

1 **Optimising sampling design and sequencing strategy for the**
2 **genomic analysis of quantitative traits in natural populations**

3

4 Jefferson F. Paril¹, David J. Balding^{1,2,3}, Alexandre Fournier-Level^{1,2}

5 ¹School of Biosciences, ²Melbourne Integrative Genomics, and ³School of Mathematics and
6 Statistics, The University of Melbourne, Parkville 3010, Australia

7

8 **Author for correspondence:** Alexandre Fournier-Level (afournier@unimelb.edu.au)

9

10 **Keywords:** genome-wide association, individual sequencing, landscape genomics, pool
11 sequencing, genomic prediction, simulation

12

13 **Abstract**

14 Mapping the genes underlying ecologically-relevant traits in natural populations is fundamental
15 to develop a molecular understanding of species adaptation. Current sequencing technologies
16 enable the characterisation of a species' genetic diversity across the landscape or even over its
17 whole range. The relevant capture of the genetic diversity across the landscape is critical for a
18 successful genetic mapping of traits and there are no clear guidelines on how to achieve an
19 optimal sampling and which sequencing strategy to implement. Here we determine through
20 simulation, the sampling scheme that maximises the power to map the genetic basis of a
21 complex trait in an outbreeding species across an idealised landscape and draw genomic
22 predictions for the trait, comparing individual and pool sequencing strategies. Our results show
23 that QTL detection power and prediction accuracy are higher when more populations over the
24 landscape are sampled and this is more cost-effectively done with pool sequencing than with
25 individual sequencing. Additionally, we recommend sampling populations from areas of high
26 genetic diversity. As progress in sequencing enables the integration of trait-based functional
27 ecology into landscape genomics studies, these findings will guide study designs allowing direct
28 measures of genetic effects in natural populations across the environment.

29

30 **Introduction**

31 Understanding how the molecular variation within species supports the evolution of functional
32 traits is a central goal in ecology through the determination of so-called genotype-to-phenotype
33 map. Genomic information for a species can then be leveraged to understand and eventually
34 predict population fitness under a range of eco-evolutionary scenarios. Unfortunately, this
35 genotype-to-phenotype map is only available for a handful of traits, and primarily in model
36 organisms. With the improved accessibility of sequencing technologies, genome-wide
37 association studies (GWAS) and genomic prediction (GP) are becoming straightforward
38 approaches to understand and predict complex traits (Gondro et al., 2013). We thus coined the
39 term GPAS (Genomic Prediction and Association Studies) to denote genome-wide association
40 studies designed to both identify quantitative trait loci (QTL) and predict traits from genomic
41 data. These studies rely on cost-effective, high-throughput sequencing and share the same well-
42 established linear modelling framework. However, how to sample natural populations to train
43 accurate GPAS models that are representative of the genetic diversity of a species is far from
44 obvious. More insights are needed to develop an optimal strategy and move away from *ad hoc*
45 field sampling.

46

47 Research in ecological genomics deals with the challenge of characterising the genetic basis of
48 traits across multiple natural populations. This requires collecting sufficient genotype and
49 experimental phenotype data to represent the species' diversity, which in practice is often
50 performed with limited resources. This raises the problem of how to sample across a landscape
51 to capture representative genetic variation while constrained by the total sequencing depth

52 attainable for a given budget. Here, we address the question: how do we allocate a fixed
53 sequencing capacity so that the genetic information captured over the landscape leads to an
54 optimal GPAS performance?

55

56 It is becoming easier to genotype genome-wide markers for large numbers of individuals, either
57 through whole-genome sequencing or complexity reduction approaches such as restriction site-
58 associated DNA sequencing (RADseq) (Baird et al., 2008) in the case of large and complex
59 genomes. Increasing the density and number of markers for GPAS has the potential to increase
60 QTL detection power and prediction accuracy (de Roos et al., 2009; Long & Langley, 1999).
61 Despite the declining costs of sequencing, genotyping every individual of every population
62 across a landscape is usually not feasible, and phenotyping remains resource-consuming. As a
63 cost-effective alternative to sequencing individuals (Indi-seq), sequencing pools of individuals
64 (Pool-seq) (Schlötterer et al., 2014) has gained popularity in ecology (Bastide et al., 2013;
65 Cheng et al., 2012; Nielsen et al., 2018), evolution (Boitard et al., 2012; Fournier-Level et al.,
66 2019; Fracassetti et al., 2015), and breeding (Beissinger et al., 2014; Bélanger et al., 2016)
67 supported by developments in quantitative genetics (Fournier-Level et al., 2017; Guo et al.,
68 2018; Knight, Saccone et al., 2009; Macgregor et al., 2006; Micheletti & Narum, 2018; Jinliang
69 Yang et al., 2015).

70

71 Indi-seq generates high-resolution genomic data of a population; while Pool-seq yields low-
72 resolution data in favour of cost reduction. Indi-seq yields individual allele information after
73 variant calling, while Pool-seq generates allele frequency estimates for a group of individuals.
74 Identifying when best to use one over the other is important. Pool-seq was shown to be at least

75 as accurate as Indi-seq in estimating genome-wide allele frequencies (Fracassetti et al., 2015;
76 Gautier et al., 2013; Rellstab et al., 2013; Zhu et al., 2012), but it is also prone to biases in
77 genome representation (i.e. unequal representation of different genomic regions) when sample
78 size and depth of coverage are low (i.e. <40 individuals per pool and <50X depth) (Cutler &
79 Jensen, 2010; Schlötterer et al., 2014). Pool-seq also loses haplotype and linkage disequilibrium
80 (LD) information (Fariello et al., 2017) which limits the number of quantitative and population
81 genetics models that can be used and requires the design of novel analysis methods (Cutler &
82 Jensen, 2010). Pool-seq is more cost-effective than Indi-seq since it requires less sequencing
83 effort to generate the same genome-wide allele frequency data (Futschik & Schlötterer, 2010;
84 Gautier et al., 2013), particularly for non-model organisms where individuals cannot be
85 maintained indefinitely and used in multiple experiments. Additionally, Pool-seq can include
86 more individuals, grouped into one or a few pools and sequenced at a high depth.

87

88 There is no research on the optimal sampling strategy across a landscape for GPAS, comparing
89 Indi-seq and Pool-seq. Quantitative genetics studies have established that large sample sizes
90 and diverse mapping and training populations improve the power of GWAS (Visscher et al.,
91 2012; Visscher et al., 2017; Gurdasani et al., 2019; Wojcik et al., 2019) and the accuracy of GP
92 (Asoro et al., 2011; Bernal-Vasquez et al., 2014; Bustos-Korts et al., 2016; Rincant et al., 2017;
93 Akdemir & Isidro-Sánchez, 2019; Edwards et al., 2019). However, the spatial component of
94 sampling across the landscape is lacking in these studies. In landscape genomics, in addition to
95 the recommendation to sample as many populations as possible (Santos & Gaiotto, 2020), the
96 spatial extent of the species' dispersal should be accounted for (Riginos et al., 2016) and
97 stratified sampling is commonly performed to represent the different environmental clines (Hoel,

98 1943; Li et al., 2017; Williams & Brown, 2019; Pais et al., 2020). Landscape genomics
99 associates genotype with the environment (Li et al., 2017); while GPAS associates genotype
100 with phenotype (Gondro et al. 2013). It follows that in GPAS, sampling every environmental cline
101 may not be needed as long as the genetic diversity is well-represented in the samples.
102 Additionally, landscapes and populations with low differentiation or structure allow for less
103 samples to represent the total diversity sufficiently. Finally, it is not known whether the cost-
104 effectiveness of Pool-seq, being the preferred genotyping approach in landscape genomics over
105 Indi-seq (Santos & Gaiotto, 2020), retains the same cost-effectiveness in landscape-wide GPAS.

