1   **Comparing methods for detecting multilocus adaptation with multivariate genotype-**

2   **environment associations**

3

4   Brenna R. Forester[1,2], Jesse R. Lasky[3], Helene H. Wagner[4], Dean L. Urban[1]

5

6   1 – Duke University, Nicholas School of the Environment, Durham, NC 27708, USA.

7   2 – Current location: Colorado State University, Department of Biology, Fort Collins, CO 80523,

8   USA.

9   3 – Pennsylvania State University, Department of Biology, University Park, PA 16802, USA.

10  4 – University of Toronto Mississauga, Department of Biology, Mississauga, ON, Canada

11

12  *Keywords*: constrained ordination, landscape genomics, natural selection, random forest,

13  redundancy analysis, simulations

14

15  *Corresponding Author*: Brenna R. Forester, Colorado State University, Department of Biology,

16  1878 Campus Delivery, Fort Collins, CO 80523, Phone: (970) 491-7011, Fax: (970) 491-0649,

17  Email: brenna.forester@colostate.edu

18

19  *Running title*: Detecting multilocus adaptation

1

**Abstract**

Identifying adaptive loci can provide insight into the mechanisms underlying local adaptation. Genotype-environment association (GEA) methods, which identify these loci based on correlations between genetic and environmental data, are particularly promising. Univariate methods have dominated GEA, despite the high dimensional nature of genotype and environment. Multivariate methods, which analyze many loci simultaneously, may be better suited to these data since they consider how sets of markers covary in response to environment. These methods may also be more effective at detecting adaptive processes that result in weak, multilocus signatures. Here, we evaluate four multivariate methods, and five univariate and differentiation-based approaches, using published simulations of multilocus selection. We found that Random Forest performed poorly for GEA. Univariate GEAs performed better, but had low detection rates for loci under weak selection. Constrained ordinations showed a superior combination of low false positive and high true positive rates across all levels of selection. These results were robust across the demographic histories, sampling designs, sample sizes, and levels of population structure tested. The value of combining detections from different methods was variable, and depended on study goals and knowledge of the drivers of selection. Reanalysis of genomic data from gray wolves highlighted the unique, covarying sets of adaptive loci that could be identified using redundancy analysis, a constrained ordination. Although additional testing is needed, this study indicates that constrained ordinations are an effective means of detecting adaptation, including signatures of weak, multilocus selection, providing a powerful tool for investigating the genetic basis of local adaptation.

**Introduction**

Analyzing genomic data for loci underlying local adaptation has become common practice in evolutionary and ecological studies (Hoban *et al.*, 2016). These analyses can help identify mechanisms of local adaptation and inform management decisions for agricultural, natural resources, and conservation applications. Genotype-environment association (GEA) approaches are particularly promising for detecting these loci (Rellstab *et al.* 2015). Unlike differentiation outlier methods, which identify loci with strong allele frequency differences among populations, GEA approaches identify adaptive loci based on associations between genetic data and environmental variables hypothesized to drive selection. Benefits of GEA include the option of using individual-based (as opposed to population-based) sampling and the ability to make explicit links to the ecology of organisms by including relevant predictors. The inclusion of predictors can also improve power and allows for the detection of selective events that do not produce high genetic differentiation among populations (De Mita *et al.*, 2013; de Villemereuil *et al.*, 2014; Rellstab *et al.*, 2015).

Univariate statistical methods have dominated GEA since their first appearance (Mitton *et al.*, 1977). These methods test one locus and one predictor variable at a time, and include generalized linear models (e.g. Joost et al. 2007; Stucki et al. 2016), variations on linear mixed effects models (e.g. Coop et al. 2010; Frichot et al. 2013; Yoder et al. 2014; Lasky et al. 2014), and non-parametric approaches (e.g. partial Mantel, Hancock et al. 2011). While these methods perform well, they can produce elevated false positive rates in the absence of correction for multiple comparisons, an issue of increased importance with large genomic data sets. Corrections such as Bonferroni can be overly conservative (potentially removing true positive detections), while alternative correction methods, such as false discovery rate (FDR, Benjamini & Hochberg 1995), rely on an assumption of a null distribution of $p$-values, which may often be violated for empirical data sets. While these issues should not discourage the use of univariate methods (though corrections should be chosen carefully, see François *et al.* (2016) for a recent overview), other analytical approaches may be better suited to the high dimensionality of modern genomic data sets.

In particular, multivariate approaches, which analyze many loci simultaneously, are well suited to data sets comprising hundreds of individuals sampled at many thousands of genetic

3

71    markers. Compared to univariate methods, these approaches are thought to more effectively

72    detect multilocus selection since they consider how groups of markers covary in response to

73    environmental predictors (Rellstab et al. 2015). This is important because many adaptive

74    processes are expected to result in weak, multilocus molecular signatures due to selection on

75    standing genetic variation, recent/contemporary selection that has not yet led to allele fixation,

76    and conditional neutrality (Yeaman & Whitlock, 2011; Le Corre & Kremer, 2012; Savolainen *et*

77    *al.*, 2013; Tiffin & Ross-Ibarra, 2014). Identifying the relevant patterns (e.g., coordinated shifts

78    in allele frequencies across many loci) that underlie these adaptive processes is essential to both

79    improving our understanding of the genetic basis of local adaptation, and advancing applications

80    of these data for management, such as conserving the evolutionary potential of species

81    (Savolainen *et al.*, 2013; Harrisson *et al.*, 2014; Lasky *et al.*, 2015). While multivariate methods

82    may, in principle, be better suited to detecting these shared patterns of response, they have not

83    yet been tested on common data sets simulating multilocus adaptation, limiting confidence in

84    their effectiveness on empirical data.

85           Here we evaluate a set of these methods, using published simulations of multilocus

86    selection (Lotterhos & Whitlock, 2014, 2015). We compare power using empirical *p*-values, and

87    evaluate false positive rates based on cutoffs used in empirical studies. We follow up with a test

88    of three of these methods on their ability to detect weak multilocus selection, as well as an

89    assessment of the common practice of combining detections across multiple tests. We investigate

90    the effects of correction for population structure in one ordination method, and follow up with an

91    application of this test to an empirical data set from gray wolves. We find that the constrained

92    ordinations we tested maintain the best balance of true and false positive rates across a range of

93    demographies, sampling designs, sample sizes, and selection levels, and can provide unique

94    insight into the processes driving selection and the multilocus architecture of local adaptation.

95

96    **Methods**

97    Multivariate approaches to GEA:

98    Multivariate statistical techniques, including ordinations such as principal components analysis

99    (PCA), have been used to analyze genetic data for over fifty years (Cavalli-Sforza, 1966).

100   Indirect ordinations like PCA (which do not use predictors) use patterns of association within

101  genetic data to find orthogonal axes that fully decompose the genetic variance. Constrained

102  ordinations extend this analysis by restricting these axes to combinations of supplied predictors

103  (Jombart *et al.*, 2009; Legendre & Legendre, 2012). When used as a GEA, a constrained

104  ordination is essentially finding orthogonal sets of loci that covary with orthogonal multivariate

105  environmental patterns. By contrast, a univariate GEA is testing for single locus relationships

106  with single environmental predictors. The use of constrained ordinations in GEA goes back as

107  far as Mulley et al. (1979), with more recent applications to genomic data sets in Lasky et al.

108  (2012), Forester et al. (2016), and Brauer et al. (2016). In this analysis, we test two promising

109  constrained ordinations, redundancy analysis (RDA) and distance-based redundancy analysis

110  (dbRDA). We also test an extension of RDA that uses a preliminary step of summarizing the

111  genetic data into sets of covarying markers (Bourret *et al.*, 2014). We do not include canonical

112  correspondence analysis, a constrained ordination that is best suited to modeling unimodal

113  responses, although this method has been used to analyze microsatellite data sets (e.g. Angers et

114  al. 1999; Grivet et al. 2008).

115  Random Forest (RF) is a machine learning algorithm that is designed to identify structure

116  in complex data and generate accurate predictive models. It is based on classification and

117  regression trees (CART), which recursively partition data into response groups based on splits in

118  predictors variables. CART models can capture interactions, contingencies, and nonlinear

119  relationships among variables, differentiating them from linear models (De'ath & Fabricius,

120  2000). RF reduces some of the problems associated with CART models (e.g. overfitting and

121  instability) by building a "forest" of classification or regression trees with two layers of

122  stochasticity: random bootstrap sampling of the data, and random subsetting of predictors at each

123  node (Breiman, 2001). This provides a built-in assessment of predictive accuracy (based on data

124  left out of the bootstrap sample) and variable importance (based on the change in accuracy when

125  covariates are permuted). For GEA, variable importance is the focal statistic, where the predictor

126  variables used at each split in the tree are molecular markers, and the goal is to sort individuals

127  into groups based on an environmental category (classification) or to predict home

128  environmental conditions (regression). Markers with high variable importance are best able to

129  sort individuals or predict environments. RF has been used in a number of recent GEA and

130  GWAS studies (e.g. Holliday et al. 2012; Brieuc et al. 2015; Pavey et al. 2015; Laporte et al.

131  2016), but has not yet been tested in a GEA simulation framework.

132  We compare these multivariate methods to the two differentiation-based and three

133  univariate GEA methods tested by Lotterhos & Whitlock (2015): the $X^TX$ statistic from Bayenv2

134  (Günther & Coop, 2013), PCAdapt (Duforet-Frebourg *et al.*, 2014), latent factor mixed models

135  (LFMM, Frichot *et al.* 2013), and two GEA-based statistics (Bayes factors and Spearman's ρ)

136  from Bayenv2. We also include generalized linear models (GLM), a regression-based GEA that

137  does not use a correction for population structure.

138

139  GEA implementation:

140  *Constrained ordinations*:

141  We tested RDA and dbRDA as implemented by Forester et al. (2016). RDA is a two-step process

142  in which genetic and environmental data are analyzed using multivariate linear regression,

143  producing a matrix of fitted values. Then PCA of the fitted values is used to produce canonical

144  axes, which are linear combinations of the predictors. We centered and scaled genotypes for

145  RDA (i.e., mean = 0, s = 1; see Jombart *et al.* 2009 for a discussion of scaling genetic data for

146  ordinations). Distance-based redundancy analysis is similar to RDA but allows for the use of

147  non-Euclidian dissimilarity indices. Whereas RDA can be loosely considered as a PCA

148  constrained by predictors, dbRDA is analogous to a constrained principal coordinate analysis

149  (PCoA, or a PCA on a non-Euclidean dissimilarity matrix). For dbRDA, we calculated the

150  distance matrix using Bray-Curtis dissimilarity (Bray & Curtis, 1957), which quantifies the

151  dissimilarity among individuals based on their multilocus genotypes (equivalent to one minus the

152  proportion of shared alleles between individuals). For both methods, SNPs are modeled as a

153  function of predictor variables, producing as many constrained axes as predictors. We identified

154  outlier loci on the constrained ordination axes based on the "locus score", which represent the

155  coordinates/loading of each locus in the ordination space. We use *rda* for RDA and *capscale* for

156  dbRDA in the vegan, v. 2.3-5 package (Oksanen *et al.*, 2013) in R v. 3.2.5 (R Development Core

157  Team, 2015) for this and all subsequent analyses.

158     *Redundancy analysis of components:*

159     This method, described by Bourret *et al.* (2014), differs from the approaches described above in

160     using a preliminary step that summarizes the genotypes into sets of covarying markers, which are

161     then used as the response in RDA. The idea is to identify from these sets of covarying loci only

162     the groups that are most strongly correlated with environmental predictors. We began by

163     ordinating SNPs into principal components (PCs) using *prcomp* in R on the scaled data,

164     producing as many axes as individuals. Following Bourret et al. (2014), we used parallel analysis

165     (Horn, 1965) to determine how many PCs to retain. Parallel analysis is a Monte Carlo approach

166     in which the eigenvalues of the observed components are compared to eigenvalues from

167     simulated data sets that have the same size as the original data. We used 1,000 random data sets

168     to generate the distribution under the null hypothesis and retained components with eigenvalues

169     greater than the 99[th] percentile of the eigenvalues of the simulated data (i.e., a significance level

170     of 0.01), using the hornpa package, v. 1.0 (Huang, 2015).

171        Next, we applied a varimax rotation to the PC axes, which maximizes the correlation

172     between the axes and the original variables (in this case, the SNPs). Note that once a rotation is

173     applied to the PC axes, they are no longer "principal" components (i.e. axes associated with an

174     eigenvalue/variance), but simply components. We then used the retained components as

175     dependent variables in RDA, with environmental variables used as predictors. Next, components

176     that were significantly correlated with the constrained axis were retained. Significance was based

177     on a cutoff (alpha = 0.05) corrected for sample sizes using a Fisher transformation as in Bourret

178     et al. (2014). Finally, SNPs were correlated with these retained components to determine

179     outliers. We call this approach redundancy analysis of components (cRDA).

