Computation Approach. *PLoS Genet*: 858–865.

581    Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, *et al.* (2013).
582        Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate
583        species. *Gigascience* **2**: 10.

584    Brisbin A, Bryc K, Byrnes J, Zakharia F, Omberg L, Degenhardt J, *et al.* (2012). PCAdmix:
585        Principal Components-Based Assignment of Ancestry along Each Chromosome in
586        Individuals with Admixed Ancestry from Two or More Populations. *Hum Biol* **84**: 343–
587        364.

588    Browning BL, Browning SR (2011). A fast, powerful method for detecting identity by
589        descent. *Am J Hum Genet* **88**: 173–182.

590    Bryant D, Bouckaert R, Felsenstein J, Rosenberg NA, Roychoudhury A (2012). Inferring
591        species trees directly from biallelic genetic markers: Bypassing gene trees in a full
592        coalescent analysis. *Mol Biol Evol* **29**: 1917–1932.

593    Buerkle CA, Gompert Z (2013). Population genomics based on low coverage sequencing:
594        how low should we go? *Mol Ecol* **22**: 3028–35.

595    Cariou M, Duret L, Charlat S (2013). Is RAD-seq suitable for phylogenetic inference? An in

596        silico assessment and optimization. *Ecol Evol* **3**: 846–852.

597    Catchen JM, Amores A, Hohenlohe P, Cresko W, Postlethwait JH (2011). Stacks: building

598        and genotyping Loci de novo from short-read sequences. *G3 (Bethesda)* **1**: 171–82.

599    Charlesworth D (2006). Balancing selection and its effects on sequences in nearby genome

600        regions. *PLoS Genet* **2**: e64.

601    Chikhi L, Sousa VC, Luisi P, Goossens B, Beaumont MA (2010). The confounding effects of

602        population structure, genetic diversity and the sampling scheme on the detection and

603        quantification of population size changes. *Genetics* **186**: 983–995.

604    Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M (2005). Blast2GO: A

605        universal tool for annotation, visualization and analysis in functional genomics research.

606        *Bioinformatics* **21**: 3674–3676.

607    Cornuet J-M, Santos F, Beaumont M a, Robert CP, Marin J-M, Balding DJ, *et al.* (2008).

608        Inferring population history with DIY ABC: a user-friendly approach to approximate

609        Bayesian computation. *Bioinformatics* **24**: 2713–9.

610    Cruickshank TE, Hahn MW (2014). Reanalysis suggests that genomic islands of speciation

611        are due to reduced diversity, not reduced gene flow. *Mol Ecol* **23**: 3133–3157.

612    Csilléry K, Blum MGB, Gaggiotti OE, François O (2010). Approximate Bayesian

613        Computation (ABC) in practice. *Trends Ecol Evol* **25**: 410–8.

614    Csilléry K, François O, Blum MGB (2012). abc: an R package for approximate Bayesian

615        computation (ABC). *Methods Ecol Evol* **3**: 475–479.

616    Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, *et al.* (2011). The

617        variant call format and VCFtools. *Bioinformatics* **27**: 2156–2158.

618    Davey JW, Hohenlohe P a, Etter PD, Boone JQ, Catchen JM, Blaxter ML (2011). Genome-

619        wide genetic marker discovery and genotyping using next-generation sequencing. *Nat*

620        *Rev Genet* **12**: 499–510.

621    Degiorgio M, Huber CD, Hubisz MJ, Hellmann I, Nielsen R (2016). Genetics and population

622        analysis SWEEPFINDER 2□: Increased sensitivity , robustness , and flexibility.

623        *Bioinformatics*.

624  DeGiorgio M, Lohmueller KE, Nielsen R (2014). A model-based approach for identifying
625      signatures of ancient balancing selection in genetic data. *PLoS Genet* **10**: e1004561.

626  Derjusheva S, Kurganova A, Habermann F, Gaginskaya E (2004). High chromosome
627      conservation detected by comparative chromosome painting in chicken, pigeon and
628      passerine birds. *Chromosome Res* **12**: 715–23.