106

107 Field researchers often adapt techniques initially developed for model organisms or crops in
108 highly controlled environments; however, with natural populations having evolved in natural
109 environments, devising the optimal sampling strategy becomes non-trivial. Individuals and pools
110 can be sampled from one to a few populations or from a large number of populations. Identifying
111 which populations warrant higher resolution (i.e. individual-resolution genotype data with Indi-
112 seq instead of group-resolution genotype data with Pool-seq), requires some prior knowledge of
113 the spatial distribution of genetic variability across the landscape. To address this question, we
114 simulated landscapes under different trait architectures and population genetics scenarios with
115 the aim of providing recommendations on the optimal sampling strategies. Specifically, we aim
116 to answer the following three questions. How many populations do we need to sample to yield
117 optimal GPAS performance? Under which landscape-specific circumstances should we use Indi-
118 seq or Pool-seq? And which populations to select under different landscape scenarios?

119

120 **Materials and methods**

121 **Workflow overview**

122 We first simulated landscapes inhabited by multiple populations of an outbreeding species. For
123 every simulated landscape, migration between adjacent pairs of populations is uniform, and
124 there is one trait of interest which is positively correlated with fitness. This trait is controlled by
125 QTL with purely additive effects. The distribution across the landscape of the corresponding
126 favourable alleles is affected by their population of origin, selection intensity and migration rate.
127 We then simulated a stratified sampling strategy and performed Indi-seq and Pool-seq. Next,
128 GPAS models were built to identify the QTL and predict the phenotypes of individuals and pools.
129 Finally, we assessed how the number of populations sampled, the genotyping strategy (Indi-seq
130 and Pool-seq), selection intensity, migration rate, and the distribution of the favourable alleles
131 across the landscape in relation to the sampled populations affected QTL detection and genomic
132 prediction accuracies.

133

134 **Landscape simulations**

135 The landscapes were simulated using quantiNemo2 (Neuenschwander et al., 2018), the
136 variables are listed in Table 1, and the fixed parameters are listed in Table S1. We simulated the
137 simple and common mating system consisting of a hermaphrodite species capable of both self-
138 and cross-fertilisation ($1/n$ and $1-1/n$ probabilities, respectively; n is the population size) with
139 discrete non-overlapping generations. This sufficiently represents allogamous and mixed mating
140 systems, including that of plants with considerable outcrossing, and many animals species.

141 Variation for a quantitative phenotype over a landscape was simulated as a function of migration
142 rate, number of QTL controlling the trait, causal allele diffusion gradient (the distribution of the
143 favourable alleles across the landscape as it migrates from the populations of origin), and
144 selection intensity with 3 levels for each of these variables (Table 1).

145

146 Each landscape consisted of 100 populations arrayed in a uniform square lattice without
147 barriers. Migration was modelled using a 2-dimensional stepping-stone model with bidirectional
148 gene flow with a uniform rate into the 8 adjacent populations and absorbing boundaries.

149

150 The quantitative trait was determined by additive QTL with effects sampled from a χ^2 distribution
151 with 1 degree of freedom to generate a cumulative heritability of 0.5. At the initial step of the
152 simulation, all the causal alleles had a frequency (q_0) of 0.01 in the populations of origin and 0
153 elsewhere. Under the uniform allele diffusion gradient, all populations had $q_0=0.01$. Under the
154 unidirectional gradient, one boundary row had $q_0=0.01$ for all 10 populations in that row and 0
155 elsewhere. Under the bidirectional gradient, two opposite boundary rows had $q_0=0.01$ in each of
156 the populations and 0 elsewhere. For clarity, these causal allele diffusion gradients are
157 illustrated in Figure 1 Panel A.

158

159 Selection was simulated using a generalised logistic model (Richards, 1959) as

$$1/w = 1 + e^{y_c - y}, \quad \text{eq. 1}$$

161 where w is fitness, y is the quantitative trait ($y \in \mathbb{R}$), $y_c = y_{min} + s(y_{max} - y_{min})$, with y_{max} and y_{min} the
162 maximum and minimum possible trait values, and s is the selection intensity. This selection
163 intensity variable is the minimum relative phenotype value (ranging from 0 to 1) for which

164 survival rate is at least 50%, assuming that the phenotype is positively correlated with fitness. It
165 follows that a selection intensity of 0.50 or 0.95 means that individuals at the top 50th or 95th
166 percentile, respectively, contribute more to the next generation than those below. We used high
167 selection intensities, i.e. $s \in \{0.5, 0.90, 0.95\}$, to simulate pressing anthropogenic selection
168 pressures, e.g. herbicide or insecticide pressure.

169

170 The 10,000 biallelic loci were randomly distributed across a large genome with 7 chromosomes
171 and a total length of 2×10^9 base-pairs and 750 centimorgans. Two hundred generations were
172 simulated to allow the causal alleles to migrate across the landscape and generate the allele
173 frequency distributions specific to each of the three causal allele diffusion gradients. Phenotype
174 values in the final generation were scaled in the 0 to 1 interval and used for the GPAS
175 experiments.

176

177 **Genome-wide association and trait prediction based on polygenic scores**

178 GPAS was performed on all the populations. Indi-seq and Pool-seq data were assumed to have
179 been generated without genotyping error. We used established tools for Indi-seq data, and
180 developed a suite of tools for Pool-seq data. For Indi-seq data, 384 individuals were sampled per
181 population, simulating four 96-well sample plates or a single 384-well plate commonly used in
182 high-throughput molecular biology workflows. For Pool-seq data, 5 pools per population were
183 sampled, where each pool consisted of 100 individuals. This corresponds to a high power design
184 that was shown to be optimal to capture QTL association (Fournier-Level et al., 2017). GPAS
185 models were trained within each population sampled and cross-validated on all other
186 populations to assess prediction accuracies. For each population, 384 individuals and 5 pools

187 were used to train and validate the models. This corresponds to an external cross-validation with
188 equally-sized mutually exclusive training and validation sets.

189

190 Allele effects were estimated using 6 Indi-seq-based GPAS (Indi-GPAS) and 3 Pool-seq-based
191 GPAS (Pool-GPAS) models. The 6 Indi-GPAS models consisted of: efficient mixed-model
192 association expedited model (EMMAX; Kang et al., 2010), genome-wide complex trait analysis
193 (GCTA; Jiang et al., 2019), and genome-wide efficient mixed-model analysis (GEMMA; Zhou &
194 Stephens, 2012), in combination with 2 types of genetic relationship matrices: GCTA-derived
195 sparse genetic relationship matrix (GRM; off diagonals <0.05 were set to zero; Zaitlen et al.,
196 2013) and GEMMA-derived standardised relatedness matrix (STD; Zhou & Stephens, 2012) (i.e.
197 EMMAX(GRM), EMMAX(STD), GCTA(GRM), GCTA(STD), GEMMA(GRM), and
198 GEMMA(STD)). The 3 Pool-GPAS models consisted of the genome-wide estimation of additive
199 effects based on trait quantile distribution from Pool-seq data (GWA α ; Fournier-Level et al.,
200 2017), and linear mixed models (LMM) with random pairwise genetic covariance matrix
201 determined by F_{ST} derived using either Hivert's (Hivert et al., 2018) or Weir and Cockerham's
202 method (Weir & Cockerham, 1984) (i.e. GWA α , LMM (Fst Hivert), and LMM (Fst Weir &
203 Cockerham)-). The variance components of the LMM were estimated using restricted maximum
204 likelihood.

205

206 Phenotype predictions were derived from polygenic scores, i.e. the sum of the products of
207 estimated allele effects and allele dosages for Indi-GPAS or allele frequencies for Pool-GPAS.
208 For the Indi-GPAS models and GWA α , this involved a two-step approach. For each training
209 set, the polygenic scores of the training set (s_{train}) were calculated as:

210
$$s_{\text{train}} = X_{\text{train}} \beta, \quad \text{eq. 2}$$

211 where X_{train} is the dosage or frequency of alleles in the training set, and β is the vector of
212 estimated SNP effects. The polygenic scores and actual phenotype values have a linear
213 relationship (Figure S1; mean adjusted $R^2=0.97\pm 0.0031$ with 1,000 individuals per population for
214 Indi-seq and mean adjusted $R^2=0.99\pm 0.0006$ with 5 pools per population for Pool-seq) as
215 expected under the additive model used to simulate the phenotypes. These polygenic scores
216 were regressed against the actual phenotype values of the training set (y_{train}),

217
$$y_{\text{train}} = \alpha_0 + \alpha_1 s_{\text{train}}, \quad \text{eq. 3}$$

218 where α_0 is the intercept, and α_1 is the slope. The polygenic scores of the validation set,
219 $s_{\text{valid}} = X_{\text{valid}} \beta$, were transformed into the predicted phenotype values ($y_{\text{predicted}}$) using

220
$$y_{\text{predicted}} = \alpha_0 + \alpha_1 s_{\text{valid}}. \quad \text{eq. 4}$$

221

222 For the Pool-GPAS linear mixed models, the predicted phenotypes ($y_{\text{predicted}}$) were calculated as:

223
$$y_{\text{predicted}} = X_{\text{valid}} \beta, \quad \text{eq. 5}$$

224 where X_{valid} is the matrix of allele frequencies of the validation set, and β is the estimated allelic
225 effects from the GPAS model built using the training set. The trained models were validated on
226 all populations in the landscape.