180

181     *Random Forest:*

182     The Random Forest approach implemented here builds off of work by Goldstein et al. (2010),

183     Holliday et al. (2012), and Brieuc et al. (2015). This three-step approach is implemented

184     separately for each predictor variable. The environmental variable used in this study was

185     continuous, so RF models were built as regression trees. For categorical predictors (e.g. soil

186     type) classification trees would be used, which require a different parameterization (important

187     recommendations for this case are provided in Goldstein et al. 2010).

188    First, we tuned the two main RF parameters, the number of trees (*ntrees*) and the number

189    of predictors sampled per node (*mtry*). We tested a range of values for *ntrees* in a subset of the

190    simulations, and found that 10,000 trees were sufficient to stabilize variable importance (note

191    that variable importance requires a larger number of trees for convergence than error rates,

192    Goldstein et al. 2010). We used the default value of *mtry* for regression (number of predictors/3,

193    equivalent to ~3,330 SNPs in this case) after checking that increasing *mtry* did not substantially

194    change variable importance or the percent variance explained. In a GEA/GWAS context, larger

195    values of *mtry* reduce error rates, improve variable importance estimates, and lead to greater

196    model stability (Goldstein *et al.* 2010).

197    Because RF is a stochastic algorithm, it is best to use multiple runs, particularly when

198    variable importance is the parameter of interest (Goldstein *et al.*, 2010). We begin by building

199    three full RF models using all SNPs as predictors, saving variable importance as mean decrease

200    in accuracy for each model. Next, we sampled variable importance from each run with a range of

201    cutoffs, pulling the most important 0.5%, 1.0%, 1.5%, and 2.0% of loci. These values correspond

202    to approximately 50/100/150/200 loci that have the highest variable importance. For each cutoff,

203    we then created three additional RF models, using the average percent variance explained across

204    runs to determine the best starting number of important loci for step 3. This step removes clearly

205    unimportant loci from further consideration (i.e. "sparsity pruning", Goldstein et al. 2010).

206    Third, we doubled the best starting number of loci from step 2; this is meant to

207    accommodate loci that may have low marginal effects (Goldstein et al. 2010). We then built

208    three RF models with these loci, and recorded the mean variance explained. We removed the

209    least important locus in each model, and recalculated the RF models and mean variance

210    explained. This procedure continues until two loci remain. The set of loci that explain the most

211    variance are the final candidates. Candidates are then combined across runs to identify outliers.

212

213    *Differentiation-based and univariate GEA methods:*

214    For the two differentiation-based and the Bayenv2-based GEA methods, we compared power

215    directly from the results provided in Lotterhos & Whitlock (2015). PCAdapt is a differentiation-

216    based method that concurrently identifies outlier loci and population structure using latent factors

217    (Duforet-Frebourg *et al.*, 2014). The $X^TX$ statistic from Bayenv2 (Günther & Coop, 2013) is an

218    $F_{ST}$ analog that uses a covariance matrix to control for population structure. The two Bayenv2

219    GEA statistics (Bayes factors and Spearman's ρ) also use the covariance matrix to control for

220    population structure, while identifying candidate loci based on log-transformed Bayes factors

221    and nonparametric correlations, respectively. Details on these methods and their implementation

222    are provided in Lotterhos & Whitlock (2015).

223         We reran latent factor mixed models, a GEA approach that controls for population

224    structure using latent factors, using updated parameters as recommended by the authors (O.

225    François, pers. comm.). We tested values of *K* (the number of latent factors) ranging from one to

226    25 using a sparse nonnegative matrix factorization algorithm (Frichot *et al.*, 2014), implemented

227    as function *snmf* in the package LEA, v. 1.2.0 (Frichot & François, 2015). We plotted the cross-

228    entropy values and selected *K* based on the inflection point in these plots; when the inflection

229    point was not clear, we used the value where additional cross-entropy loss was minimal. We

230    parameterized LFMM models with this best estimate of *K*, and ran each model ten times with

231    5,000 iterations and a burnin of 2,500. We used the median of the squared z-scores to rank loci

232    and calculate a genomic inflation factor (GIF) to assess model fit (Frichot & François, 2015;

233    François *et al.*, 2016). The GIF is used to correct for inflation of z-scores at each locus, which

234    can occur when population structure or other confounding factors are not sufficiently accounted

235    for in the model (François *et al.* 2016). The GIF is calculated by dividing the median of the

236    squared z-scores by the median of the chi-squared distribution. We used the LEA and qvalue, v.

237    2.2.2 (Storey *et al.*, 2015) packages in R. Full K and GIF results are presented in Table S1.

238    Finally, we ran generalized linear models (GLM) on individual allele counts using a binomial

239    family and logistic link function for comparison with LFMM; GIF results are presented in Table

240    S1.

241

242    Simulations:

243    We used a subset of simulations published by Lotterhos & Whitlock (2014, 2015). Briefly, four

244    demographic histories are represented in these data, each with three replicated environmental

245    surfaces (Fig. S1): an equilibrium island model (IM), equilibrium isolation by distance (IBD),

246    and nonequilibrium isolation by distance with expansion from one (1R) or two (2R) refugia. In

247    all cases, demography was independent of selection strength, which is analogous to simulating

248    soft selection (Lotterhos & Whitlock, 2014). Haploid, biallelic SNPs were simulated

249    independently, with 9,900 neutral loci and 100 under selection. Note that haploid SNPs will yield

250    half the information content of diploid SNPs (Lotterhos & Whitlock 2015). The mean of the

251    environmental/habitat parameter had a selection coefficient equal to zero and represented the

252    background across which selective habitat was patchily distributed (Fig. S1). Selection

253    coefficients represent a proportional increase in fitness of alleles in response to habitat, where

254    selection is increasingly positive as the environmental value increases from the mean, and

255    increasingly negative as the value decreases from the mean (Lotterhos & Whitlock 2014, Fig.

256    S1). This landscape emulates a weak cline, with a north-south trend in the selection surface. Of

257    the 100 adaptive loci, most were under weak selection. For the IBD scenarios, selection

258    coefficients were 0.001 for 40 loci, 0.005 for 30 loci, 0.01 for 20 loci, and 0.1 for 10 loci. For the

259    1R, 2R, and IM scenario, selection coefficients were 0.005 for 50 loci, 0.01 for 33 loci, and 0.1

260    for 17 loci. Note that realized selection varied across demographies, so results across

261    demographic histories are not directly comparable (Lotterhos & Whitlock 2015).

262            We used the following sampling strategies and sample sizes from Lotterhos & Whitlock

263    (2015): random, paired, and transect strategies, with 90 demes sampled, and 6 or 20 individuals

264    sampled per deme. Paired samples (45 pairs) were designed to maximize environmental

265    differences between locations while minimizing geographic distance; transects (nine transects

266    with ten locations) were designed to maximize environmental differences at transect ends

267    (Lotterhos & Whitlock 2015). Overall, we used 72 simulations for testing. We assessed trend in

268    neutral loci using linear models of allele frequencies within demes as a function of coordinates.

269    We evaluated the strength of local adaptation using linear models of allele frequencies within

270    demes as a function of environment. Note that the Lotterhos & Whitlock (2014, 2015)

271    simulations assigned SNP genotypes to individuals within a population sequentially (i.e., the first

272    few individuals would all get the same allele until its target frequency was reached, the

273    remaining individuals would get the other allele). This creates artifacts (e.g., artificially low

274    observed heterozygosity) and may affect statistical error rates when subsampling individuals or

275    performing analyses at the individual level. As recommended by K. Lotterhos (pers. comm.), we

276    avoided these problems by randomizing allele counts for each SNP among individuals within

277     each population. The habitat surface, which imposed a continuous selective gradient on non-

278     neutral loci, was used as the environmental predictor.

279

280     <u>Evaluation statistics</u>:

281     In order to equitably compare power (true positive detections out of the number of loci under

282     selection) across these methods, we calculated empirical *p*-values using the method of Lotterhos

283     & Whitlock (2015). In this approach, we first built a null distribution based on the test statistics

284     of all neutral loci, and then generated a *p*-value for each selected locus based on its cumulative

285     frequency in the null distribution. We then converted empirical *p*-values to *q*-values to assess

286     significance, using the same *q*-value cutoff (0.01) as Lotterhos & Whitlock (2015). We used

287     code provided by K. Lotterhos to calculate empirical *p*-values (code provided in Supplemental

288     Information.

289           Because false positive rates (FPRs) are not very informative for empirical *p*-values (rates

290     are universally low, see Lotterhos & Whitlock 2015 for a discussion), we applied cutoffs (e.g.

291     thresholds for statistical significance) to assess both true and false positive rates across methods.

292     While power is important, determining FPRs is also an essential component of assessing method

293     performance, since high power achieved at the cost of high FPRs is problematic. Because cutoffs

294     differ across methods, we tested a range of commonly used thresholds for each method and

295     chose the approach that performed the best (i.e., best balance of TPR and FPR). Note that cutoffs

296     can be adjusted for empirical studies based on research goals and the tolerance for TP and FP

297     detections. For each cutoff tested, we calculated the TPR as the number of correct positive

298     detections out of the number possible, and the FPR as the number of incorrect positive detections

299     out of 9900 possible. For the main text, we present results from the best cutoff for each method;

300     full results for all cutoffs tested are presented in the Supplemental Information. For constrained

301     ordinations (RDA and dbRDA) we identified outliers as SNPs with a locus score +/- 2.5 and 3

302     SD from the mean score of each constrained axis. For cRDA, we used cutoffs for SNP-

303     component correlations of alpha = 0.05, 0.01, and 0.001, corrected for sample sizes using a

304     Fisher transformation as in Bourret et al. (2014). For GLM and LFMM, we compared two

305     Bonferroni-corrected cutoffs (0.05 and 0.01) and a FDR cutoff of 0.1.

11

306     Weak selection:

307     We compared the best-performing multivariate methods (RDA, dbRDA, and cRDA) for their

308     ability to detect signals of weak selection (s = 0.005 and s = 0.001). All tests were performed as

309     described above, after removing loci under strong (s = 0.1) and moderate (s = 0.01) selection

310     from the simulation data sets. The number of loci under selection in these cases ranged from 43

311     to 76.

312

313     Combining detections:

314     We compared the effects of combining detections (i.e., looking for overlap) using cutoff results

315     from two of the best-performing methods, RDA and LFMM. We also included a scenario in

316     which a second, uninformative predictor (the x-coordinate of each individual) is included in the

317     RDA and LFMM tests. This predictor is analogous to including an environmental variable

318     hypothesized to drive selection that covaries with longitude.

319

320     Correction for population structure in RDA:

321     To determine how explicit modeling of population structure affects the performance of the best-

322     performing multivariate method, RDA, we accounted for population structure using three

323     approaches: (1) partialling out significant spatial eigenvectors not correlated with the habitat

324     predictor, (2) partialling out all significant spatial eigenvectors, and (3) partialling out ancestry

325     coefficients. The spatial eigenvector procedure uses Moran eigenvector maps (MEM) as spatial

326     predictors in a partial RDA. MEMs provide a decomposition of the spatial relationships among

327     sampled locations based on a spatial weighting matrix (Dray *et al.*, 2006). We used spatial

328     filtering to determine which MEMs to include in the partial analyses (Dray *et al.*, 2012). Briefly,

329     this procedure begins by applying a principal coordinate analysis (PCoA) to the genetic distance

330     matrix, which we calculated using Bray-Curtis dissimilarity. We used the broken-stick criterion

331     (Legendre & Legendre, 2012) to determine how many genetic PCoA axes to retain. Retained

332     axes were used as the response in a full RDA, where the predictors included all MEMs. Forward

333     selection (Blanchet *et al.*, 2008) was used to reduce the number of MEMs, using the full RDA

334     adjusted $R^2$ statistic as the threshold. In the first approach, retained MEMs that were significantly

335     correlated with environmental predictors were removed (alpha = 0.05/number of MEMs), and the

336    remaining set of significant MEMs were used as conditioning variables in RDA. Note that this

337    approach will be liberal in removing MEMs correlated with environment. In the second

338    approach, all significant MEMs were used as conditioning variables, the most conservative use

339    of MEMs. We used the spdep, v. 0.6-9 (Bivand *et al.*, 2013) and adespatial, v. 0.0-7 (Dray *et al.*,

340    2016) packages to calculate MEMs. For the third approach, we used individual ancestry

341    coefficients as conditioning variables. We used function *snmf* in the LEA package to estimate

342    individual ancestry coefficients, running five replicates using the best estimate of $K$, and

343    extracting individual ancestry coefficients from the replicate with the lowest cross-entropy.