629  Doran AG, Creevey CJ (2013). Snpdat: easy and rapid annotation of results from de novo snp
630      discovery projects for model and non-model organisms. *BMC Bioinformatics* **14**: 45.

631  Drummond AJ, Rambaut A (2007). BEAST: Bayesian evolutionary analysis by sampling
632      trees. *BMC Evol Biol* **7**: 214.

633  Durand EY, Patterson N, Reich D, Slatkin M (2011). Testing for ancient admixture between
634      closely related populations. *Mol Biol Evol* **28**: 2239–2252.

635  Eaton DAR (2014). PyRAD: Assembly of de novo RADseq loci for phylogenetic analyses.
636      *Bioinformatics* **30**: 1844–1849.

637  Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, *et al.* (2012). The
638      genomic landscape of species divergence in Ficedula flycatchers. *Nature* **491**: 756–60.

639  Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, *et al.* (2011). A
640      Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species.
641      *PLoS One* **6**: e19379.

642  Ewing G, Hermisson J (2010). MSMS: A coalescent simulation program including
643      recombination, demographic structure and selection at a single locus. *Bioinformatics* **26**:
644      2064–2065.

645  Excoffier L, Dupanloup I, Huerta-Sanchez E, Sousa VC, Foll M (2013). Robust Demographic
646      Inference from Genomic and SNP Data. *PLoS Genet* **9**.

647  Excoffier L, Foll M (2011). Fastsimcoal: a Continuous-Time Coalescent Simulator of
648      Genomic Diversity Under Arbitrarily Complex Evolutionary Scenarios. *Bioinformatics*
649      **27**: 1332–4.

650  Excoffier L, Lischer HEL (2010). Arlequin suite ver 3.5: a new series of programs to perform
651      population genetics analyses under Linux and Windows. *Mol Ecol Resour* **10**: 564–7.

652 Fay JC, Wu CI (2000). Hitchhiking under positive Darwinian selection. *Genetics* **155**: 1405–
653     13.

654 Fields PD, Reisser C, Dukic M, Haag CR, Ebert D (2015). Genes mirror geography in
655     Daphnia magna. *Mol Ecol* **24**: 4521–4536.

656 Fijarczyk A, Babik W (2015). Detecting balancing selection in genomes: limits and prospects.
657     *Mol Ecol* **24**: 3529–3545.

658 Foll M, Gaggiotti O (2008). A genome-scan method to identify selected loci appropriate for
659     both dominant and codominant markers: a Bayesian perspective. *Genetics* **180**: 977–93.

660 François O, Martins H, Caye K, Schoville SD (2015). Controlling False Discoveries in
661     Genome Scans for Selection. *Mol Ecol* **55**: in press.

662 Frichot E, Schoville SD, Bouchard G, François O (2013). Testing for associations between
663     loci and environmental gradients using latent factor mixed models. *Mol Biol Evol* **30**:
664     1687–1699.

665 Futschik A, Schlötterer C (2010). The next generation of molecular markers from massively
666     parallel sequencing of pooled DNA samples. *Genetics* **186**: 207–18.

667 Gardner SN, Hall BG (2013). When whole-genome alignments just won't work: KSNP v2
668     software for alignment-free SNP discovery and phylogenetics of hundreds of microbial
669     genomes. *PLoS One* **8**.

670 Garrigan D (2013). POPBAM: Tools for evolutionary analysis of short read sequence
671     alignments. *Evol Bioinforma* **2013**: 343–353.

672 Garrison E, Marth G (2012). Haplotype-based variant detection from short-read sequencing.
673     *arXiv Prepr arXiv12073907*: 9.

674 Garud NR, Messer PW, Buzbas EO, Petrov DA (2015). Recent Selective Sweeps in North
675     American Drosophila melanogaster Show Signatures of Soft Sweeps. *PLoS Genet* **11**: 1–
676     32.

677 Gautier M (2015). Genome-Wide Scan for Adaptive Divergence and Association with
678     Population-Specific Covariates. *Genetics* **XXX**: XXXX–XXXX.