227

228 GPAS performance was measured using three GWAS metrics, and one phenotype prediction
229 metric. The GWAS metrics were:

230 1 area under the receiver operating curve (AUC) (Fawcett, 2006),

231 2 true positive rate (TPR) which was defined as the fraction of causal QTL with a
232 significantly associated SNP within 1 kbp, and
233 3 false positive rate (FPR) which was defined as the fraction of the significantly associated
234 SNPs with no causal QTL within 1 kbp, unless it tags a true QTL through a chain of
235 associated SNPs each less than 1kb apart; multiple associated SNPs within 1kbp were
236 counted as one.

237 The family-wise type I error rate was set at $\alpha=0.05$. The metric for phenotype prediction is the
238 root mean square error (RMSE) between actual and predicted phenotype values:

$$239 \text{ RMSE} = \sqrt{\sum \frac{(y - \hat{y})^2}{n}}, \quad \text{eq. 6}$$

240 where y is the actual phenotypes, \hat{y} is the predicted phenotypes and n is the number of
241 observations.

242

243 **Sampling strategy optimisation**

244 A total of 405 landscapes were simulated, corresponding to all combinations of the 4 landscape
245 variables with 3 levels each and 5 replicates (Table 1). For each landscape, Indi-GPAS and
246 Pool-GPAS experiments were performed for each population independently. This constitutes the
247 intra-population dataset. The landscape was divided into equally sized rectangular regions, and
248 the approximately central population was selected from each region to simulate a stratified
249 sampling strategy. This is illustrated in Figure 1 Panel B. This constitutes the inter-population
250 dataset. AUC and RMSE were averaged across the populations sampled. TPR and FPR were
251 calculated using the cumulative number of true and false positive candidate loci across the

252 populations sampled. AUC was used to measure the accuracy of QTL detection per population,
253 while TPR and FPR were used to measure QTL detection accuracy of multiple populations.

254

255 The single best performing modelling framework was identified for each genotyping scheme
256 (Indi-GPAS and Pool-GPAS) based on AUC and RMSE for independent populations tests using
257 Tukey's honest significant difference (HSD mean comparison) at $\alpha=0.05$.

258

259 How many populations do we need to sample to yield optimal GPAS performance?

260 To determine how many populations to sample to yield optimal GPAS performance, we used the
261 inter-population dataset. We assessed the suitability of the four metrics (i.e. mean AUC, mean
262 RMSE, TPR and FPR) to address this question by visualising their relationships with the number
263 of populations sampled. Additionally, we compared the expected performance of Indi-GPAS and
264 Pool-GPAS under the same sequencing capacity constraint. The Indi-GPAS experiments we
265 simulated included 384 individuals per population, while Pool-GPAS included only 5 pools per
266 population. Assuming a 5X sequencing depth per individual for Indi-seq (Brouard et al., 2017)
267 and the recommended 50X depth per pool for Pool-seq (Schlötterer et al., 2014), these equate
268 to a sequencing depth of 1,920X per base per population for Indi-seq and only 250X for Pool-
269 seq. This means that for the sequencing capacity required to characterise one population
270 through Indi-seq, approximately 7 populations ($\lfloor 1920/250 \rfloor$) can be characterised through Pool-
271 seq.

272

273 Under which landscape-specific circumstances should we use Indi-seq or Pool-seq?

274 The second question we addressed was which landscape-specific circumstances warrant Indi-
275 GPAS or Pool-GPAS? Specifically, if we were to perform GPAS on one population, which
276 sequencing strategy (Indi-seq or Pool-seq) is better, and how does the optimal choice vary with
277 the polygenicity of the trait, selection intensity, and gene flow? The intra-population dataset was
278 analysed using AUC and RMSE as the GPAS performance metrics.

279

280 Which populations to select under different landscape scenarios?

281 To determine which populations to select to best capture the genetic basis of a trait and yield
282 accurate trait predictions under the 3 causal allele diffusion gradients, we analysed the intra-
283 population dataset using AUC and RMSE as the GPAS performance metrics. The populations
284 were classified into 10 groups, where each group represents a row perpendicular to the causal
285 allele diffusion gradient (refer to Figure 1A). The top row corresponds to populations 1 to 10, the
286 second row to populations 11 to 20, and so on. The general and landscape variable-specific
287 trends in GPAS performance across the landscape were visualised using violin plots and means
288 compared using Tukey's HSD ($\alpha=0.05$). Linear mixed models fitted linear and quadratic
289 relationships (using second degree polynomial fit) between GPAS performance and the row
290 groups. The row group was treated as a numeric variable, and nested within each level of the
291 variables: number of QTL, selection intensity, migration rate, and GPAS model.

292

293 **Implementation**

294 The landscapes were simulated using quantiNemo2 (Neuenschwander et al., 2018). The
295 genome and QTL information were simulated in R (R Core Team, 2018). The quantiNemo2
296 outputs were parsed using R and Julia (Nardelli et al., 2018). GEMMA (Zhou & Stephens, 2012),
297 EMMAX (Kang et al., 2010), GCTA (Jiang et al., 2019), and Plink (Purcell et al., 2007) were
298 used for Indi-GPAS. [GWAlpha.jl](#) was used for Pool-GPAS. The R package [violinplotter](#) was used
299 to generate violin plots with HSD mean comparison grouping. The GNU shell (Free Software
300 Foundation, 2016), Spartan (Lafayette & Wiebelt, 2017), Slurm (Yoo et al., 2003), and GNU
301 parallel (Tange, 2011) were used extensively. The workflow is available in the github repository:
302 <https://github.com/jeffersonparil/GPAS-landscape-simulation.git>.

303

304 **Results**

305 **GPAS model selection and the effects of landscape and sampling variables**

306 GEMMA (STD) and GWAlpha showed the best GPAS performances, with >79% AUC and
307 <5.9% RMSE (Table S2). Therefore, these two frameworks were selected as the representatives
308 of Indi-GPAS and Pool-GPAS models, respectively. Overall, Indi-GPAS performed better than
309 Pool-GPAS.

310

311 Factors increasing statistical power to identify causal loci through GPAS included a lower
312 number of QTL controlling the trait, more intense selection, higher migration among populations,
313 and more populations sampled (Figure S2). Accuracy in phenotype predictions improved as the
314 number of QTL controlling the trait increases, as selection intensity decreases, and as migration

315 rate increases (Figure S3). Accuracy is unaffected by the number of populations sampled since
316 each model was trained independently for each population. In addition, power and accuracy are
317 higher when QTL diffuses across the landscape uniformly.

318

319 **How many populations to sample and when to use Indi-seq or Pool-seq?**

320 TPR and FPR increase logarithmically as the number of populations sampled increases (Figure
321 2), so there is no optimum based on these metrics.

322

323 Indi-GPAS achieves greater power than Pool-GPAS at the cost of a higher false positive rate
324 (Figure 2). However, Pool-GPAS can outperform Indi-GPAS under the assumptions detailed in
325 the materials and methods section, where for every population characterised with Indi-seq,
326 approximately 7 populations can be characterised with Pool-seq. Under this 1:7 ratio, Indi-GPAS
327 on 10 populations yield an average TPR of 0.388 and FPR of 0.0150; for the same sequencing
328 capacity Pool-GPAS can be performed on 70 populations, yielding an average TPR of 0.418 and
329 FPR of 0.0115. We explored a range of ratios deviating from this 1:7 ratio. This is because the
330 5X depth requirement for variant calling in Indi-seq and 50X depth for allele frequency estimation
331 in Pool-seq depend on the species of interest and the resources available. Lower ratios, e.g. 1:8
332 to 1:10, mean even more populations can be characterised with Pool-seq for every population
333 characterised with Indi-seq. Using our simulated data to explore various ratios, we find that there
334 exists a range where Pool-GPAS can outperform Indi-GPAS, i.e. TPR is higher and FPR is lower
335 for Pool-GPAS than Indi-GPAS (Figure 3). This shows that characterising more of the landscape
336 at low resolution can be better than characterising a small portion of the landscape at high
337 resolution.