344

345    Empirical data set:

346    To provide an example of the use and interpretation of RDA as a GEA, we reanalyzed data from

347    94 North American gray wolves (*Canis lupus*) sampled across Canada and Alaska at 42,587

348    SNPs (Schweizer *et al.*, 2016). These data show similar global population structure to the

349    simulations analyzed here: wolf data Fst = 0.09; average simulation Fst = 0.05. We reduced the

350    number of environmental covariates originally used by Schweizer *et al.* (2016) from 12 to eight

351    to minimize collinearity among them (e.g., |r| < 0.7). One predictor, land cover, was removed

352    because the distribution of cover types was heavily skewed toward two of the ten types. Missing

353    data levels were low (3.06%). Because RDA requires complete data frames, we imputed missing

354    values by replacing them with the most common genotype across individuals. We identified

355    candidate adaptive loci as SNPs loading +/- 3 SD from the mean loading of significant RDA axes

356    (significance determined by permutation, $p < 0.05$). We then identified the covariate most

357    strongly correlated with each candidate SNP (i.e., highest correlation coefficient), to group

358    candidates by potential driving environmental variables. Annotated code for this example is
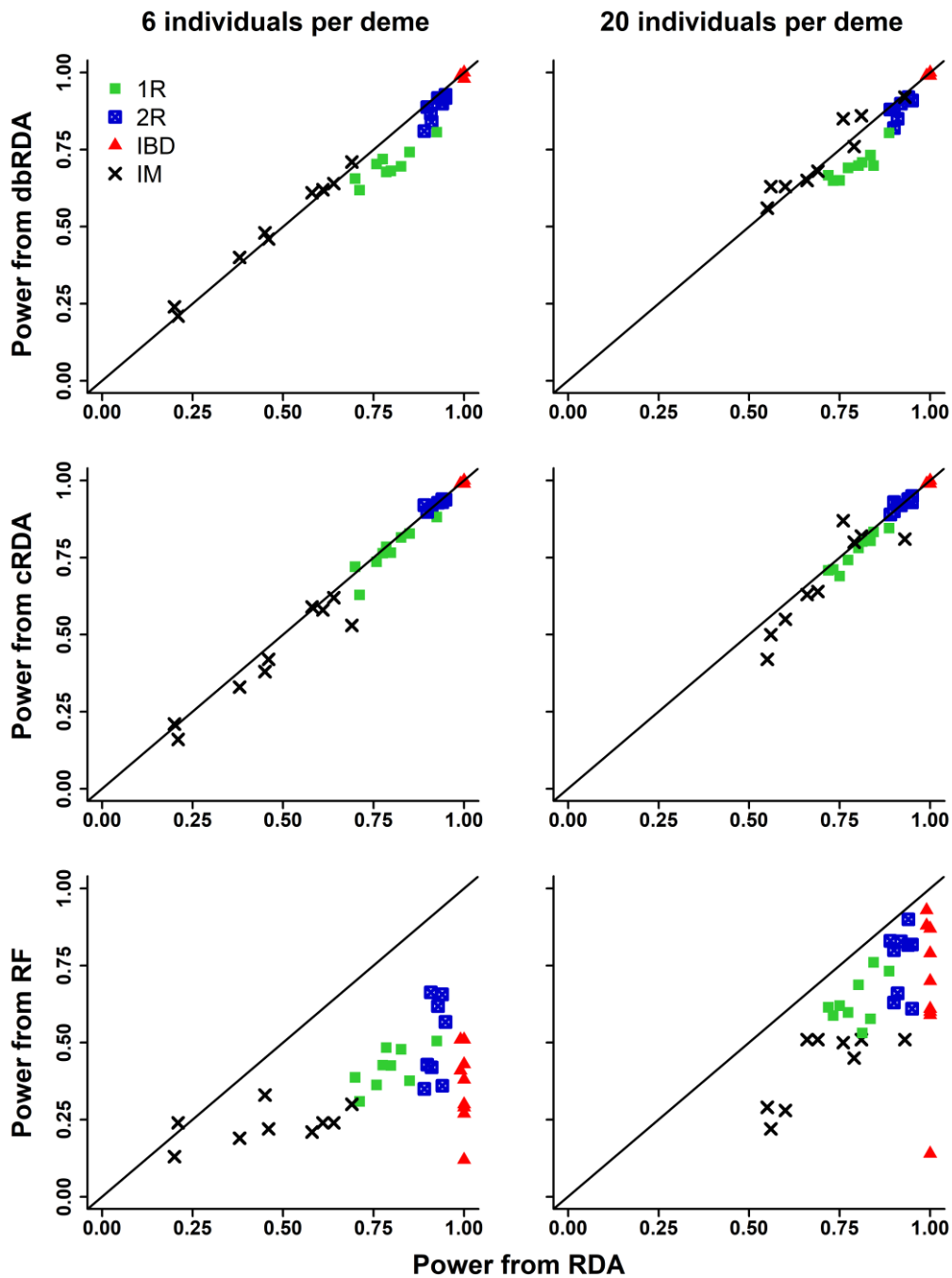
359    provided in the Supplementary Information.

360 **Results**

361 Empirical *p*-value results:

362 Power across the three ordination techniques was comparable, while power for RF was relatively

363 low (Fig. 1). Ordinations performed best in IBD, 1R, and 2R demographies, with the larger

364 sample size improving power for the IM demography. Within ordination techniques, RDA and

365 cRDA had slightly higher detection rates compared to dbRDA; subsequent comparisons are
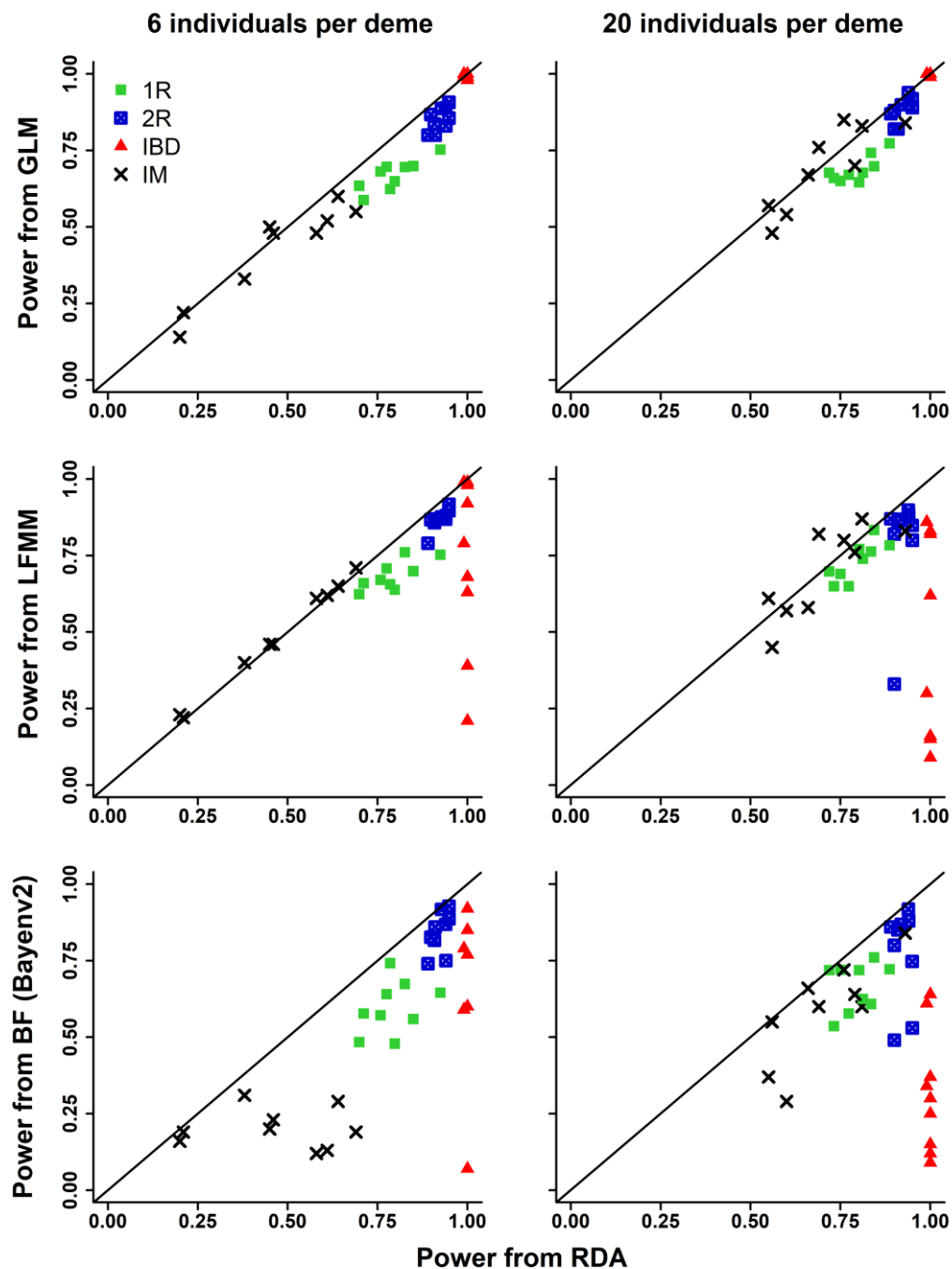
366 made using RDA results.

367 Except for a few cases in the IM demography, the power of RDA was generally higher

368 than univariate GEAs (Fig. 2). Of the univariate methods, GLM had the highest overall power,

369 while LFMM had reduced power for the IBD demography. Power from the Bayes Factor

370 (Bayenv2) was generally lower than RDA across all demographies. Finally, RDA had overall

371 higher power than the two differentiation-based methods (Fig. 3), with the exception of the IBD

372 demography, where power was high for all methods.

373 Among the methods with the highest overall power, all performed well at detecting loci

374 under strong selection (Fig. 4 and S2). Detection rates for loci under moderate and weak

375 selection were highest for ordination methods, with RDA and cRDA having the overall highest

376 detection rates. Detection of moderate and weakly selected loci was lower and more variable for

377 univariate methods, especially LFMM, where detection was dependent on demography and
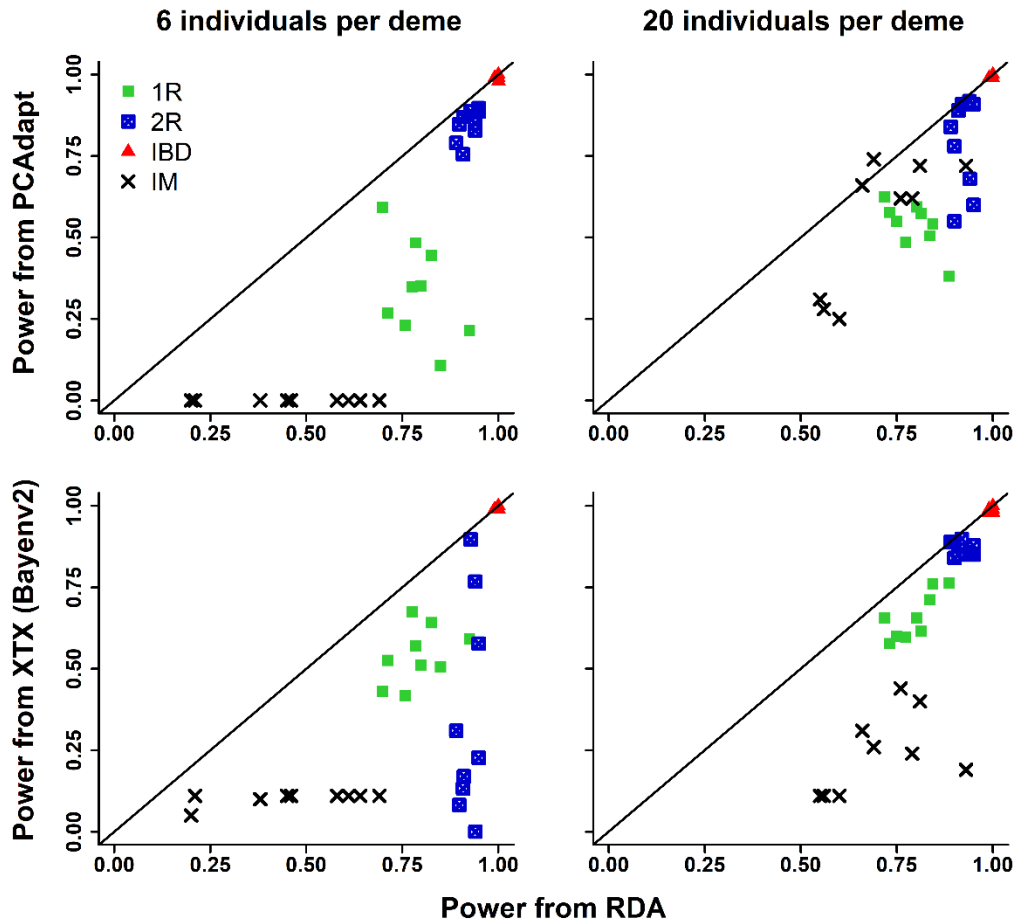
378 sampling scheme.

14

379

**Figure 1.** Comparison of power (from empirical *p*-values) from RDA (x-axis) and three other multivariate GEAs (y-axes, rows) for two sample sizes (columns). Points reflect demographies: 1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model. Some variation within demographies comes from sampling design.