679 Gautier M, Vitalis R (2012). Rehh An R package to detect footprints of selection in genome-

22

680  wide SNP data from haplotype structure. *Bioinformatics* **28**: 1176–1177.

681  Goecks J, Nekrutenko A, Taylor J (2010). Galaxy: a comprehensive approach for supporting
682      accessible, reproducible, and transparent computational research in the life sciences.
683      *Genome Biol* **11**: R86.

684  Guillot G, Renaud S, Ledevin R, Michaux J, Claude J (2012). A unifying model for the
685      analysis of phenotypic, genetic, and geographic data. *Syst Biol* **61**: 897–911.

686  Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O (2010). New
687      algorithms and methods to estimate maximum-likelihood phylogenies: Assessing the
688      performance of PhyML 3.0. *Syst Biol* **59**: 307–321.

689  Günther T, Coop G (2013). Robust identification of local adaptation from allele frequencies.
690      *Genetics* **195**: 205–220.

691  Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD (2009). Inferring the joint
692      demographic history of multiple populations from multidimensional SNP frequency data.
693      *PLoS Genet* **5**.

694  Habel J, Zachos F, Dapporto L, Rödder D, Radespiel U, Tellier A, *et al.* (2015). Population
695      genetics revisited – towards a multidisciplinary research field. *Biol J Linn Soc* **115**: 1–12.

696  Hahn C, Bachmann L, Chevreux B (2013). Reconstructing mitochondrial genomes directly
697      from genomic next-generation sequencing reads - A baiting and iterative mapping
698      approach. *Nucleic Acids Res* **41**.

699  Hedrick PW (2006). Genetic Polymorphism in Heterogeneous Environments: The Age of
700      Genomics. *Annu Rev Ecol Evol Syst* **37**: 67–93.

701  Heller R, Chikhi L, Siegismund HR (2013). The Confounding Effect of Population Structure
702      on Bayesian Skyline Plot Inferences of Demographic History. *PLoS One* **8**.

703  Hey J, Nielsen R (2007). Integration within the Felsenstein equation for improved Markov
704      chain Monte Carlo methods in population genetics. *Proc Natl Acad Sci U S A* **104**: 2785–
705      90.

706  Hohenlohe P a (2014). Ecological genomics in full colour. *Mol Ecol* **23**: 5129–31.

707  Huber W, Carey VJ, Gentleman R, Anders S, Carlson M, Carvalho BS, *et al.* (2015).

Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Methods* **12**: 115–121.

Hudson RR (2002). Generating samples under a Wright–Fisher neutral model of genetic variation. *Bioinformatics* **18**: 337–338.

Huson DH, Bryant D (2006). Application of phylogenetic networks in evolutionary studies. *Mol Biol Evol* **23**: 254–267.

Jenner RA, Wills MA (2007). The choice of model organisms in evo-devo. *Nat Rev Genet* **8**: 311–319.

Jensen JD (2014). On the unfounded enthusiasm for soft selective sweeps. *Nat Commun* **5**: 5281.

Kemppainen P, Knight CG, Sarma DK, Hlaing T, Prakash A, Maung Maung YN, *et al.* (2015). Linkage disequilibrium network analysis (LDna) gives a global view of chromosomal inversions, local adaptation and geographic structure. *Mol Ecol Resour*: 1031–1045.

Kent WJ (2002). BLAT — The BLAST -Like Alignment Tool. *Genome Res* **12**: 656–664.

Kingman JFC (1982). The coalescent. *Stoch Process their Appl* **13**: 235–248.

Kofler R, Orozco-terWengel P, De Maio N, Pandey RV, Nolte V, Futschik A, *et al.* (2011). PoPoolation: a toolbox for population genetic analysis of next generation sequencing data from pooled individuals. *PLoS One* **6**: e15925.

Kofler R, Pandey RV, Schlötterer C (2011). PoPoolation2: identifying differentiation between populations using sequencing of pooled DNA samples (Pool-Seq). *Bioinformatics* **27**: 3435–6.