338

339 If we were to perform GPAS on one population, Indi-GPAS is better than Pool-GPAS. However
340 in cases where selection intensity is high (i.e. 0.90 to 0.95) or migration rate is high (i.e. 0.01)
341 Pool-GPAS performance is not significantly different from Indi-GPAS in terms of prediction
342 accuracy (Figure 4).

343

344 **Which populations to select under different landscape scenarios?**

345 GPAS performance is maximised in populations with high genetic variability which at the
346 landscape level, means sampled close to the place of origin of the causal allele (Figure 5). This
347 area of high genetic variability is characterised by intermediate causal allele frequencies which
348 translate into populations with high phenotypic variability. In the absence of a causal allele
349 diffusion gradient (i.e. uniform causal allele distribution), no row seems to be optimal for
350 sampling, except for some slightly better performance from populations in the middle rows.
351 Under unidirectional gradient (i.e. causal alleles originated from the top row and diffused
352 downwards hence a single diffusion front) and in terms of QTL detection accuracy, sampling the
353 populations from the top row is optimal; however, in terms of prediction accuracy, the
354 populations in the middle rows appear to be better. Under bidirectional gradient (i.e. causal
355 alleles originated from the top and bottom rows hence two diffusion fronts) both QTL detection
356 and prediction accuracies are optimal in the populations from the top and bottom rows. These
357 trends across the landscape correlate with the trends in the mean number of polymorphic QTL
358 per population and causal allele frequencies.

359

360 In the presence of causal allele diffusion gradients, the relationship of the sampling location
361 (defined as rows perpendicular to the diffusion gradient) with both QTL detection and prediction
362 accuracies appears to be quadratic, except for QTL detection accuracy under unidirectional
363 causal allele diffusion, for which the relationship is linear (Figure 5). In terms of GWAS accuracy
364 as measured by AUC, sampling near the diffusion fronts becomes less important (i.e. slope
365 under unidirectional gradient and curvature under bidirectional gradient are reduced) as the
366 number of QTL increases, as selection intensity decreases, and as migration rate increases
367 (Figure 6 columns 1-3). In addition, sampling near the diffusion fronts is more important for Pool-
368 GPAS than Indi-GPAS (Figure 6 column 4), in other words, power diminishes quicker for Pool-
369 GPAS than Indi-GPAS as we move away from areas of high diversity.

370

371 In terms of prediction accuracy as measured by RMSE, the degree to which the middle rows (i.e.
372 areas of high genetic and phenotypic variability) are the optimal sampling locations under
373 unidirectional diffusion decreases (i.e. curvature becomes less severe) as the number of QTL
374 increases, as selection intensity decreases, and as migration rate increases (Figure 7 top
375 graphs). Also, sampling from the middle rows under unidirectional diffusion is slightly more
376 important for Indi-GPAS than Pool-GPAS. On the other hand, the degree to which the top and
377 bottom rows are optimal under bidirectional diffusion decreases (i.e. curvature becomes less
378 severe) as the number of QTL, selection intensity, and migration rate increase (Figure 7 bottom
379 graphs). Also, sampling from the top and bottom rows under bidirectional diffusion is more
380 important for Pool-GPAS than Indi-GPAS, in other words, similar to that of power, accuracy
381 diminishes quicker for Pool-GPAS than Indi-GPAS as we move away from areas of high
382 diversity.

383

384 The trends in GPAS performance across the landscape correlate with the trends in genetic
385 variability (expressed in terms of causal allele frequency, i.e. frequencies closer to 0.5 indicates
386 higher diversity; Figures S4 to S9). Opposite trends are observed between causal allele diffusion
387 gradients for RMSE as selection intensity increases, i.e. in the middle of the landscape, RMSE is
388 minimised for unidirectional causal diffusion but maximised for bidirectional causal allele
389 diffusion (Figure S5 and S8).

390

391 **Discussion**

392 GPAS has the potential to extend the scope of genomic studies in ecology and evolution beyond
393 environment association and niche modelling (Dormann et al., 2012; Exposito-Alonso et al.,
394 2018; Fournier-Level et al., 2011). The predicted environmental range of individual genotypes
395 determined through genome-environment associations (Manel et al., 2018) can be
396 complemented by the phenotype predictions of GPAS (Cotto et al., 2017). This can be
397 transformational for the way we monitor invasive species or assess the adaptive potential of
398 endangered ones. We provide recommendations for optimising the sampling strategy to
399 maximise the power to detect QTL and the accuracy of quantitative phenotype prediction in
400 natural populations of any outbreeding species. We stress the importance of capturing sufficient
401 representation of the genetic variability present over the landscape by sampling populations from
402 areas of high genetic diversity. On a per population basis or if only one population were to be
403 sampled, we recommend using Indi-seq over Pool-seq. We have not considered phenotyping
404 costs here, but if it is high, then it would increase the attractiveness of Indi-seq to maximise

405 information per unit cost. However, similar to a study on the estimation of population
406 differentiation (Goudet & Büchi, 2006) and a meta analysis of several landscape genomics
407 studies (Santos & Gaiotto, 2020), we demonstrated the value of shallow but extensive genotyping
408 with Pool-seq in maximising the number of populations that can be analysed without
409 compromising power. This is especially true if the aim is to predict phenotypes of some future
410 populations for the rapid and timely monitoring of invasive and threatened species.

411

412 **How many populations to sample?**

413 Our results emphasise the need to sample as many populations as possible from regions of high
414 genetic diversity. The power to detect QTL is maximised if all the populations in the landscape
415 were included in the study. This is possible for endangered species with a small number of
416 populations in the wild. However this is not feasible for species with a healthier number of
417 populations. The best populations to sample are located in areas of high genetic diversity which
418 manifests as areas with high trait variability where the causal alleles are at intermediate
419 frequencies.

420

421 Our results show a diminishing return in terms of GPAS power when increasing the number of
422 populations sampled. This is consistent with a study by Selmoni et al. (2020) on sampling
423 strategy optimisation for landscape genomics which found that sampling an intermediate number
424 of sites can perform as well as maximising the number of sites sampled. This study only
425 considered Indi-seq and tested different sample sizes per population, and our approach is
426 comparable because using Indi-seq equates to a high-resolution characterisation of the
427 landscape and Pool-seq to a low-resolution one. Our analysis extends this result further because

428 it is independent of the number of individuals sampled per population. The cost-effectiveness of
429 Pool-seq allows for more populations to be sampled and included in the study than Indi-seq.

430

431 **When to use Indi-seq or Pool-seq?**

432 The number of individuals per population selected for our Indi-GPAS simulations (384
433 individuals) exceeds the sample size of most ecological studies (e.g. 100-200 individual samples
434 per population in birds (Hansson et al., 2018; Perrier et al., 2018), <100 samples per population
435 in trees (Cappa et al., 2013; Holliday et al., 2010), ~100 samples per population in mammals
436 (Johnston et al., 2011; Pallares et al., 2014), and <20 samples per population in fish (Willing et
437 al., 2010)). Thus for most experiments, the power of Indi-GPAS is expected to be lower than in
438 our simulations. On the contrary, the power of Pool-GPAS is expected to remain the same since
439 five pools per population was found to be optimal (Fournier-Level et al., 2017) and each pool can
440 include a non-limiting number of individuals. The sequencing capacity required for Indi-seq is
441 always higher than for Pool-seq, and more populations can be characterised with Pool-seq than
442 with Indi-seq under the same budgetary constraints (Schlötterer et al., 2014). Therefore, the
443 range of the number of populations sampled where Pool-GPAS outperforms Indi-GPAS is likely
444 to be even broader than reported here, as long as the genomic characterisation approach yields
445 accurate genomic data.