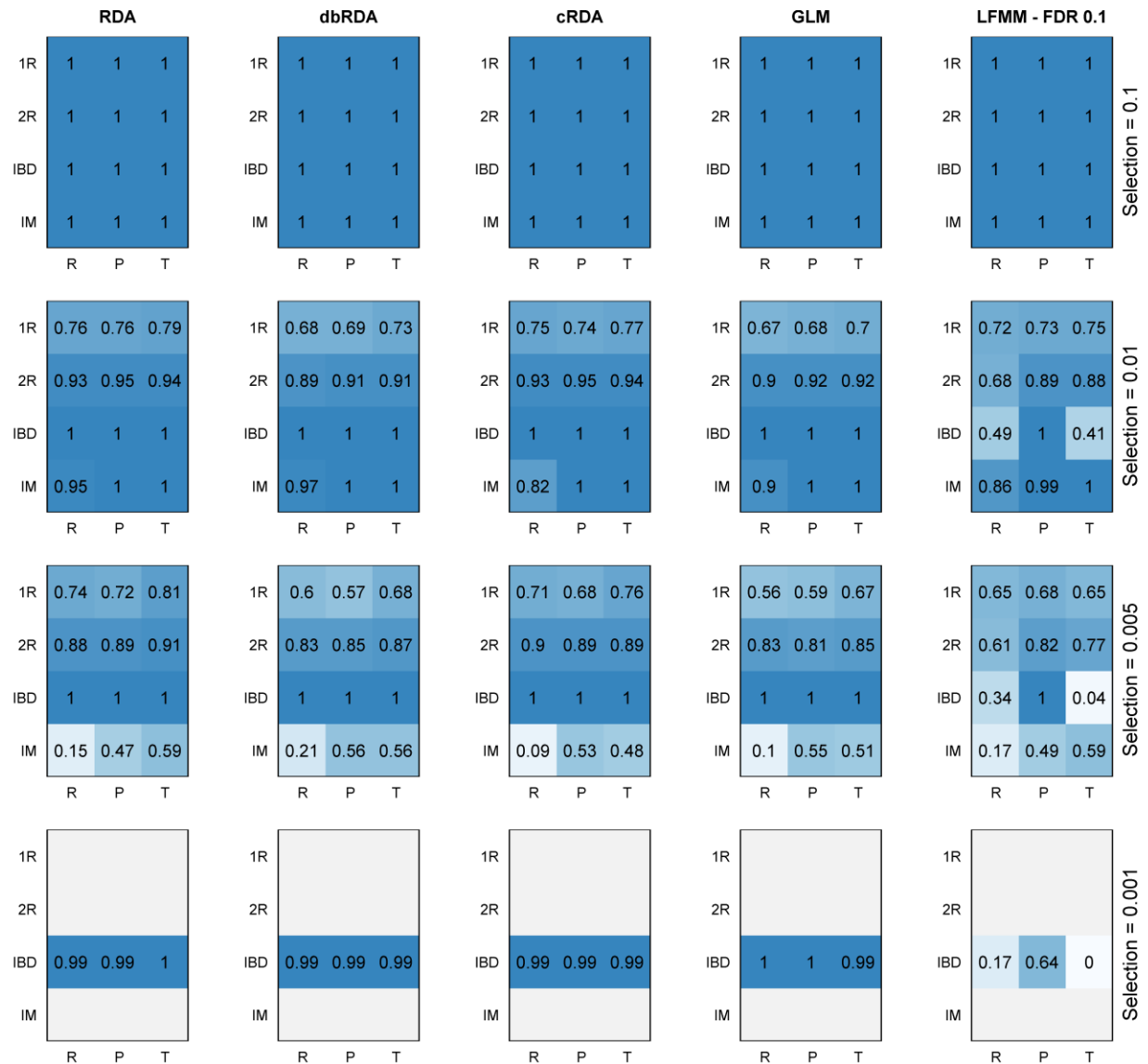**Figure 2.** Comparison of power (from empirical *p*-values) from RDA (x-axis) and three univariate GEAs (y-axes, rows) for two sample sizes (columns). Points reflect demographies: 1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model. Some variation within demographies comes from sampling design.

**Figure 3.** Comparison of power (from empirical *p*-values) from RDA (x-axis) and two differentiation-based outlier detection methods (y-axes, rows) for two sample sizes (columns). Points reflect demographies: 1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model. Some variation within demographies comes from sampling design.

17

**Figure 4.** Average power (from empirical *p*-values) for different levels of selection (rows) from five methods (columns) using a sample size of 20 individuals per deme. Each method shows results for different sampling strategies (R = random, P = pairs, T = transects) and demographies (1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model). Only the IBD demography included very weak selection (s=0.001).
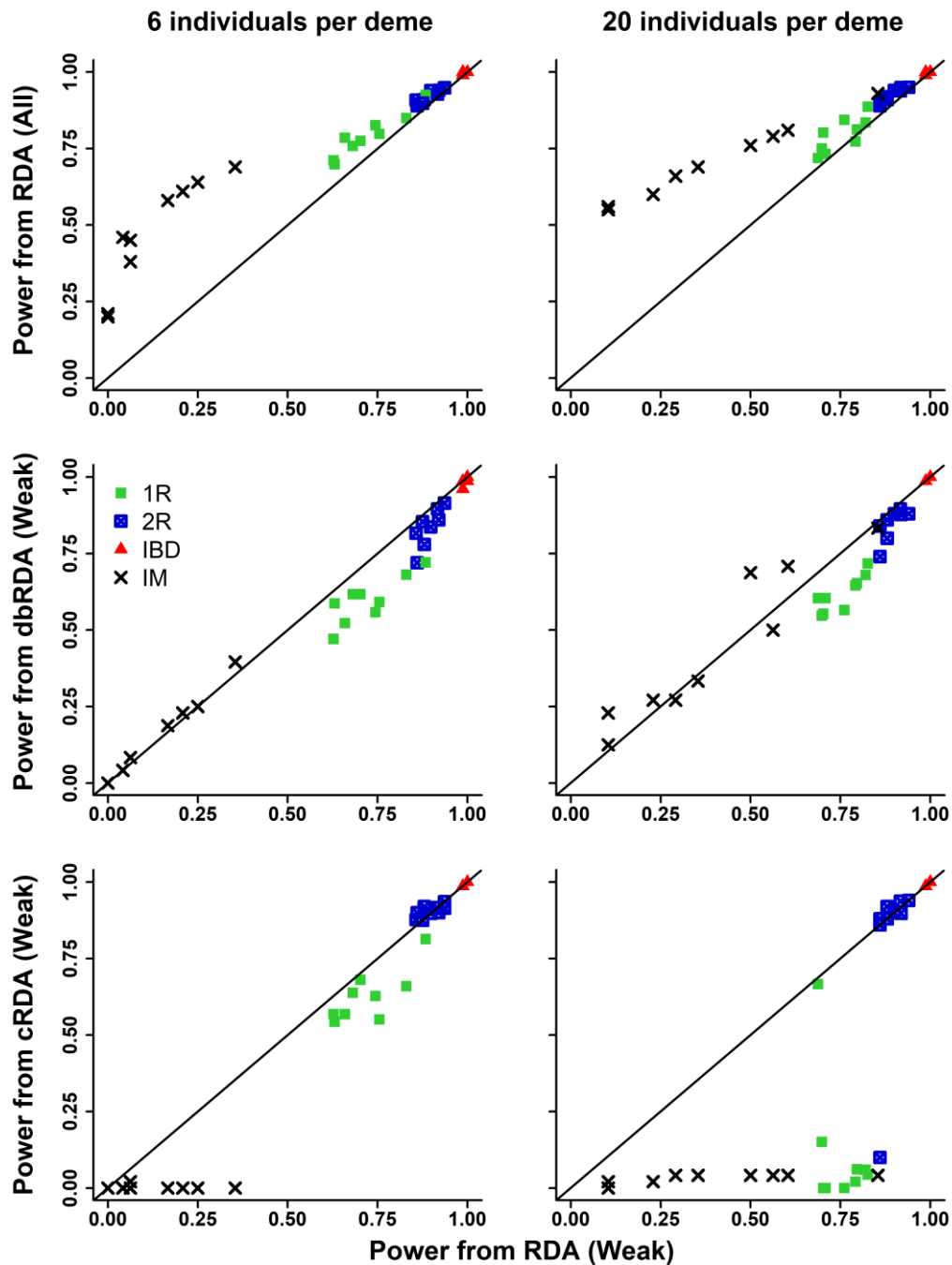
18

401  Weak selection:

402  We compared the three ordination methods for their power to detect only weak loci in the

403  simulations (Fig. 5). Power from RDA was higher when all selected loci were included,

404  especially for the IM demography. Power using only weakly selected loci was comparable

405  between RDA and dbRDA, with power slightly higher for RDA in most cases. cRDA was

406  comparable to RDA for the IBD and 2R demographies, but had very low to no power in the IM

407  demography, and the 1R demography with the larger sample size.
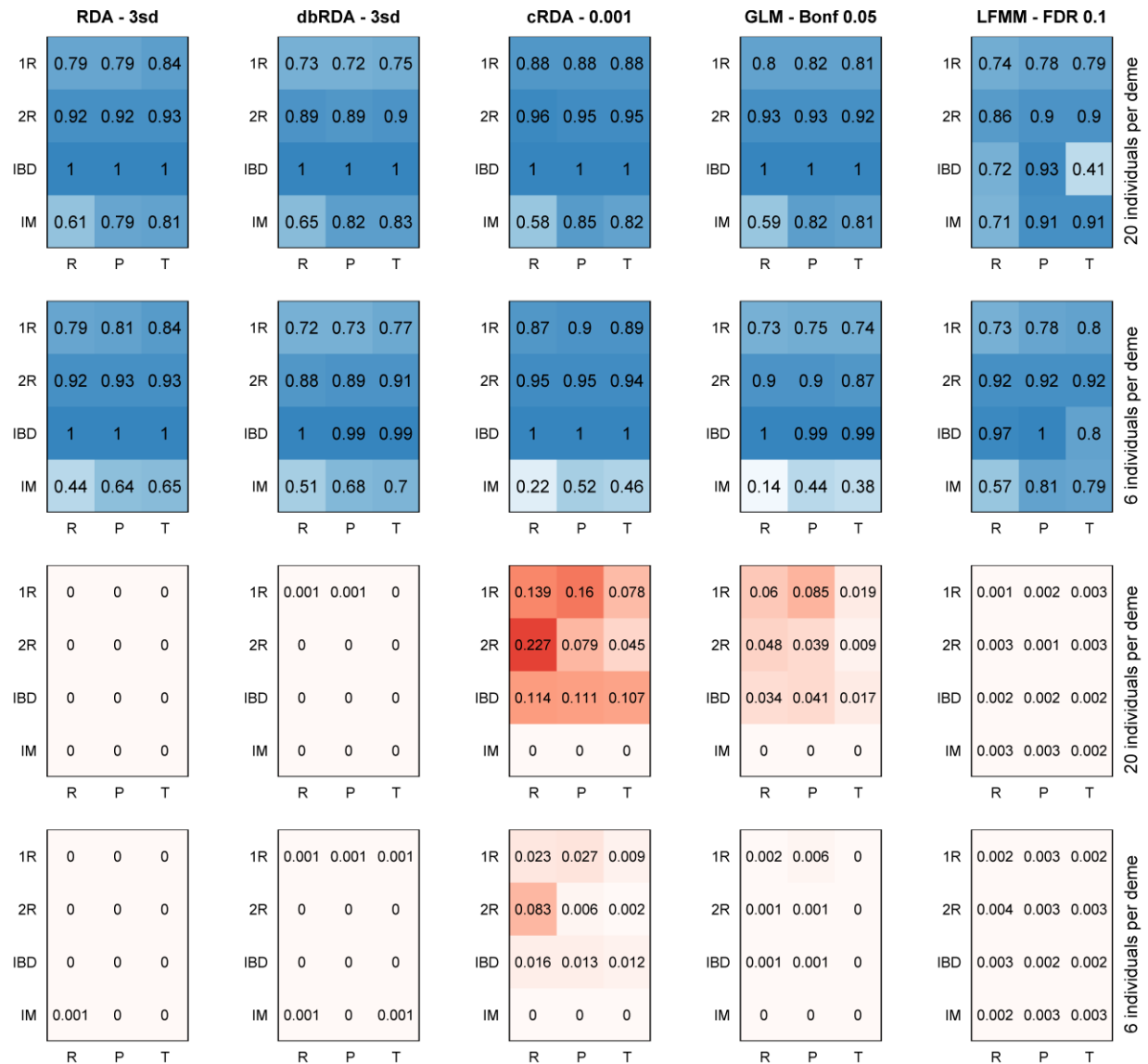
408

409  Cutoff results:

410  We compared cutoff results for the methods with the highest overall power: RDA, dbRDA,

411  cRDA, GLM, and LFMM. The best performing cutoffs were: RDA/dbRDA, +/- 3 SD; cRDA,

412  alpha = 0.001; GLM, Bonferroni = 0.05, and LFMM, FDR = 0.1. We did not choose the FDR

413  cutoff for GLMs since GIFs indicated that the test $p$-values were not appropriately calibrated

414  (i.e., GIFs > 1, Table S1). For some scenarios, LFMM GIFs were less than one (indicating a

415  conservative correction for population structure, Table S1). We reran LFMM models with the

416  best estimate of K minus one (i.e., K-1) to determine if a less conservative correction would

417  influence LFMM results. Because there was no consistent improvement in power or TPR/FPRs

418  using K-1 (Tables S2-S3), all subsequent results refer to LFMM runs using the best estimate of

419  K.

420      Full cutoff results for each method are presented in the Supplementary Information (Fig.

421  S3-S6). Cutoff FPRs were highest for cRDA and GLM (Fig. 6). By contrast, RDA and dbRDA

422  had mostly zero FPRs, with slightly higher FPRs for LFMM. Within these three low-FPR

423  methods, RDA maintained the highest TPRs, except in the IM demography, where LFMM

424  maintained higher power. LFMM was more sensitive to sampling design than the other methods,

425  with more variation in TPRs across designs.