Kolaczkowski B, Kern AD, Holloway AK, Begun DJ (2011). Genomic differentiation between temperate and tropical Australian populations of Drosophila melanogaster. *Genetics* **187**: 245–60.

Kubota S, Iwasaki T, Hanada K, Nagano AJ, Fujiyama A, Toyoda A, *et al.* (2015). A Genome Scan for Genes Underlying Microgeographic-Scale Local Adaptation in a Wild Arabidopsis Species. *PLoS Genet* **11**: 1–26.

736 Kuhner MK (2009). Coalescent genealogy samplers: windows into population history. *Trends*
737     *Ecol Evol* **24**: 86–93.

738 Layer RM, Chiang C, Quinlan AR, Hall IM (2014). LUMPY: a probabilistic framework for
739     structural variant discovery. *Genome Biol* **15**: R84.

740 Leache AD, Banbury BL, Felsenstein J, De Oca ANM, Stamatakis A (2015). Short tree, long
741     tree, right tree, wrong tree: New acquisition bias corrections for inferring SNP
742     phylogenies. *Syst Biol* **64**: 1032–1047.

743 Lee T-H, Guo H, Wang X, Kim C, Paterson AH (2014). SNPhylo: a pipeline to construct a
744     phylogenetic tree from huge SNP data. *BMC Genomics* **15**: 162.

745 Legrand D, Tenaillon MI, Matyot P, Gerlach J, Lachaise D, Cariou M-L (2009). Species-wide
746     genetic variation and demographic history of Drosophila sechellia, a species lacking
747     population structure. *Genetics* **182**: 1197–206.

748 Lewontin RC, Krakauer J (1973). Distribution of gene frequency as a test of the theory of the
749     selective neutrality of polymorphisms. *Genetics* **74**: 175–195.

750 Lexer C, Mangili S, Bossolini E, Forest F, Stölting KN, Pearman PB, *et al.* (2013). 'Next
751     generation' biogeography: towards understanding the drivers of species diversification
752     and persistence (M Carine, Ed.). *J Biogeogr* **40**: 1013–1022.

753 Li H, Durbin R (2009). Fast and accurate short read alignment with Burrows-Wheeler
754     transform. *Bioinformatics* **25**: 1754–60.

755 Li H, Durbin R (2011). Inference of human population history from individual whole-genome
756     sequences. *Nature* **475**: 493–496.

757 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, *et al.* (2009). The Sequence
758     Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–9.

759 Malé PJG, Bardon L, Besnard G, Coissac E, Delsuc F, Engel J, *et al.* (2014). Genome
760     skimming by shotgun sequencing helps resolve the phylogeny of a pantropical tree
761     family. *Mol Ecol Resour* **14**: 966–975.

762 Mandoli DF, Olmstead R (2000). The importance of emerging model systems in plant
763     biology. *J Plant Growth Regul* **19**: 249–252.

764   Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen W-M (2010). Robust
765       relationship inference in genome-wide association studies. *Bioinformatics* **26**: 2867–
766       2873.

767   Martin SH, Davey JW, Jiggins CD (2015). Evaluating the use of ABBA-BABA statistics to
768       locate introgressed loci. *Mol Biol Evol* **32**: 244–257.

769   Mazet O, Rodriguez W, Chikhi L (2015). Demographic inference using genetic data from a
770       single individual: Separating population size variation from population structure. *Theor*
771       *Popul Biol* **104**: 46–58.

772   McCarroll SA, Sabeti PC, Frazer KA, Varilly P, Fry B, Ballinger DG, *et al.* (2007). Genome-
773       wide detection and characterization of positive selection in human populations. *Nature*
774       **449**: 913–8.

775   McDonald JH, Kreitman M (1991). Adaptive protein evolution at the Adh locus in
776       Drosophila. *Nature* **351**: 652–4.

777   McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, *et al.* (2010). The
778       genome analysis toolkit: A MapReduce framework for analyzing next-generation DNA
779       sequencing data. *Genome Res* **20**: 1297–1303.