446

447 Sufficient depth of coverage is required to correct for sequencing errors inherent to current high-
448 throughput sequencing technologies. Indi-seq can provide high-resolution genomic information
449 for a population, but comes at a high cost. A given genomic region needs to be sequenced at
450 least 5 times for each individual to correct for sequencing errors and yield accurate basecalling

451 information (Brouard et al., 2017). Additionally, many individuals are required to accurately
452 represent a population. This is only resource-effective when the individuals are part of an
453 association panel and the genomic information can be leveraged across several research
454 projects (Robin et al., 2019). On the other hand, Pool-seq generates low-resolution genomic
455 information on a population that is cost-effective while maintaining high power. To maximise the
456 accuracy of allele frequency estimation (i.e. to correct for sequencing errors and sufficiently
457 represent the pool), Pool-seq guidelines and recommendations have been proposed (Schlötterer
458 et al. 2014; Fracassetti et al. 2015; Anand et al. 2016). Pool sizes of at least 25 (Gautier et al.
459 2013; Fracassetti et al. 2015) to 40 individuals (Schlötterer et al. 2014), and depths greater than
460 50X (Zhu et al. 2012) to 65X (Gautier et al. 2013) have been recommended. Hundreds of
461 individuals can be pooled to yield accurate allele frequency data (Schlötterer et al., 2014). This
462 means that in an outbreeding species, a few pools consisting of hundreds of individuals each
463 can represent a population better than a few individuals.

464

465 Pool-seq is more widely used than Indi-seq in ecological and evolutionary studies because the
466 focus is generally on populations rather than individuals, and because of its cost-effectiveness
467 (Futschik & Schlötterer, 2010). In contrast, Cutler and Jensen (2010) concluded that Indi-seq
468 should be preferred over Pool-seq for many applications due to the loss of haplotype and LD
469 information. They focused on applications in human and model organisms, whereas Pool-seq
470 has its highest impact for high-throughput data acquisition in non-model species of critical
471 ecological and economical importance.

472

473 **Which populations to select under different landscape scenarios?**

474 We have shown that sampling from genetically diverse populations maximises GPAS
475 performance. Capturing greater genetic diversity was shown to increase the power to detect
476 causal loci (Alqudah et al., 2020; Rosenberg et al., 2010; Wojcik et al., 2019). Similarly,
477 populations which represent the overall diversity found in the landscape or are similar to the
478 validation populations, improve prediction accuracies (Akdemir & Isidro-Sánchez, 2019; Asoro et
479 al., 2011; Edwards et al., 2019). Populations with high genetic diversity were found along the
480 diffusion fronts, i.e. the areas where the causal alleles migrate from their site of origin into the
481 neighbouring populations. The rate at which GPAS performance decreases as we sample
482 farther away from the diffusion fronts correlates with the decrease in genetic diversity at the
483 causal loci. In the absence of prior genomic information, the areas of high genetic diversity
484 coincide with regions of high phenotypic diversity. Gaining prior information on the location of
485 these areas of high genetic diversity and causal allele diffusion fronts or more broadly the
486 landscape of adaptive genetic diversity (Eckert & Dyer, 2012) is key to an optimal sampling
487 strategy.

488

489 When the causal allele diffusion gradient is unknown and a uniform causal allele distribution is
490 assumed, there is a small advantage in choosing populations in the middle of the landscape.
491 This can be explained by the absorbing boundaries used in the migration model which simulates
492 a restricted range whereby alleles going beyond the border are lost. This reflects the
493 phenomenon in fringe populations where migration regularly occurs beyond the suitable
494 environmental niche of the species and the migrants fail to survive (Sexton et al., 2009). This is
495 expected to apply to organisms with restricted range such as corals (Guan et al., 2015), 24% of

496 coral reef fish species (Hawkins et al., 2006), and organisms in hydrothermal vents (Mullineaux
497 et al., 2018). However, in cases where non-uniform causal allele distribution is assumed (e.g.
498 temperature response in lodgepole pine and interior spruce (Liepe et al., 2015); and white
499 spruce (Hornoy et al., 2015)), these populations along the borders of the species range are
500 important to sample because they may carry unique advantageous variants, especially in the
501 context of climate change adaptation (Geber & Eckhart, 2005; Bridle & Vines, 2007; le Roux ,
502 2009; Pais et al., 2020).

503

504 When considering the genetic architecture of a trait, there is less power to detect QTL but higher
505 prediction accuracy for highly polygenic traits than for traits controlled by fewer loci. As a
506 consequence, if we expect the polygenic trait to be controlled by relatively few loci, and we are
507 more focused on mapping its genetic basis than on phenotype prediction, then sufficient power
508 to detect QTL can be achieved with less populations sampled. Biotic stress resistance traits
509 were often shown to be oligogenic, for example resistance to *Pseudomonas syringae* in
510 *Arabidopsis thaliana* (Atwell et al., 2010), and resistance to *Cronartium ribicola* in sugar pine
511 (Weiss et al., 2020). On the other hand, if we expect the trait to be highly polygenic, and we are
512 more focused on phenotype prediction for monitoring purposes, then we have to sample as
513 many populations as possible to sufficiently represent the variation across the landscape and
514 maximise prediction accuracy. This is typically the case for abiotic stress resistance traits such
515 as aluminium and proton tolerance in *Arabidopsis thaliana* (Nakano, et al., 2020), and coral
516 bleaching resistance in *Acropora millepora* (Fuller et al., 2020) As the number of loci controlling
517 the trait increases, the selection pressure acting on each locus decreases (Walsh & Lynch,
518 2018). This reduces the power to detect QTL since the individual contribution of each QTL

519 decreases as more loci control the trait (Wang & Xu, 2019). This in turn reduces the proportion
520 of polymorphic QTL within populations: if the majority of the QTL have small effects, they have a
521 higher chance of getting lost due to drift than QTL with large effects. On the other hand,
522 prediction accuracy increases since the rate at which genetic variance decreases due to
523 directional selection is reduced as the number of QTL increases. Genetic variance should
524 eventually become zero under constant stabilising or directional selection, but the rate of this
525 reduction becomes slower with an increased number of loci controlling the trait (Crow & Kimura,
526 1970). Hence, GWAS and GP complement each other to achieve high QTL detection accuracy
527 or high prediction accuracy for quantitative traits controlled by any number of loci.

528

529 When considering selection intensity, there is higher power to detect QTL but lower prediction
530 accuracy in populations under intense selection. As a consequence, populations under high
531 selection pressure can be selected for high-power GWAS: for example, weed and insect pest
532 populations in agricultural areas including herbicide-resistant *Lolium rigidum* (Powles et al.,
533 1998), and Bt-resistant *Helicoverpa armigera* (Jin et al., 2018). On the other hand, populations
534 under low selection pressure can be selected for GP. For example, non-target organisms
535 including *Drosophila melanogaster* populations which are resistant to Bt toxins (Babin et al.,
536 2020) and imidacloprid (Fournier-Level et al., 2019) Increasing the selection intensity increases
537 QTL detection power, since the effect of individual QTL becomes greater within each population
538 (Wang & Xu, 2019) and less QTL alleles are lost due to drift. However, the predictive ability will
539 be reduced by the Bulmer effect (Bulmer, 1971), where covariance between loci (partially
540 explained by linkage disequilibrium (Walsh & Lynch, 2018)) is reduced after selection. Increasing
541 selection intensity magnifies this reduction resulting in diminished additive genetic variance and

542 less predictive models. This further solidifies the complementary nature of GWAS and GP and
543 the utility of performing both with GPAS.

544

545 When considering migration rate, there is more power to detect QTL and greater prediction
546 accuracy in populations experiencing high migration than in reproductively isolated ones. As a
547 consequence, if gene flow is high, then less populations need to be sampled to sufficiently
548 represent the variation across the landscape. This is the case for highly mobile organisms
549 including birds (Pulido, 2007); and fishes (Brodersen et al., 2008), On the other hand, if gene
550 flow is low which results in highly structured landscapes, then more populations need to be
551 sampled to sufficiently represent the landscape variation. This can be the case in landscapes
552 with considerable natural or artificial barriers to migration, and in species with reproductive
553 structures impeding outbreeding, e.g. cleistogamous plants including *Crotalaria micans*
554 (Etcheverry et al., 2003) and *Vigna caracalla* (Etcheverry et al., 2008). Increasing migration rate
555 decreases differentiation between populations allowing for more causal alleles to be shared This
556 leads to higher additive genetic variance per population, resulting in higher power and more
557 accurate predictions (Liu et al., 2020).