19

**Figure 5.** Comparison of power (from empirical *p*-values) from RDA tested on weak selection only (x-axis) and RDA tested on all loci under selection (first row), as well as dbRDA and cRDA tested on weak selection only (second and third rows) for two sample sizes (columns). Points reflect demographies: 1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model. Some variation within demographies comes from sampling design.

20

**Figure 6.** Average true positive (top two rows, in blue) and false positive (bottom two rows, in red) rates from five methods (columns) using the best cutoff for each method. Each method shows results for different sampling strategies (R = random, P = pairs, T = transects), demographies (1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model), and sample sizes (rows).
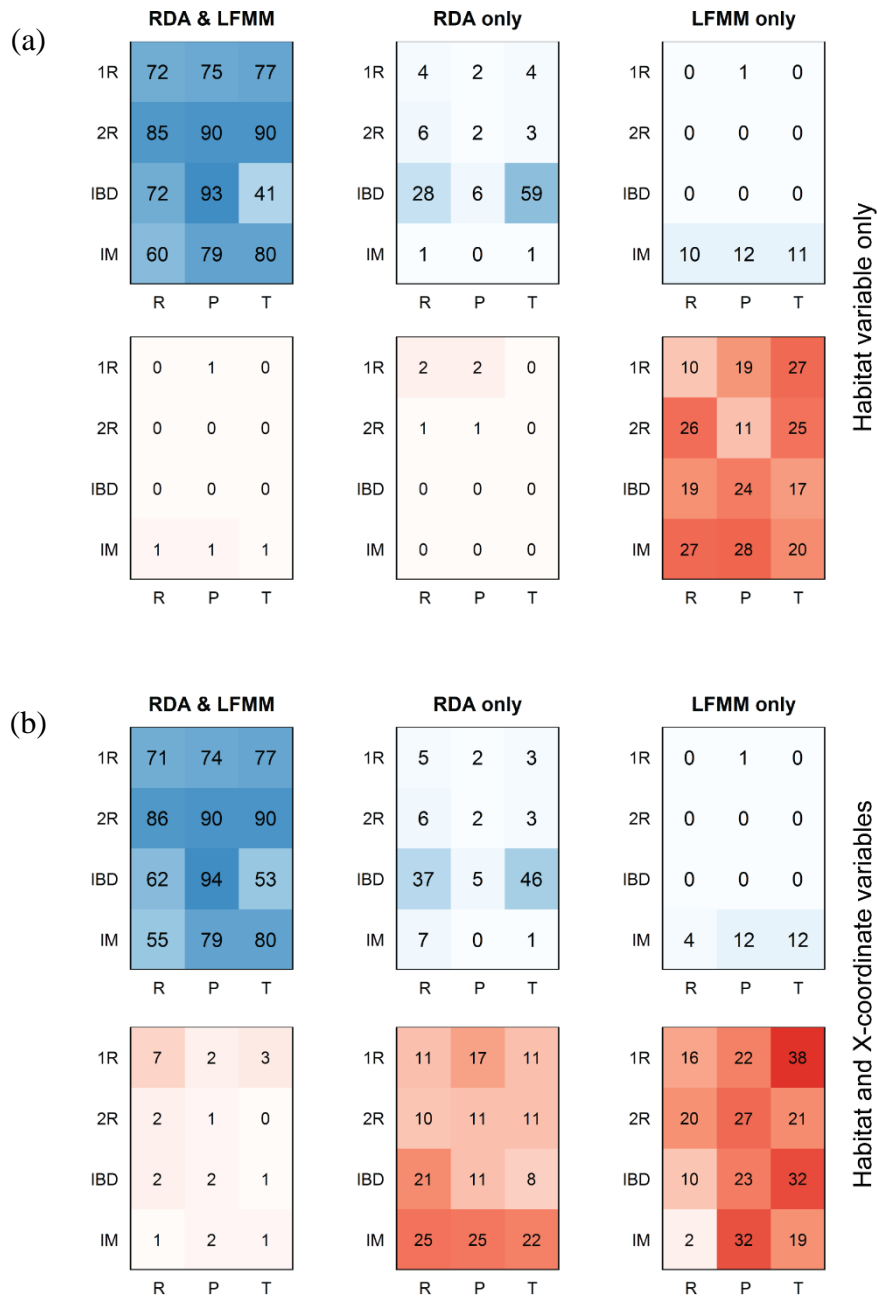
440     Combining detections:

441     We compared the univariate LFMM and multivariate RDA cutoff results for overlap and

442     differences in their detections using both the habitat predictor only, and the habitat and

443     (uninformative) x-coordinate predictor (Figs. 7 and S7). When the driving environmental

444     predictor is known, RDA detections alone are the best choice, since FPRs are very low and RDA

445     detects a large number of selected loci that are not identified by LFMM (except in the IM

446     demography, Fig. 7a). However, when a noninformative environmental predictor is included,

447     combining test results yields greater overall benefits, since the tests show substantial

448     commonality in TP detections, but show very low commonality in FP detections (Fig. 7b). By

449     retaining only overlapping loci, FPRs are substantially reduced at some loss of power due to

450     discarded RDA (and LFMM in the IM demography) detections.

451

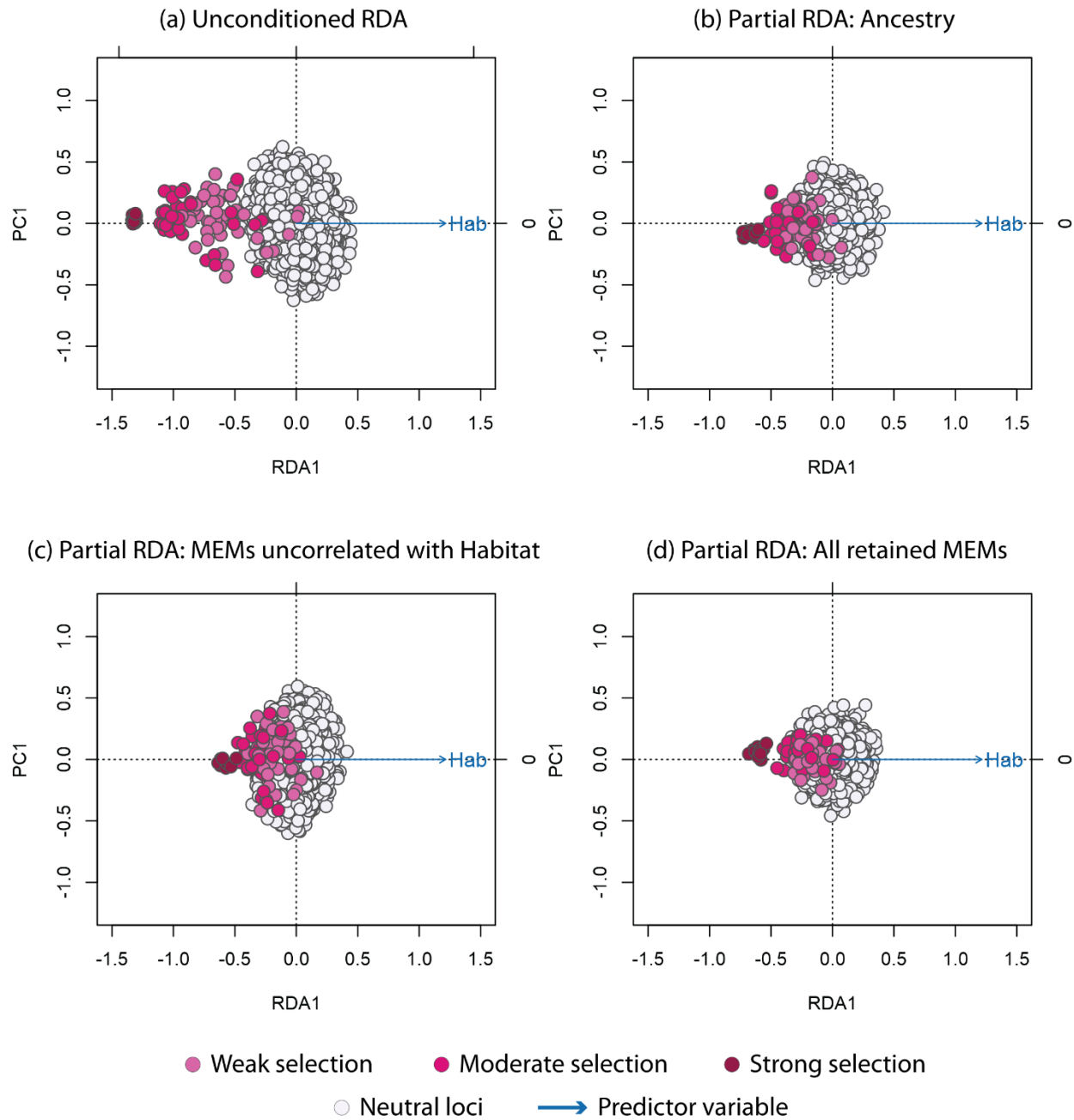452     Correction for population structure in RDA:

453     No MEM-based corrections for RDA were applied to IM scenarios, due to low spatial structure

454     (i.e., no PCoA axes were retained based on the broken-stick criterion). The more liberal approach

455     to correction using MEMs (removing retained MEMs significantly correlated with environment),

456     resulted in removal of MEMs with correlation coefficients ranging from 0.07 to 0.72. Ancestry-

457     based corrections were only applied to IM scenarios with 20 individuals since 6 individual

458     samples had K=1. All approaches that correct for population structure in RDA resulted in

459     substantial loss of power across all scenarios, both in terms of empirical $p$-values and cutoff

460     TPRs (Table 1 and Table S4). False positive rates (which were already very low for RDA)

461     increased slightly when correcting for population structure. There were only two scenarios where

462     FPRs improved (one and two fewer FP detections); however, these scenarios saw a reduction in

463     TPR of 81% and 92%, respectively (Table S4).

**Figure 7.** Average counts of true positive (top rows of a and b, in blue) and false positive (bottom rows of a and b, in red) detections for two methods, RDA and LFMM, using their best cutoffs and a sample size of 20 individuals per deme. The first column shows the average number of loci detected by both methods. The second and third columns show the average number of detections that are unique to RDA and LFMM, respectively. (a) Results for GEAs using Habitat as the only predictor. (b) Results for GEAs using Habitat and the (uninformative) X-coordinate predictor. Results are presented for different sampling strategies (R = random, P = pairs, T = transects), demographies (1R and 2R = refugial expansion, IBD = equilibrium isolation by distance, IM = equilibrium island model), and sample sizes (rows).

23

474     **Table 1.** Average change in power (from empirical *p*-values) and true and false positive rates

475     (from cutoffs) for RDA using three different approaches for partialling out population structure.

476     All approaches led to an overall loss of power and an increase in false positive rates. There are

477     no MEM corrections for the IM demography, which has no significant spatial structure. Ancestry

478     corrections apply only to 20 individual runs, where K ≠ 1.