780   Messer PW, Petrov DA (2013). Population genomics of rapid adaptation by soft selective
781       sweeps. *Trends Ecol Evol* **28**: 659–669.

782   Mills JA, Teplitsky C, Arroyo B, Charmantier A, Becker PH, Birkhead TR, *et al.* (2015).
783       Archiving Primary Data: Solutions for Long-Term Studies. *Trends Ecol Evol* **30**: 581–
784       589.

785   Molinari NA, Petrov DA, Price HJ, Smith JD, Gold JR, Vassiliadis C, *et al.* (2008). Synteny
786       and Collinearity in Plant Genomes. *Science (80- )*: 486–489.

787   Nei M, Li WH (1979). Mathematical model for studying genetic variation in terms of
788       restriction endonucleases. *Proc Natl Acad Sci U S A* **76**: 5269–73.

789   Nekrutenko A, Taylor J (2012). Next-generation sequencing data interpretation: enhancing
790       reproducibility and accessibility. *Nat Rev Genet* **13**: 667–72.

791   Nielsen R, Hellmann I, Hubisz M, Bustamante C, Clark AG (2007). Recent and ongoing
792       selection in the human genome. *Nat Rev Genet* **8**: 857–868.

793 Nielsen R, Williamson S, Kim Y, Nielsen R, Williamson S, Kim Y, *et al.* (2005). Genomic
794      scans for selective sweeps using SNP data Genomic scans for selective sweeps using
795      SNP data. *Genome Res*: 1566–1575.

796 O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, *et al.* (2014). A General
797      Approach for Haplotype Phasing across the Full Spectrum of Relatedness. *PLoS Genet*
798      **10**.

799 Orozco-terWengel P (2016). The devil is in the details: the effect of population structure on
800      demographic inference. *Heredity (Edinb)* **116**: 349–350.

801 Patterson N, Price AL, Reich D (2006). Population structure and eigenanalysis. *PLoS Genet* **2**:
802      2074–2093.

803 Pfeifer B, Wittelsburger U, Ramos-Onsins SE, Lercher MJ (2014). PopGenome: An efficient
804      swiss army knife for population genomic analyses in R. *Mol Biol Evol* **31**: 1929–1936.

805 Pickrell JK, Pritchard JK (2012). Inference of population splits and mixtures from genome-
806      wide allele frequency data. *PLoS Genet* **8**: e1002967.

807 Poelstra JW, Vijay N, Bossu CM, Lantz H, Ryll B, Baglione V, *et al.* (2014). The genomic
808      landscape underlying phenotypic integrity in the face of gene flow in crows. *Science (80-*
809      *)* **344**: 1410–1414.

810 Prins P, de Ligt J, Tarasov A, Jansen RC, Cuppen E, Bourne PE (2015). Toward effective
811      software solutions for big biology. *Nat Biotech* **33**: 686–687.

812 Pritchard JK, Pickrell JK, Coop G (2010). The Genetics of Human Adaptation: Hard Sweeps,
813      Soft Sweeps, and Polygenic Adaptation. *Curr Biol* **20**: R208–R215.

814 Pritchard JK, Di Rienzo A (2010). Adaptation – not by sweeps alone. *Nat Rev Genet* **11**: 665–
815      667.

816 Pritchard JK, Stephens M, Donnelly P (2000). Inference of population structure using
817      multilocus genotype data. *Genetics* **155**: 945–959.

818 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, *et al.* (2007).
819      PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage
820      Analyses. *Am J Hum Genet* **81**: 559–575.

821    Puritz JB, Matz M V., Toonen RJ, Weber JN, Bolnick DI, Bird CE (2014). Demystifying the

822         RAD fad. *Mol Ecol* **23**: 5937–5942.

823    Raj A, Stephens M, Pritchard JK (2014). FastSTRUCTURE: Variational inference of

824         population structure in large SNP data sets. *Genetics* **197**: 573–589.

825    Rasmussen MD, Hubisz MJ, Gronau I, Siepel A (2014). Genome-Wide Inference of Ancestral

826         Recombination Graphs. *PLoS Genet* **10**.

827    Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO (2012). DELLY: Structural

828         variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* **28**:

829         333–339.