558

559 **Conclusion**

560 Genome-wide association studies and genomic prediction (GPAS) are poised to complement
561 existing methodologies in ecology in evolution. GPAS provides powerful tools to dissect the
562 genetic basis of ecologically important quantitative traits including fitness and to rapidly monitor
563 natural populations including invasive and threatened species. Understanding how the number

564 of population samples and the different landscape properties affect the QTL detection power and
565 phenotype prediction accuracies is integral to planning population collections for GPAS
566 experiments. We recommend sampling as many populations as possible from areas of high
567 genetic diversity. We also recommend Pool-seq whenever Indi-seq is too costly; since sampling
568 more populations at the cost of lower resolution can be better than characterising a small
569 number of populations at high resolution. The complementary nature of GWAS and GP allows
570 good QTL detection power or prediction accuracy under low to high trait polygenicity and
571 selection intensity. In the absence of prior information on the areas of high genetic diversity, we
572 recommend against sampling populations at the border of the species' range.

573

574 **Acknowledgements**

575 The authors wish to thank Uli Felzmann from IT services, Faculty of Science, The University of
576 Melbourne for his outstanding support. We acknowledge the extensive use of the Spartan High
577 Performance Computing and the Melbourne Research Cloud systems. Part of this work was
578 supported through the Computational Biology Research Initiative Seed Fund awarded to AF-L,
579 and the David Hay Fund.

580

581 **Data accessibility**

582 Codes to reproduce the simulation data and the subsequent analysis are publicly available on
583 github: <https://github.com/jeffersonparil/GPAS-landscape-simulation>.

584

585 **Author Contributions**

586 JFP, DJB and AFL designed the study. JFP analysed the data. JFP and AFL wrote the
587 manuscript with input from DJB.

588 **References**

- 589
590 Akdemir, D., & Isidro-Sánchez, J. (2019). Design of training populations for selective
591 phenotyping in genomic prediction. *Scientific Reports*, 9(1), 1446. doi: 10.1038/s41598-
592 018-38081-6
593
- 594 Alqudah, A. M., Sallam, A., Stephen Baenziger, P., & Börner, A. (2020). GWAS: Fast-forwarding
595 gene identification and characterization in temperate Cereals: lessons from Barley – A
596 review. *Journal of Advanced Research*, 22, 119–135. doi: 10.1016/j.jare.2019.10.013
597
- 598 Asoro, F. G., Newell, M. A., Beavis, W. D., Scott, M. P., & Jannink, J.-L. (2011). Accuracy and
599 Training Population Design for Genomic Selection on Quantitative Traits in Elite North
600 American Oats. *The Plant Genome*, 4(2), 132–144. doi:
601 10.3835/plantgenome2011.02.0007
602
- 603 Atwell, Susanna, Yu S. Huang, Bjarni J. Vilhjálmsson, Glenda Willems, Matthew Horton, Yan Li,
604 Dazhe Meng, et al. (2010). “Genome-Wide Association Study of 107 Phenotypes in
605 Arabidopsis Thaliana Inbred Lines.” *Nature* 465, no. 7298: 627–31.
606 <https://doi.org/10.1038/nature08800>.
607
- 608 Baird, N. A., Etter, P. D., Atwood, T. S., Currey, M. C., Shiver, A. L., Lewis, Z. A., ... Johnson, E.
609 A. (2008). Rapid SNP Discovery and Genetic Mapping Using Sequenced RAD Markers.
610 *PLoS ONE*, 3(10), e3376. doi: 10.1371/journal.pone.0003376
611
- 612 Bastide, H., Betancourt, A., Nolte, V., Tobler, R., Stöbe, P., Futschik, A., & Schlotterer, C.
613 (2013). A Genome-Wide, Fine-Scale Map of Natural Pigmentation Variation in *Drosophila*
614 *melanogaster*. *PLoS Genetics*, 9(6), e1003534. doi: 10.1371/journal.pgen.1003534
615
- 616 Beissinger, T. M., Hirsch, C. N., Vaillancourt, B., Deshpande, S., Barry, K., Buell, C. R., ... de
617 Leon, N. (2014). A genome-wide scan for evidence of selection in a maize population
618 under long-term artificial selection for ear number. *Genetics*, 196(3), 829–40. doi:
619 10.1534/genetics.113.160655
620
- 621 Bélanger, S., Esteves, P., Clermont, I., Jean, M., & Belzile, F. (2016). Genotyping-by-
622 Sequencing on Pooled Samples and its Use in Measuring Segregation Bias during the
623 Course of Androgenesis in Barley. *The Plant Genome*, 9(1), 0. doi:
624 10.3835/plantgenome2014.10.0073
625
- 626 Bernal-Vasquez, AM., Möhring, J., Schmidt, M., Schönleben, M., Schön, S.-C., & Piepho, H.-P.
627 (2014). The importance of phenotypic data analysis for genomic prediction - a case study
628 comparing different spatial models in rye. *BMC Genomics* 15, 646 .
629 <https://doi.org/10.1186/1471-2164-15-646>
630
- 631 Boitard, S., Schlotterer, C., Nolte, V., Pandey, R. V., & Futschik, A. (2012). Detecting Selective
632 Sweeps from Pooled Next-Generation Sequencing Samples. *Molecular Biology and*

633 *Evolution*, 29(9), 2177–2186. doi: 10.1093/molbev/mss090
634

635 Bridle, Jon R., & Vines, Timothy H., (2007). Limits to evolution at range margins: when and why
636 does adaptation fail?. *Trends in Ecology & Evolution*. Volume 22, Issue 3, pp 140-147.
637 ISSN 0169-5347. <https://doi.org/10.1016/j.tree.2006.11.002>.
638

639 Brodersen, Jakob, P. Anders Nilsson, Lars-Anders Hansson, Christian Skov, & Christer
640 Brönmark. (2008). "Condition-Dependent Individual Decision-Making Determines
641 Cyprinid Partial Migration." *Ecology* 89, no. 5: 1195–1200. <https://doi.org/10.1890/07-1318.1>.
642
643

644 Brouard, J.-S., Boyle, B., Ibeagha-Awemu, E. M., & Bissonnette, N. (2017). Low-depth
645 genotyping-by-sequencing (GBS) in a bovine population: strategies to maximize the
646 selection of high quality genotypes and the accuracy of imputation. *BMC Genetics*, 18.
647 doi: 10.1186/s12863-017-0501-y
648

649 Bulmer, M. G. (1971). The Effect of Selection on Genetic Variability. *The American Naturalist*.
650 Retrieved from <https://www.journals.uchicago.edu/doi/10.1086/282718>
651

652 Bustos-Korts, D., Malosetti, M., Chapman S., Biddulph, B., & van Eeuwijk, F. (2016).
653 Improvement of Predictive Ability by Uniform Coverage of the Target Genetic Space. *G3*
654 *Genes|Genomes|Genetics*, Volume 6, Issue 11, 1 November 2016, Pages 3733–3747,
655 <https://doi.org/10.1534/g3.116.035410>

656 Cappa, E. P., El-Kassaby, Y. A., Garcia, M. N., Acuña, C., Borralho, N. M. G., Grattapaglia, D.,
657 & Poltri, S. N. M. (2013). Impacts of Population Structure and Analytical Models in
658 Genome-Wide Association Studies of Complex Traits in Forest Trees: A Case Study in
659 *Eucalyptus globulus*. *PLOS ONE*, 8(11), e81267. doi: 10.1371/journal.pone.0081267
660

661 Cheng, C., White, B. J., Kamdem, C., Mockaitis, K., Costantini, C., Hahn, M. W., & Besansky, N.
662 J. (2012). Ecological genomics of *Anopheles gambiae* along a latitudinal cline: A
663 population-resequencing approach. *Genetics*, 190(4), 1417–1432. doi:
664 10.1534/genetics.111.137794
665