479

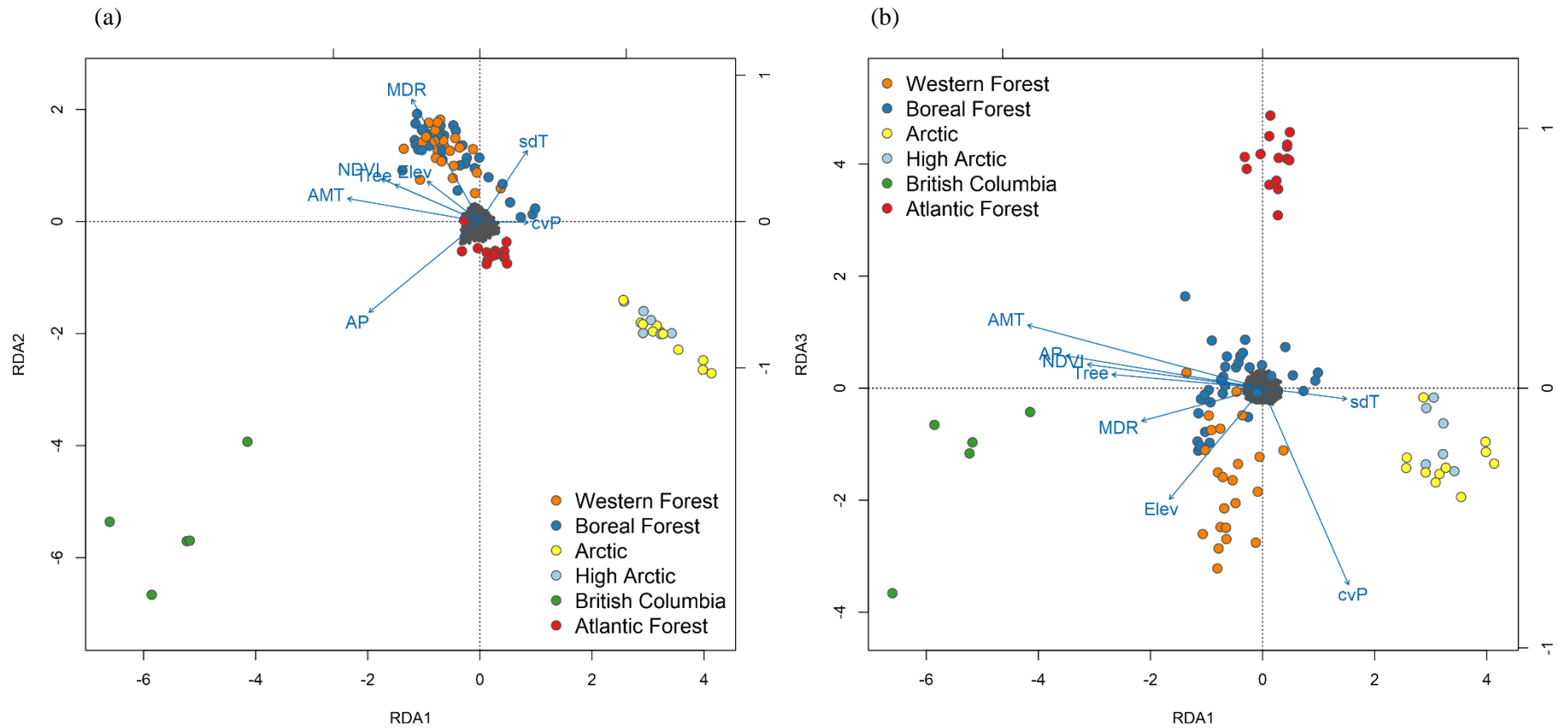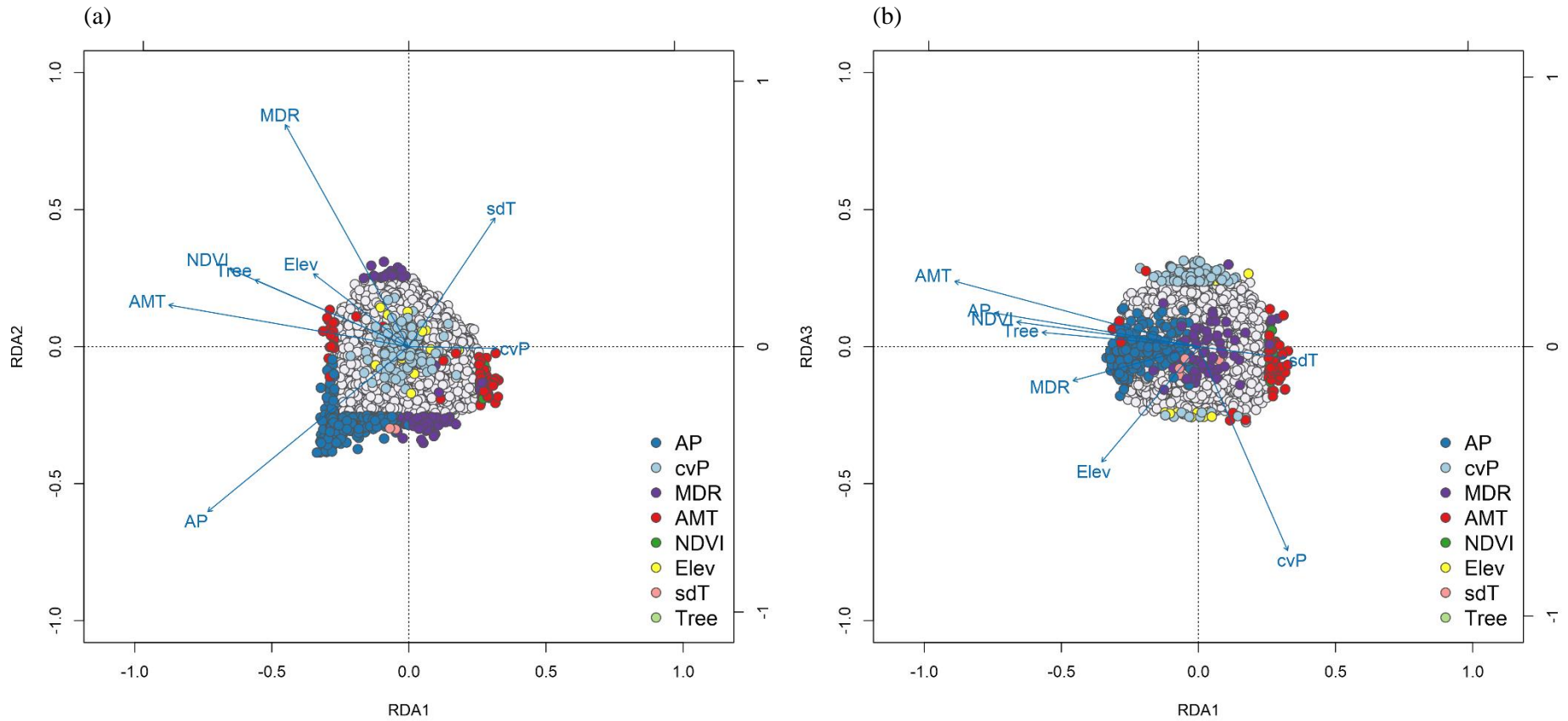| Indiv./ deme | Demo-graphy | Ancestry | MEMs uncorr. Habitat | All retained MEMs |
|---|---|---|---|---|
| | | **Change in power (empirical *p*-values)** | | |
| 6 | 1R | -0.53 | -0.59 | -0.72 |
| | 2R | -0.81 | -0.53 | -0.84 |
| | IBD | -0.94 | -0.75 | -0.96 |
| | IM | - | - | - |
| 20 | 1R | -0.26 | -0.14 | -0.58 |
| | 2R | -0.64 | -0.12 | -0.70 |
| | IBD | -0.93 | -0.69 | -0.93 |
| | IM | -0.70 | - | - |
| Average | | -0.69 | -0.47 | -0.79 |
| | | **Change in TPR (cutoffs)** | | |
| 6 | 1R | -0.39 | -0.43 | -0.69 |
| | 2R | -0.70 | -0.40 | -0.76 |
| | IBD | -0.93 | -0.69 | -0.94 |
| | IM | - | - | - |
| 20 | 1R | -0.16 | -0.16 | -0.47 |
| | 2R | -0.47 | -0.17 | -0.51 |
| | IBD | -0.92 | -0.60 | -0.90 |
| | IM | -0.71 | - | - |
| Average | | -0.61 | -0.41 | -0.71 |
| | | **Change in FPR (cutoffs)** | | |
| 6 | 1R | 0.0011 | 0.0013 | 0.0020 |
| | 2R | 0.0021 | 0.0011 | 0.0021 |
| | IBD | 0.0025 | 0.0017 | 0.0023 |
| | IM | - | - | - |
| 20 | 1R | 0.0005 | 0.0003 | 0.0014 |
| | 2R | 0.0014 | 0.0003 | 0.0015 |
| | IBD | 0.0021 | 0.0010 | 0.0021 |
| | IM | 0.0023 | - | - |
| Average | | 0.0017 | 0.0009 | 0.0019 |

**Figure 8.** Redundancy analysis biplots for simulation 1R, paired sampling, environmental surface 453, and 6 individuals per deme. Distribution of loci using: (a) unconditioned RDA (no correction for population structure); (b) partial RDA using ancestry values; (c) partial RDA using retained MEMs that are not significantly correlated with Habitat; (d) partial RDA using all retained MEMs.

486    Empirical data set:

487    There were four significant RDA axes in the ordination of the wolf data set (Fig. 9), which

488    returned 556 unique candidate loci that loaded +/- 3 SD from the mean loading on each axis: 171

489    SNPs detected on RDA axis 1, 222 on RDA axis 2, and 163 on RDA axis 3 (Fig. 10). Detections

490    on axis 4 were all redundant with loci already identified on axes 1-3. The majority of detected

491    SNPs were most strongly correlated with precipitation covariates: 231 SNPs correlated with

492    annual precipitation (AP) and 144 SNPs correlated with precipitation seasonality (cvP). The

493    number of SNPs correlated with the remaining predictors were: 72 with mean diurnal

494    temperature range (MDR); 79 with annual mean temperature (AMT); 13 with NDVI; 12 with

495    elevation; 4 with temperature seasonality (sdT); and 1 with percent tree cover (Tree).

**Figure 9.** Triplots of wolf data for (a) RDA axes 1 and 2, and (b) axes 1 and 3. The dark gray cloud of points at the center of each plot represent the SNPs, colored points represent individual wolves with coding by ecotype. Blue vectors represent environmental predictors (see text for abbreviations). Triplot scaling is symmetrical (both SNP and individual scores are scaled symmetrically by the square root of the eigenvalues).

**Figure 10.** Magnification of wolf data triplots from Figure 9 to highlight SNP loadings on (a) RDA axes 1 and 2, and (b) axes 1 and 3. Candidate SNPs are shown as colored points with coding by most highly correlated environmental predictor. SNPs not identified as candidates (neutral SNPs) are shown in light gray. Blue vectors represent environmental predictors (see text for abbreviations).

504   **Discussion**

505   Multivariate genotype-environment association (GEA) methods have been noted for their ability

506   to detect multilocus selection (Rellstab *et al.*, 2015; Hoban *et al.*, 2016), although there has been

507   no controlled assessment of the effectiveness of these methods in detecting multilocus selection

508   to date. Since these approaches are increasingly being used in empirical analyses (e.g. Bourret et

509   al. 2014; Brieuc et al. 2015; Pavey et al. 2015; Hecht et al. 2015; Laporte et al. 2016; Brauer et

510   al. 2016), it is important that these claims are evaluated to ensure that the most effective GEA

511   methods are being used, and that their results are being appropriately interpreted.

512        Here we compare a suite of methods for detecting selection in a simulation framework to

513   assess their ability to correctly detect multilocus selection under different demographic and

514   sampling scenarios. We found that constrained ordinations had the best overall performance

515   across the demographies, sampling designs, sample sizes, and selection levels tested here. The

516   univariate LFMM method also performed well, though power was scenario-dependent and was

517   reduced for loci under weak selection (in agreement with findings by de Villemereuil *et al.*,

518   2014). Random Forest, by contrast, had lower detection rates overall. In the following sections

519   we discuss the performance of these methods and provide suggestions for their use on empirical

520   data sets.

521

522   Random Forest:

523   Random Forest performed relatively poorly as a GEA. This poor performance is caused by the

524   sparsity of the genotype matrix (i.e., most SNPs are not under selection), which results in

525   detection that is dominated by strongly selected loci (i.e., loci with strong marginal effects). This

526   issue has been documented in other simulation and empirical studies (Goldstein *et al.*, 2010;

527   Winham *et al.*, 2012; Wright *et al.*, 2016) and indicates that RF is not suited to identifying weak

528   multilocus selection or interaction effects in these large data sets. Empirical studies that have

529   used RF as a GEA have likely identified a subset of loci under strong selection, but are unlikely

530   to have identified loci underlying more complex genetic architectures. Note that the amount of

531   environmental variance explained by the RF model can be high (i.e., overall percent variance

532   explained by the detected SNPs, which ranged from 79-91% for these simulations, Table S5),

533   while still failing to identify most of the loci under selection. Removing strong associations from

29

534 the genotypic matrix can potentially help with the detection of weaker effects (Goldstein *et al.*,

535 2010), but this approach has not been tested on large matrices. Combined with the computational

536 burden of this method (taking ~10 days on a single core for the larger data sets), as well as the

537 availability of fast and accurate alternatives such as RDA (which takes ~3 minutes on the same

538 data), it is clear that RF is not a viable option for GEA analysis of genomic data.