830    Rimmer A, Phan H, Mathieson I, Iqbal Z, Twigg SRF, Wilkie AOM, *et al.* (2014). Integrating

831         mapping-, assembly- and haplotype-based approaches for calling variants in clinical

832         sequencing applications. *Nat Genet* **46**: 912–918.

833    Rockman M V (2012). The QTN program and the alleles that matter for evolution: all that's

834         gold does not glitter. *Evolution (N Y)* **66**: 1–17.

835    Roux C, Fraisse C, Castric V, Vekemans X, Pogson GH, Bierne N (2014). Can we continue to

836         neglect genomic variation in introgression rates when inferring the history of speciation?

837         A case study in a Mytilus hybrid zone. *J Evol Biol* **27**: 1662–1675.

838    Roux C, Pauwels M, Ruggiero M-V, Charlesworth D, Castric V, Vekemans X (2013). Recent

839         and ancient signature of balancing selection around the S-locus in Arabidopsis halleri

840         and A. lyrata. *Mol Biol Evol* **30**: 435–47.

841    Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, Schaffner SF, *et al.* (2002).

842         Detecting recent positive selection in the human genome from haplotype structure. **419**.

843    Scheet P, Stephens M (2006). A fast and flexible statistical model for large-scale population

844         genotype data: applications to inferring missing genotypes and haplotypic phase. *Am J*

845         *Hum Genet* **78**: 629–44.

846    Schiffels S, Durbin R (2014). Inferring human population size and separation history from

847         multiple genome sequences. *Nat Genet* **46**: 919–25.

848    Schrider DR, Mendes FK, Hahn MW, Kern AD (2015). Soft shoulders ahead: Spurious

849         signatures of soft and partial selective sweeps result from linked hard sweeps. *Genetics*

850        **200**: 267–284.

851    Schubert M, Jónsson H, Chang D, Der Sarkissian C, Ermini L, Ginolhac A, *et al.* (2014).

852        Prehistoric genomes reveal the genetic foundation and cost of horse domestication. *Proc*

853        *Natl Acad Sci* **111**: 201416991.

854    Schwartz S, Kent W, Smit A (2003). Human–mouse alignments with BLASTZ. *Genome Res*:

855        103–107.

856    Shafer AB a., Wolf JBW, Alves PC, Bergström L, Bruford MW, Brännström I, *et al.* (2015).

857        Genomics and the challenging translation into conservation practice. *Trends Ecol Evol*

858        **30**: 78–87.

859    Sheehan S, Harris K, Song YS (2013). Estimating Variable Effective Population Sizes from

860        Multiple Genomes □: A Sequentially Markov. *Genetics* **194**: 647–662.

861    Soltis PS, Marchant DB, Van de Peer Y, Soltis DE (2015). Polyploidy and genome evolution

862        in plants. *Curr Opin Genet Dev* **35**: 119–125.

863    Stamatakis A (2014). RAxML version 8: A tool for phylogenetic analysis and post-analysis of

864        large phylogenies. *Bioinformatics* **30**: 1312–1313.

865    Tajima F (1989). Statistical method for testing the neutral mutation hypothesis by DNA

866        polymorphism. *Genetics* **123**: 585–95.

867    Takezaki N, Nei M, Tamura K (2010). POPTREE2: Software for constructing population

868        trees from allele frequency data and computing other population statistics with windows

869        interface. *Mol Biol Evol* **27**: 747–752.

870    Tang K, Thornton KR, Stoneking M (2007). A new approach for using genome scans to

871        detect recent positive selection in the human genome. *PLoS Biol* **5**: 1587–1602.

872    Vitti JJ, Grossman SR, Sabeti PC (2013). Detecting Natural Selection in Genomic Data. *Annu*

873        *Rev Genet* **47**: 97–120.

874    Wang M, Huang X, Li R, Xu H, Jin L, He Y (2014). Detecting recent positive selection with

875        high accuracy and reliability by conditional coalescent tree. *Mol Biol Evol* **31**: 3068–

876        3080.