666 Cotto, O., Wessely, J., Georges, D., Klöner, G., Schmid, M., Dullinger, S., Thuiller, W., &
667 Guillaume, F.. (2017). A dynamic eco-evolutionary model predicts slow response of
668 alpine plants to climate warming. *Nat Commun* 8, 15399.
669 <https://doi.org/10.1038/ncomms15399>
670

671 Crow, J. F., & Kimura, M. (1970). *An Introduction to Population Genetics Theory*. doi:
672 10.2307/1529706
673

674 Cutler, D. J., & Jensen, J. D. (2010). To pool, or not to pool? *Genetics*, 186(1), 41–3. doi:
675 10.1534/genetics.110.121012
676

677 de Roos, A. P. W., Hayes, B. J., & Goddard, M. E. (2009). Reliability of genomic predictions
678 across multiple populations. *Genetics*, 183(4), 1545–53. doi:

679 10.1534/genetics.109.104935
680
681 Dormann, C. F., Schymanski, S. J., Cabral, J., Chuine, I., Graham, C., Hartig, F., ... Singer, A.
682 (2012). Correlation and process in species distribution models: bridging a dichotomy.
683 *Journal of Biogeography*, 39(12), 2119–2131. doi: 10.1111/j.1365-2699.2011.02659.x
684
685 Eckert, A. J., & Dyer, R. J. (2012). Defining the landscape of adaptive genetic diversity.
686 *Molecular Ecology*, 21(12), 2836–2838. doi: 10.1111/j.1365-294X.2012.05615.x
687
688 Edwards, S. M., Buntjer, J. B., Jackson, R., Bentley, A. R., Lage, J., Byrne, E., ... Hickey, J. M.
689 (2019). The effects of training population design on genomic prediction accuracy in
690 wheat. *Theoretical and Applied Genetics*, 132(7), 1943–1952. doi: 10.1007/s00122-019-
691 03327-y
692
693 Etcheverry, Angela Virginia, Maria Mercedes Alemán, and Trinidad Figueroa Fleming. (2008).
694 “Flower Morphology, Pollination Biology and Mating System of the Complex Flower of
695 *Vigna caracalla* (Fabaceae: Papilionoideae).” *Annals of Botany* 102, no. 3 (September 1,
696 2008): 305–16. <https://doi.org/10.1093/aob/mcn106>.
697
698 Etcheverry, A. V., J. J. Protomastro, and C. Westerkamp. (2003). “Delayed Autonomous Self-
699 Pollination in the Colonizer *Crotalaria Micans* (Fabaceae: Papilionoideae): Structural and
700 Functional Aspects.” *Plant Systematics and Evolution* 239, no. 1/2 (2003): 15–28.
701
702 Exposito-Alonso, M., Vasseur, F., Ding, W., Wang, G., Burbano, H. A., & Weigel, D. (2018).
703 Genomic basis and evolutionary potential for extreme drought adaptation in *Arabidopsis*
704 *thaliana*. *Nature Ecology & Evolution*, 2(2), 352–358. doi: 10.1038/s41559-017-0423-0
705
706 Fariello, M. I., Boitard, S., Mercier, S., Robelin, D., Faraut, T., Arnould, C., ... SanCristobal, M.
707 (2017). Accounting for linkage disequilibrium in genome scans for selection without
708 individual genotypes: The local score approach. *Molecular Ecology*, 26(14), 3700–3714.
709 doi: 10.1111/mec.14141
710
711 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8), 861–
712 874. doi: 10.1016/j.patrec.2005.10.010
713
714 Fournier-Level, A., Good, R.T., Wilcox, S.A. Rane, R.V., Schiffer, M., Chen, W., Battlay, P.,
715 Perry, T., Batterham, P., Hoffmann, A.A., & Robin, C. (2019). The spread of resistance to
716 imidacloprid is restricted by thermotolerance in natural populations of *Drosophila*
717 *melanogaster*. *Nat Ecol Evol* 3, 647–656. <https://doi.org/10.1038/s41559-019-0837-y>
718
719 Fournier-Level, A., Korte, A., Cooper, M. D., Nordborg, M., Schmitt, J., & Wilczek, A. M. (2011).
720 A Map of Local Adaptation in *Arabidopsis thaliana*. *Science*, 334(6052), 86–89. doi:
721 10.1126/science.1209271
722
723 Fournier-Level, Alexandre, Robin, C., & Balding, D. J. (2017). GWAlpha: Genome-wide
724 estimation of additive effects (alpha) based on trait quantile distribution from pool-

725 sequencing experiments. *Bioinformatics*. doi: 10.1093/bioinformatics/btw805
726

727 Fracassetti, M., Griffin, P. C., & Willi, Y. (2015). Validation of Pooled Whole-Genome Re-
728 Sequencing in *Arabidopsis lyrata*. *PloS One*, *10*(10), e0140462. doi:
729 10.1371/journal.pone.0140462
730

731 Free Software Foundation. (2016). *GNU bash*. Retrieved from
732 <https://www.gnu.org/software/bash/>
733

734 Fuller, Zachary L., Veronique J. L. Mocellin, Luke A. Morris, Neal Cantin, Jihanne Shepherd,
735 Luke Sarre, Julie Peng, et al. (2020) "Population Genetics of the Coral *Acropora*
736 *Millepora*: Toward Genomic Prediction of Bleaching." *Science* 369, no. 6501 (July 17,
737 2020). <https://doi.org/10.1126/science.aba4674>.
738

739 Futschik, A., & Schlötterer, C. (2010). The next generation of molecular markers from massively
740 parallel sequencing of pooled DNA samples. *Genetics*, *186*(1), 207–18. doi:
741 10.1534/genetics.110.114397
742

743 Gautier, M., Foucaud, J., Gharbi, K., Cézard, T., Galan, M., Loiseau, A., ... Estoup, A. (2013).
744 Estimation of population allele frequencies from next-generation sequencing data: pool-
745 versus individual-based genotyping. *Molecular Ecology*, *22*(14), 3766–3779. doi:
746 10.1111/mec.12360
747

748 Geber, M.A. & Eckhart, V.M. (2005). Experimental studies of adaptation in *Clarkia xantiana*: II.
749 Fitness variation across a subspecies border. *Evolution*, *59*: 521-531.
750 <https://doi.org/10.1111/j.0014-3820.2005.tb01012.x>
751

752 Gondro, C., van der Werf, J., & Hayes, B. (Eds.). (2013). *Genome-wide association studies and*
753 *genomic prediction* (1st ed.). Humana Press. Doi: 10.1007/978-1-62703-447-0_20
754

755 Goudet, J., & Büchi. (2006). The effects of dominance, regular inbreeding and sampling design
756 on QST, an estimator of population differentiation for quantitative traits. *Genetics* *172*,
757 1337-1347. <https://doi.org/10.1534/genetics.105.050583>.
758

759 Guan, Y., Hohn, S., & Merico, A. (2015). Suitable environmental ranges for potential coral reef
760 habitats in the tropical ocean. *PloS one*, *10*(6), e0128831.
761 <https://doi.org/10.1371/journal.pone.0128831>
762

763 Guo, X., Cericola, F., Fè, D., Pedersen, M. G., Lenk, I., Jensen, C. S., ... Janss, L. L. (2018).
764 Genomic Prediction in Tetraploid Ryegrass Using Allele Frequencies Based on
765 Genotyping by Sequencing. *Frontiers in Plant Science*, *9*, 1165. doi:
766 10.3389/fpls.2018.01165
767

768 Gurdasani, D., Barroso, I., Zeggini, E. et al. (2019). Genomics of disease risk in globally diverse
769 populations. *Nat Rev Genet* *20*, 520–535. <https://doi.org/10.1038/s41576-019-0144-0>