539  Random Forest does hold promise for the detection of interaction effects in much smaller

540 data sets (e.g., tens of loci, Holliday *et al*. 2012). However, this is an area of active research, and

541 the capacity of RF models in their current form to both capture and identify SNP interactions has

542 been disputed (Winham *et al.*, 2012; Wright *et al.*, 2016). New modifications of RF models are

543 being developed to more effectively identify interaction effects (e.g. Li et al. 2016), but these

544 models are computationally demanding and are not designed for large data sets. Overall,

545 extensions of RF show potential for identifying more complex genetic architectures on small sets

546 of loci, but caution is warranted in using them on empirical data prior to rigorous testing on

547 realistic simulation scenarios.

548

549 <u>Constrained ordinations</u>:

550 The three constrained ordination methods all performed well. RDA in particular had the highest

551 overall power across all methods tested here (Figs. 1-3). Ordinations were relatively insensitive

552 to sample size (6 vs 20 individuals sampled per deme), with the exception of the IM

553 demography, where larger sample sizes consistently improved TPRs, as previously noted by De

554 Mita *et al*. (2013) and Lotterhos & Whitlock (2015) for univariate GEAs. Power was lowest in

555 the IM demography, which is typified by a lack of spatial autocorrelation in allele frequencies

556 and a reduced signal of local adaptation (Table S6), making detection more difficult. This

557 corresponds with univariate GEA results from Lotterhos & Whitlock (2015), who found very

558 low detection rates for loci under weak selection in the IM demography. Power was highest for

559 IBD, followed by the 2R and 1R demographies. Data from natural systems likely lie somewhere

560 among these demographic extremes, and successful differentiation in the presence of IBD and

561 non-equilibrium conditions indicate that ordinations should work well across a range of natural

562 systems.

563       All three methods were relatively insensitive to sampling design, with transects

564    performing slightly better in 1R and random sampling performing worst in IM (Figs. 4, 6, and

565    S2). Otherwise results were consistent across designs, in contrast to the univariate GEAs tested

566    by Lotterhos and Whitlock (2015), most of which had higher power with the paired sampling

567    strategy. Ordinations are likely less sensitive to sampling design since they take advantage of

568    covarying signals of selection across loci, making them more robust to sampling that does not

569    maximize environmental differentiation (e.g., random or transect designs). All methods

570    performed similarly in terms of detection rates across selection strengths (Figs. 4 and S2). As

571    expected, weak selection was more difficult to detect than moderate or strong selection, except

572    for IBD, where detection levels were high regardless of selection.

573       High TPRs were maintained when using cutoffs for all three ordination methods (Fig. 6).

574    False positives were universally low for RDA and dbRDA. By contrast, cRDA showed high

575    FPRs for all demographies except IM, tempering its slightly higher TPRs. These higher FPRs are

576    a consequence of using component axes as predictors. Across all scenarios and sample sizes,

577    cRDA detected component 1, 2, or both as significantly associated with the constrained RDA

578    axes (Table S7). Most selected loci load on these components (keeping TPRs high), but neutral

579    markers also load on these axes, especially in cases where there are strong trends in neutral loci

580    (i.e., maximum trends in neutral markers reflect FPRs for cRDA, Table S6, Fig. 6). Given these

581    results, we hypothesized that it might be challenging for cRDA to detect weak selection in the

582    absence of a covarying signal from loci with stronger selection coefficients. If the selection

583    signature is weak, it may load on a lower-level component axis (i.e., an axis that explains less of

584    the genetic variance), or it may load on higher-level axes, but fail to be significantly associated

585    with the constrained axes. Note that although cRDA contains a step to reduce the number of

586    components, parallel analysis resulted in retention of all axes in every simulation tested here

587    (Table S7). This meant that cRDA could search for the signal of selection across all possible

588    components.

589       When tested on simulations with loci under weak selection only, RDA, which uses the

590    genotype matrix directly, maintained similar power as in the full data set (except in the IM

591    scenario, where power was higher when all selected loci were included), indicating that selection

592    signals can be detected with this method in the absence of loci under strong selection (Fig. 5, top

593     row). By contrast, cRDA detection was more variable, ranging from comparable detection rates

594     with the full data set, to no/poor detections under certain demographies and sample sizes. In

595     these latter cases, poor performance is reflected in the component axes detected as significant

596     (Table S7); instead of identifying the signal in the first few axes, a variable set of lower-variance

597     axes are detected (or none are detected at all). This indicates that the method is not able to "find"

598     the selected signal in the component axes in cases where that signal is not driven by strong

599     selection. This result, in addition to higher FPRs for cRDA, builds a case for using the genotype

600     matrix directly with a constrained ordination such as RDA or dbRDA, as opposed to a

601     preliminary step of data conversion with PCA.

602

603     <u>Should results from different tests be combined?</u>

604     A common approach in local adaptation studies is to run multiple tests (GEA only, or a

605     combination of GEA and differentiation methods) and look for overlapping detections across

606     methods. This ad hoc approach is thought to increase confidence in TPRs, while minimizing

607     FPRs. The problem with this approach is that it can bias detection toward strong selective sweeps

608     to the exclusion of other adaptive mechanisms which may be equally important in shaping

609     phenotypic variation (Le Corre & Kremer, 2012; François *et al.*, 2016). If the goal is to detect

610     other forms of selection such as recent selection or selection on standing genetic variation, this

611     approach will not be effective since most methods are unlikely to detect these weak signals.

612     Additionally, this approach limits detections to those of the least powerful method used, forcing

613     overall detection rates to be a function of the weakest method implemented.

614         The complexities of this issue are illustrated by comparing results across two sets of RDA

615     and LFMM results: one where the driving environmental variable is known (Fig. 7a), and

616     another where the environmental predictors represent hypotheses about the most important

617     factors driving selection (Fig. 7b). In both cases, agreement on TPs is high, and RDA has a large

618     number of true positive detections that are unique to that method, while unique detections by

619     LFMM are largely limited to the IM demography. The differences in the cases lies in FP

620     detections: when selection is well understood, and uninformative predictors are not used,

621     retaining RDA detections only is the approach that will maximize TPRs (and detection of weak

622     loci under selection) while maintaining minimal to zero FPRs (Fig. 7a). Where GEA analyses are

32

623   more exploratory (i.e., when selective gradients are unknown), combining detections can help

624   reduce FPRs (Fig. 7b). If some FP detections are acceptable, keeping only RDA detections will

625   improve TPRs at the cost of slightly increased FPRs. A third approach, keeping all detections

626   across both methods, would yield little improvement in TPRs in both cases, since LFMM has

627   high levels of unique FPs, and minimal unique TP detections.

628      The decision of whether and how to combine results from different tests will be specific

629   to the study questions, the tolerance for false negative and false positive detections, and the

630   capacity for follow-up analyses on detected markers. For example, if the goal is to detect loci

631   with strong effects while keeping false positive rates as low as possible, or GEA is being used as

632   an exploratory analysis, running multiple GEA methods and considering only overlapping

633   detections could be a suitable strategy. However, if the goal is to detect selection on standing

634   genetic variation or a recent selection event, and the most important selective agents (or close

635   correlates of them) are known, combining detections from multiple tests would likely be too

636   conservative. In this case, the best approach would be to use a single GEA method, such as

637   RDA, that can effectively detect covarying signals arising from multilocus selection, while being

638   robust to selection strength, sampling design, and sample size.

639

640   Correction for population structure:

641   All three methods used to correct for populations structure in RDA resulted in substantial loss of

642   power and, in most cases, increased FPRs (Table 1 and S4). The effect of correcting for

643   population structure can be seen in ordination biplots from an example simulation scenario (Fig.

644   8). In this 1R demographic scenario, the selection surface ("Hab") and the refugial expansion

645   gradient coincide, so any correction for population structure will also remove the signal of

646   selection from the selected loci. The correction is most conservative when using all significant

647   MEM predictors to account for spatial structure (Fig. 8d), and is less conservative when using

648   only MEMs not significantly correlated with environment (Fig. 8c), or ancestry coefficients (Fig.

649   8b). In all cases, however, the loss of the selection signal is significant (Table 1), and is visible in

650   the increasing overlap of selected loci with neutral loci.

651      While the simulations used here have overall low global Fst (average Fst = 0.05),

652   population structure is significant enough in many scenarios to result in elevated FPRs for GLMs

653    (univariate linear models which do not correct for population structure, Fig. 6). Despite this,

654    RDA and dbRDA (the multivariate analogue of GLMs) do not show elevated FPRs, even when

655    selection covaries with a range expansion front, as in the 1R and 2R demographies. This is likely

656    because only loci with extreme loadings are identified as potentially under selection, leaving

657    most neutral loci, which share a similar, but weaker, spatial signature, loading less than +/- 3 SD

658    from the mean. The generality of these results needs to be tested in a comprehensive manner

659    using an expanded simulation parameter space that includes stronger population structure and

660    metapopulation dynamics; this work is currently in progress. In the meantime, we recommend

661    that RDA be used conservatively in empirical systems with higher population structure than is

662    tested here, for example, by finding overlap between detections identified by RDA and LFMM

663    (or another GEA that accounts for population structure).

664

665    Empirical example:

666    Triplots of three of the four significant RDA axes for the wolf data show SNPs (dark gray

667    points), individuals (colored circles), and environmental variables (blue arrows, Fig. 9). The

668    relative arrangement of these items in the ordination space reflects their relationship with the

669    ordination axes, which are linear combinations of the predictor variables. For example,

670    individuals from wet and temperate British Columbia are positively related to high annual

671    precipitation (AP) and low temperature seasonality (sdT, Fig. 9a). By contrast, Artic and High

672    Arctic individuals are characterized by small mean diurnal temperature range (MDR), low

673    annual mean temperature (AMT), lower levels of tree cover (Tree) and NDVI (a measure of

674    vegetation greenness), and are found at lower elevation (Fig. 9a). Atlantic Forest and Western

675    Forest individuals load more strongly on RDA axis 3, showing weak and strong precipitation

676    seasonality (cvP) respectively (Fig. 9b), consistent with continental-scale climate in these

677    regions.

678            If we zoom into the SNPs, we can visualize how candidate SNPs load on the RDA axes

679    (Fig. 10). For example, SNPs most strongly correlated with AP have strong loadings in the lower

680    left quadrant between RDA axes 1 and 2 along the AP vector, accounting for the majority of

681    these 231 AP-correlated detections (Fig. 10a). Most candidates highly correlated with AMT and

682    MDR load strongly on axes 1 and 2, respectively. Note how candidate SNPs correlated with

683    precipitation seasonality (cvP) and elevation are located in the center of the plot, and will not be

684    detected as outliers on axes 1 or 2 (Fig. 10a). However, these loci are detected as outliers on axis

685    3 (Fig. 10b). Overall, candidate SNPs on axis 1 represent multilocus haplotypes associated with

686    annual precipitation and mean diurnal range; SNPs on axis 2 represent haplotypes associated

687    with annual precipitation and annual mean temperature; and SNPs on axis 3 represent haplotypes

688    associated with precipitation seasonality.

689          Of the 1661 candidate SNPs identified by Schweizer *et al.,* (2016) using Bayenv (Bayes

690    Factor > 3), only 52 were found in common with the 556 candidates from RDA. Of these 52

691    common detections, only nine were identified based on the same environmental predictor. If we

692    include Bayenv detections using highly correlated predictors (removed for RDA) we find nine

693    more candidates identified in common. Additionally, only 18% of the Bayenv identifications

694    were most strongly related to precipitation variables, which are known drivers of morphology

695    and population structure in gray wolves (Geffen *et al.*, 2004; O'Keefe *et al.*, 2013; Schweizer *et*

696    *al.*, 2016). By contrast, 67% of RDA detections were most strongly associated with precipitation

697    variables, providing new candidate regions for understanding local adaptation of gray wolves

698    across their North American range.

699

700    <u>Conclusions and recommendations</u>:

701    We found that constrained ordinations, especially RDA, show a superior combination of low

702    FPRs and high TPRs across weak, moderate, and strong multilocus selection. These results were

703    robust across the levels of population structure, demographic histories, sampling designs, and

704    sample sizes tested here. Additionally, RDA outperformed an alternative ordination-based

705    approach, cRDA, especially (and importantly) when the multilocus selection signature was

706    completely derived from loci under weak selection. It is important to note that population

707    structure was relatively low in these simulations. Results may differ for systems with strong

708    population structure or metapopulation dynamics, where it can be important to correct for

709    structure or combine detections with another GEA that accounts for structure. Continued testing

710    of these promising methods is needed in simulation frameworks that include more population

711    structure, multiple selection surfaces, and genetic architectures that are more complex than the

712    multilocus selection response modeled here. However, this study indicates that constrained

713    ordinations are an effective means of detecting adaptive processes that result in weak, multilocus

714    molecular signatures, providing a powerful tool for investigating the genetic basis of local

715    adaptation and informing management actions to conserve the evolutionary potential of species

716    of agricultural, forestry, fisheries, and conservation concern.