877    Weber JN, Peterson BK, Hoekstra HE (2013). Discrete genetic modules are responsible for

878         complex burrow evolution in Peromyscus mice. *Nature* **493**: 402–5.

879    Wegmann D, Leuenberger C, Neuenschwander S, Excoffier L (2010). ABCtoolbox: a

880         versatile toolkit for approximate Bayesian computations. *BMC Bioinformatics* **11**: 116.

881    White BJ, Cheng C, Simard F, Costantini C, Besansky NJ (2010). Genetic association of

882         physically unlinked islands of genomic divergence in incipient species of Anopheles

883         gambiae. *Mol Ecol* **19**: 925–939.

884    Whitlock MC, Bronstein JL, Bruna EM, Ellison AM, Fox CW, McPeek MA, *et al.* (2015). A

885         Balanced Data Archiving Policy for Long-Term Studies. *Trends Ecol Evol* **xx**: 1–2.

886    Yang Z (2007). PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**:

887         1586–1591.

888    Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Others (2010). Common SNPs

889         explain a large proportion of the heritability for human height. *Nat Genet* **42**: 565–569.

890    Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS (2012). A high-performance

891         computing toolset for relatedness and principal component analysis of SNP data.

892         *Bioinformatics* **28**: 3326–3328.

893

894

895

896

897

898

899

900

901

902

903

904   Figures

905   Figure 1: Summary of methods assessing the demographic history of populations (middle

906         panel). References are provided in the main text.

907   Figure 2: Summary of methods dedicated to the detection of various signatures of selection in

908         the genome. In this simple example, the mutation on the left is under positive selection

909         in one population (red) but not the other (black). The mutation on the right is under

910         ancient balancing selection in the two populations.

911

912

913

31

## Assessing population structure and family relationships

PCA (SNPRElate, GenAbel in R)
STRUCTURE and fastSTRUCTURE
GENELAND
PLINK, VCFTOOLS, KING (relatedness)
Arlequin (AMOVA)

## Phylogenetic relationships between individuals/populations



TREEMIX, fastSTRUCTURE (admixture)
POPTREE, PopGenome, VCFTOOLS (differentiation)
Splitstree, SNPhylo, RAxML, PhyML, BEAST2 (individual phylogeny)
POPTREE, BEAST*, SNAPP (Species/Population trees)



## Model testing



Likelihood comparison: Migrate-n, IMa2 (for small datasets)
ABC methods: DIYABC, ms, fastsimcoal2, ABCToolbox, package abc in R.
Likelihood comparison using AFS (fastsimcoal2, dadi)

## Demographic parameters inference



SNAPP, BEAST*, IMa, Migrate-n, LAMARC (for small datasets)
ABC methods: DIYABC, ms, fastsimcoal2
Inference from allele frequency spectrum (AFS): fastsimcoal2, dadi
diCal, PSMC, MSMC (for whole genome resequencing)