**Acknowledgements**

## References

724

725 Angers B, Magnan P, Plante M, Bernatchez L (1999) Canonical correspondence analysis for estimating
726     spatial and environmental effects on microsatellite gene diversity in brook charr (*Salvelinus*
727     *fontinalis*). *Molecular Ecology*, **8**, 1043–1053.

728 Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate - a practical and powerful approach.
729     *Journal of the Royal Statistical Society Series B-Methodological*, **57**, 289–300.

730 Bivand R, Hauke J, Kossowski T (2013) Computing the Jacobian in Gaussian spatial autoregressive
731     models: an illustrated comparison of available methods. *Geographical Analysis*, **45**, 150–179.

732 Blanchet FG, Legendre P, Borcard D (2008) Forward selection of explanatory variables. *Ecology*, **89**,
733     2623–2632.

734 Bourret V, Dionne M, Bernatchez L (2014) Detecting genotypic changes associated with selective
735     mortality at sea in Atlantic salmon: polygenic multilocus analysis surpasses genome scan.
736     *Molecular Ecology*, **23**, 4444–4457.

737 Brauer CJ, Hammer MP, Beheregaray LB (2016) Riverscape genomics of a threatened fish across a
738     hydroclimatically heterogeneous river basin. *Molecular Ecology*, **25**, 5093–5113.

739 Bray JR, Curtis JT (1957) An ordination of the upland forest communities of southern Wisconsin.
740     *Ecological Monographs*, **27**, 325–349.

741 Breiman L (2001) Random Forests. *Machine Learning*, **45**, 5–32.

742 Brieuc MSO, Ono K, Drinan DP, Naish KA (2015) Integration of Random Forest with population-based
743     outlier analyses provides insight on the genomic basis and evolution of run timing in Chinook
744     salmon (*Oncorhynchus tshawytscha*). *Molecular Ecology*, **24**, 2729–2746.

745 Cavalli-Sforza LL (1966) Population structure and human evolution. *Proceedings of the Royal Society B:*
746     *Biological Sciences*, **164**, 362–379.

747 Coop G, Witonsky D, Rienzo AD, Pritchard JK (2010) Using environmental correlations to identify loci
748     underlying local adaptation. *Genetics*, **185**, 1411–1423.

749 De Mita S, Thuillet A-C, Gay L, Ahmadi N, Manel S, Ronfort J, Vigouroux Y (2013) Detecting selection
750     along environmental gradients: analysis of eight methods and their effectiveness for outbreeding
751     and selfing populations. *Molecular Ecology*, **22**, 1383–1399.

752 De'ath G, Fabricius KE (2000) Classification and regression trees: a powerful yet simple technique for
753     ecological data analysis. *Ecology*, **81**, 3178–3192.

754 Dray S, Legendre P, Peres-Neto PR (2006) Spatial modelling: a comprehensive framework for principal
755     coordinate analysis of neighbour matrices (PCNM). *Ecological Modelling*, **196**, 483–493.

756 Dray S, Pélissier R, Couteron P et al. (2012) Community ecology in the age of multivariate multiscale
757     spatial analysis. *Ecological Monographs*, **82**, 257–275.

758 Dray S, Blanchet G, Borcard D et al. (2016) *adespatial: Multivariate multiscale spatial analysis*. R
759     package version 0.0-7.

Duforet-Frebourg N, Bazin E, Blum MGB (2014) Genome scans for detecting footprints of local adaptation using a Bayesian factor model. *Molecular Biology and Evolution*, **31**, 2483–2495.

Forester BR, Jones MR, Joost S, Landguth EL, Lasky JR (2016) Detecting spatial genetic signatures of local adaptation in heterogeneous landscapes. *Molecular Ecology*, **25**, 104–120.

François O, Martins H, Caye K, Schoville SD (2016) Controlling false discoveries in genome scans for selection. *Molecular Ecology*, **25**, 454–469.

Frichot E, François O (2015) LEA: An R package for landscape and ecological association studies. *Methods in Ecology and Evolution*, **6**, 925–929.

Frichot E, Schoville SD, Bouchard G, François O (2013) Testing for associations between loci and environmental gradients using latent factor mixed models. *Molecular Biology and Evolution*, **30**, 1687–1699.

Frichot E, Mathieu F, Trouillon T, Bouchard G, François O (2014) Fast and efficient estimation of individual ancestry coefficients. *Genetics*, **196**, 973–983.

Geffen E, Anderson MJ, Wayne RK (2004) Climate and habitat barriers to dispersal in the highly mobile grey wolf. *Molecular Ecology*, **13**, 2481–2490.

Goldstein BA, Hubbard AE, Cutler A, Barcellos LF (2010) An application of Random Forests to a genome-wide association dataset: methodological considerations & new findings. *BMC Genetics*, **11**, 1–13.

Grivet D, Sork VL, Westfall RD, Davis FW (2008) Conserving the evolutionary potential of California valley oak (*Quercus lobata* Née): a multivariate genetic approach to conservation planning. *Molecular Ecology*, **17**, 139–156.

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

Hancock AM, Brachi B, Faure N et al. (2011) Adaptation to climate across the *Arabidopsis thaliana* genome. *Science*, **334**, 83–86.

Harrisson KA, Pavlova A, Telonis-Scott M, Sunnucks P (2014) Using genomics to characterize evolutionary potential for conservation of wild populations. *Evolutionary Applications*, **7**, 1008–1025.

Hecht BC, Matala AP, Hess JE, Narum SR (2015) Environmental adaptation in Chinook salmon (*Oncorhynchus tshawytscha*) throughout their North American range. *Molecular Ecology*, **24**, 5573–5595.

Hoban S, Kelley JL, Lotterhos KE et al. (2016) Finding the genomic basis of local adaptation: pitfalls, practical solutions, and future directions. *The American Naturalist*, **188**, 379–397.

Holliday JA, Wang T, Aitken S (2012) Predicting adaptive phenotypes from multilocus genotypes in Sitka Spruce (*Picea sitchensis*) using Random Forest. *G3: Genes|Genomes|Genetics*, **2**, 1085–1093.

Horn JL (1965) A rationale and test for the number of factors in factor analysis. *Psychometrika*, **30**, 179–185.

798    Huang F (2015) *hornpa: Horn's (1965) Test to Determine the Number of Components/Factors.* R
799          package version 1.0.

800    Jombart T, Pontier D, Dufour A-B (2009) Genetic markers in the playground of multivariate analysis.
801          *Heredity*, **102**, 330–341.

802    Joost S, Bonin A, Bruford MW, Després L, Conord C, Erhardt G, Taberlet P (2007) A spatial analysis
803          method (SAM) to detect candidate loci for selection: towards a landscape genomics approach to
804          adaptation. *Molecular Ecology*, **16**, 3955–3969.

805    Laporte M, Pavey SA, Rougeux C et al. (2016) RAD sequencing reveals within-generation polygenic
806          selection in response to anthropogenic organic and metal contamination in North Atlantic Eels.
807          *Molecular Ecology*, **25**, 219–237.

808    Lasky JR, Des Marais DL, McKay JK, Richards JH, Juenger TE, Keitt TH (2012) Characterizing
809          genomic variation of *Arabidopsis thaliana*: the roles of geography and climate. *Molecular
810          Ecology*, **21**, 5512–5529.

811    Lasky JR, Marais D, L D et al. (2014) Natural variation in abiotic stress responsive gene expression and
812          local adaptation to climate in *Arabidopsis thaliana*. *Molecular Biology and Evolution*, **31**, 2283–
813          2296.

814    Lasky JR, Upadhyaya HD, Ramu P et al. (2015) Genome-environment associations in sorghum landraces
815          predict adaptive traits. *Science Advances*, **1**, e1400218.

816    Le Corre V, Kremer A (2012) The genetic differentiation at quantitative trait loci under local adaptation.
817          *Molecular Ecology*, **21**, 1548–1566.

818    Legendre P, Legendre L (2012) *Numerical Ecology*, 3rd edition edn. Elsevier, Amsterdam, The
819          Netherlands.

820    Li J, Malley JD, Andrew AS, Karagas MR, Moore JH (2016) Detecting gene-gene interactions using a
821          permutation-based random forest method. *BioData Mining*, **9**, 14.

822    Lotterhos KE, Whitlock MC (2014) Evaluation of demographic history and neutral parameterization on
823          the performance of FST outlier tests. *Molecular Ecology*, **23**, 2178–2192.

824    Lotterhos KE, Whitlock MC (2015) The relative power of genome scans to detect local adaptation
825          depends on sampling design and statistical method. *Molecular Ecology*, **24**, 1031–1046.

826    Mitton JB, Linhart YB, Hamrick JL, Beckman JS (1977) Observations on the genetic structure and
827          mating system of Ponderosa Pine in the Colorado front range. *Theoretical and Applied Genetics*,
828          **51**, 5–13.

829    Mulley JC, James JW, Barker JSF (1979) Allozyme genotype-environment relationships in natural
830          populations of *Drosophila buzzatii*. *Biochemical Genetics*, **17**, 105–126.

831    O'Keefe FR, Meachen J, Fet EV, Brannick A (2013) Ecological determinants of clinal morphological
832          variation in the cranium of the North American gray wolf. *Journal of Mammalogy*, **94**, 1223–
833          1236.

834    Oksanen J, Blanchet FG, Kindt R et al. (2013) *vegan: Community Ecology Package*.

835    Pavey SA, Gaudin J, Normandeau E, Dionne M, Castonguay M, Audet C, Bernatchez L (2015) RAD
836         sequencing highlights polygenic discrimination of habitat ecotypes in the panmictic American
837         Eel. *Current Biology*, **25**, 1666–1671.

838    R Development Core Team (2015) *R: a language and environment for statistical computing*. R
839         Foundation for Statistical Computing, Vienna, Austria.

840    Rellstab C, Gugerli F, Eckert AJ, Hancock AM, Holderegger R (2015) A practical guide to environmental
841         association analysis in landscape genomics. *Molecular Ecology*, **24**, 4348–4370.

842    Savolainen O, Lascoux M, Merilä J (2013) Ecological genomics of local adaptation. *Nature Reviews
843         Genetics*, **14**, 807–820.

844    Schweizer RM, vonHoldt BM, Harrigan R et al. (2016) Genetic subdivision and candidate genes under
845         selection in North American grey wolves. *Molecular Ecology*, **25**, 380–402.

846    Storey JD, Bass AJ, Dabney A, Robinson D (2015) *qvalue: Q-value estimation for false discovery rate
847         control*. R package version 2.2.2.

848    Stucki S, Orozco-terWengel P, Forester BR et al. (2016) High performance computation of landscape
849         genomic models including local indicators of spatial association. *Molecular Ecology Resources*.

850    Tiffin P, Ross-Ibarra J (2014) Advances and limits of using population genetics to understand local
851         adaptation. *Trends in Ecology & Evolution*, **29**, 673–680.

852    de Villemereuil P, Frichot É, Bazin É, François O, Gaggiotti OE (2014) Genome scan methods against
853         more complex models: when and how much should we trust them? *Molecular Ecology*, **23**, 2006–
854         2019.

855    Winham SJ, Colby CL, Freimuth RR, Wang X, Andrade M de, Huebner M, Biernacka JM (2012) SNP
856         interaction detection with Random Forests in high-dimensional genetic data. *BMC
857         Bioinformatics*, **13**, 164.

858    Wright MN, Ziegler A, König IR (2016) Do little interactions get lost in dark random forests? *BMC
859         Bioinformatics*, **17**, 145.

860    Yeaman S, Whitlock MC (2011) The genetic architecture of adaptation under migration–selection
861         balance. *Evolution*, **65**, 1897–1911.

862    Yoder JB, Stanton-Geddes J, Zhou P, Briskine R, Young ND, Tiffin P (2014) Genomic signature of
863         adaptation to climate in *Medicago truncatula*. *Genetics*, **196**, 1263–1275.

864 **Data accessibility**

865 Simulation data from Lotterhos & Whitlock (2015): Dryad: doi:10.5061/dryad.mh67v

866 Supporting simulation data (coordinate files) for Lotterhos & Whitlock (2015) data provided by

867 Wagner *et al*. (2017): Dryad: doi:10.5061/dryad.b12kk. Wolf data from Schweizer *et al.* (2016):

868 Dryad: doi.org/10.5061/dryad.c9b25.

869

870 **Author contributions**

871 BRF and DLU conceived the study. BRF performed the analyses and wrote the manuscript.

872 HHW contributed code. JRL, HHW, and DLU helped interpret the results and write the

873 manuscript.

874

875 **Supporting information**

876 Forester_Simulation_Rcode.zip: Contains R code for data preparation and all of the methods

877 tested.

878 Forester_Wolf_Rcode.zip: Contains R code for data preparation, analysis, interpretation, and

879 plotting of the wolf data set with RDA.

880 Forester_SI.pdf: Supplemental Figures S1-S7 and Tables S1-S7.