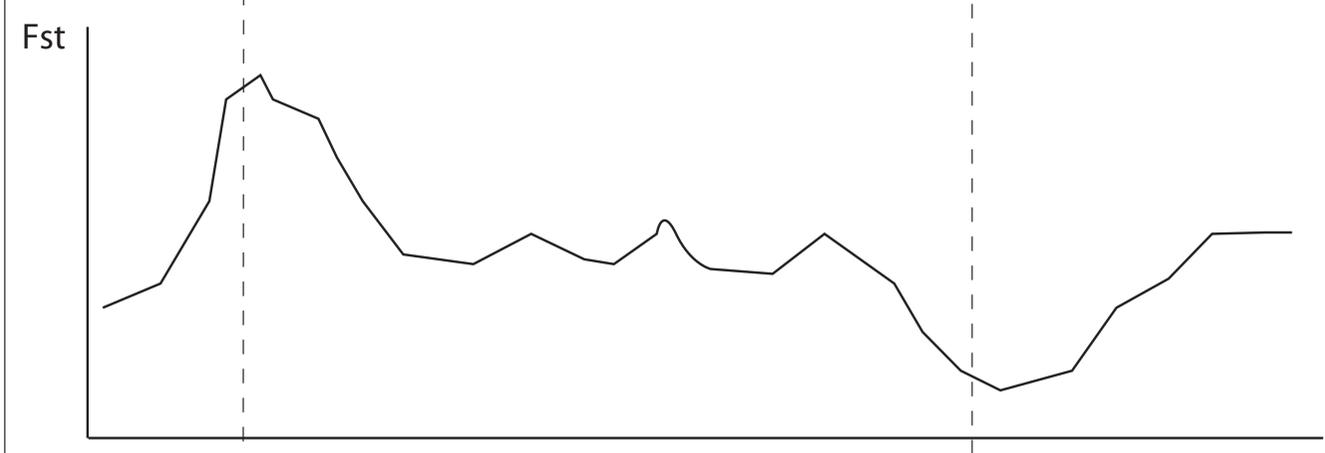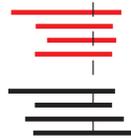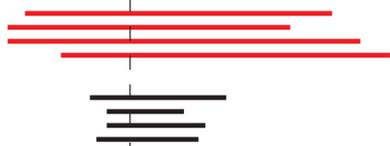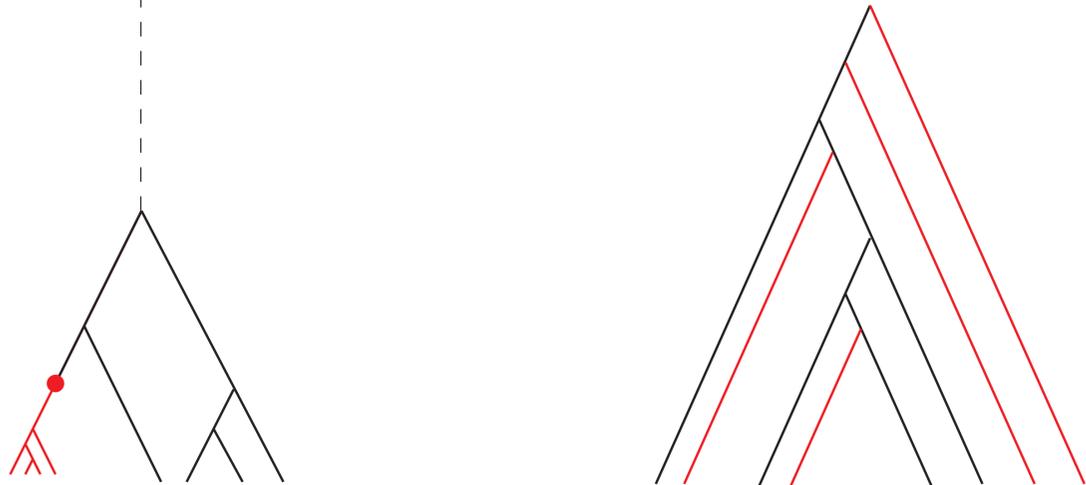| Positive selection | Balancing selection | Methods and tools available | |
|---|---|---|---|
| π | | PopGenome, VCFTOOLS, POPBAM<br><br>SweepFinder 2<br><br>Biopython<br><br>Popoolation (pooled data) | Diversity |
| Fst | | PopGenome, VCFTOOLS, Popoolation<br><br>Fst outlier methods: BAYPASS, BAYENV<br><br>Association with environment: BAYPASS, LFMM, BAYENV<br><br>Association with phenotype: GENABEL, PLINK | Differentiation/Association |
| | | Genome-wide LD: PLINK, VCFTOOLS<br><br>EHH and Rsb tests: rehh package, Sweep<br><br>LD clusters: LDna | Haplotype length and LD patterns |
| | | MSMS, PopGenome<br><br>ARGWeaver *<br><br>BALLET *<br><br>SCCT *<br><br>* Better suited for whole-genome resequencing | Coalescence and alleles history |