770
771 Hansson, B., Sigeman, H., Stervander, M., Tarka, M., Ponnikas, S., Strandh, M., ... Hasselquist,
772 D. (2018). Contrasting results from GWAS and QTL mapping on wing length in great
773 reed warblers. *Molecular Ecology Resources*, 18(4), 867–876. doi: 10.1111/1755-
774 0998.12785
775
776 Hawkins, J.P., Roberts, C.M. & Clark, V. (2000), The threatened status of restricted–range coral
777 reef fish species. *Animal Conservation*, 3: 81-88. [https://doi.org/10.1111/j.1469-
778 1795.2000.tb00089.x](https://doi.org/10.1111/j.1469-1795.2000.tb00089.x)
779
780 Hornoy, Benjamin, Nathalie Pavy, Sébastien Gérardi, Jean Beaulieu, & Jean Bousquet. (2015).
781 “Genetic Adaptation to Climate in White Spruce Involves Small to Moderate Allele
782 Frequency Shifts in Functionally Diverse Genes.” *Genome Biology and Evolution* 7, no.
783 12: 3269–85. <https://doi.org/10.1093/gbe/evv218>.
784
785 Hivert, V., Leblois, R., Petit, E. J., Gautier, M., & Vitalis, R. (2018). Measuring Genetic
786 Differentiation from Pool-seq Data. *Genetics*, 210(1), 315–330. doi:
787 10.1534/genetics.118.300900
788
789 Hoel, P. (1943). The Accuracy of Sampling Methods in Ecology on JSTOR. *The Annals of*
790 *Mathematical Statistics*, 14(3), 289–300.
791
792 Holliday, J. A., Ritland, K., & Aitken, S. N. (2010). Widespread, ecologically relevant genetic
793 markers developed from association mapping of climate-related traits in Sitka spruce
794 (*Picea sitchensis*). *New Phytologist*, 188(2), 501–514. doi: 10.1111/j.1469-
795 8137.2010.03380.x
796
797 Jiang, L., Zheng, Z., Qi, T., Kemper, K. E., Wray, N. R., Visscher, P. M., & Yang, J. (2019). A
798 resource-efficient tool for mixed model association analysis of large-scale data. *Nature*
799 *Genetics*, 51(12), 1749–1755. doi: 10.1038/s41588-019-0530-8
800
801 Johnston, S. E., McEWAN, J. C., Pickering, N. K., Kijas, J. W., Beraldi, D., Pilkington, J. G., ...
802 Slate, J. (2011). Genome-wide association mapping identifies the genetic basis of
803 discrete and quantitative variation in sexual weaponry in a wild sheep population.
804 *Molecular Ecology*, 20(12), 2555–2566. doi: 10.1111/j.1365-294X.2011.05076.x
805
806 Kang, H. M., Sul, J. H., Service, S. K., Zaitlen, N. A., Kong, S., Freimer, N. B., ... Eskin, E.
807 (2010). Variance component model to account for sample structure in genome-wide
808 association studies. *Nature Genetics*, 42(4), 348–354. doi: 10.1038/ng.548
809
810 Knight, J., Saccone, S. F., Zhang, Z., Ballinger, D. G., & Rice, J. P. (2009). A Comparison of
811 Association Statistics between Pooled and Individual Genotypes. *Human Heredity*, 67(4),
812 219–225. doi: 10.1159/000194975
813
814 Lafayette, L., & Wiebelt, B. (2017). Spartan and NEMO: Two HPC-Cloud Hybrid
815 Implementations. *2017 IEEE 13th International Conference on E-Science (e-Science)*,

816 458–459. Auckland: IEEE. doi: 10.1109/eScience.2017.70
817

818 Le Roux P. (2009) Plants at the margin. Ecological limits and climate change. *Ann Bot.* 104(7):ix.
819 doi:10.1093/aob/mcp220
820

821 Li, Y., Zhang, X.-X., Mao, R.-L., Yang, J., Miao, C.-Y., Li, Z., & Qiu, Y.-X. (2017). Ten Years of
822 Landscape Genomics: Challenges and Opportunities. *Frontiers in Plant Science*, 8,
823 2136. doi: 10.3389/fpls.2017.02136
824

825 Liepe, Katharina J., Andreas Hamann, Pia Smets, Connor R. Fitzpatrick, & Sally N. Aitken.
826 (2016). “Adaptation of Lodgepole Pine and Interior Spruce to Climate: Implications for
827 Reforestation in a Warming World.” *Evolutionary Applications* 9, no. 2: 409–19.
828 <https://doi.org/10.1111/eva.12345>.
829

830 Liu, L., Wang, Y., Zhang, D., Chen, Z., Chen, X., Su, Z., & He, X. (2020). The Origin of Additive
831 Genetic Variance Driven by Positive Selection. *Molecular Biology and Evolution*, 37(8),
832 2300–2308. doi: 10.1093/molbev/msaa085
833

834 Long, A. D., & Langley, C. H. (1999). The Power of Association Studies to Detect the
835 Contribution of Candidate Genetic Loci to Variation in Complex Traits. *Genome*
836 *Research*, 9(8), 720–731. doi: 10.1101/gr.9.8.720
837

838 Macgregor, S., Visscher, P. M., & Montgomery, G. (2006). Analysis of pooled DNA samples on
839 high density arrays without prior knowledge of differential hybridization rates. *Nucleic*
840 *Acids Research*, 34(7). doi: 10.1093/nar/gkl136
841

842 Manel, S., Andrello, M., Henry, K., Verdelet, D., Darracq, A., Guerin, P.-E., ... Devaux, P.
843 (2018). Predicting genotype environmental range from genome–environment
844 associations. *Molecular Ecology*, 27(13), 2823–2833. doi: 10.1111/mec.14723
845

846 Micheletti, S. J., & Narum, S. R. (2018). Utility of pooled sequencing for association mapping in
847 nonmodel organisms. *Molecular Ecology Resources*, 18(4), 825–837. doi: 10.1111/1755-
848 0998.12784
849

850 Mullineaux, Lauren S., Anna Metaxas, Stace E. Beaulieu, Monika Bright, Sabine Gollner,
851 Benjamin M. Grupe, Santiago Herrera, et al. (2018). “Exploring the Ecology of Deep-Sea
852 Hydrothermal Vents in a Metacommunity Framework.” *Frontiers in Marine Science* 5.
853 <https://doi.org/10.3389/fmars.2018.00049>.
854

855 Neuenschwander, S., Michaud, F., Goudet, J., & Stegle, O. (2018). quantiNemo 2: a swiss knife
856 to simulate complex demographic and genetic scenarios, forward and backward in time.
857 *Bioinformatics*. doi: 10.1093/bioinformatics/bty737
858

859 Nielsen, E. S., Henriques, R., Toonen, R. J., Knapp, I. S. S., Guo, B., & von der Heyden, S.
860 (2018). Complex signatures of genomic variation of two non-model marine species in a
861 homogeneous environment. *BMC Genomics*, 19(1), 347. doi: 10.1186/s12864-018-4721-

862 y
863

864 Pais, A.L., Whetten, R.W. & Xiang, Q.-Y. (2020). Population structure, landscape genomics,
865 and genetic signatures of adaptation to exotic disease pressure in *Cornus florida* L. -
866 Insights from GWAS and GBS data. *J. Syst. Evol.*, 58: 546-570.
867 <https://doi.org/10.1111/jse.12592>
868

869 Pallares, L. F., Harr, B., Turner, L. M., & Tautz, D. (2014). Use of a natural hybrid zone for
870 genomewide association mapping of craniofacial traits in the house mouse. *Molecular*
871 *Ecology*, 23(23), 5756–5770. doi: 10.1111/mec.12968
872

873 Perrier, C., Delahaie, B., & Charmantier, A. (2018). Heritability estimates from genomewide
874 relatedness matrices in wild populations: Application to a passerine, using a small
875 sample size. *Molecular Ecology Resources*, 18(4), 838–853. doi: 10.1111/1755-
876 0998.12886
877

878 Powles, Stephen B, Debrah F Lorraine-Colwill, James J Dellow, & Christopher Preston. (1998).
879 “Evolved Resistance to Glyphosate in Rigid Ryegrass (*Lolium rigidum*) in Australia.”
880 *Weed Science* 46, no. 5: 604–7.
881

882 Pulido, Francisco. (2007) “The Genetics and Evolution of Avian Migration.” *BioScience* 57, no. 2:
883 165–74. <https://doi.org/10.1641/B570211>.
884

885 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A. R., Bender, D., ... Sham, P.
886 C. (2007). PLINK: A Tool Set for Whole-Genome Association and Population-Based
887 Linkage Analyses. *The American Journal of Human Genetics*, 81(3), 559–575. doi:
888 10.1086/519795
889

890 R Core Team. (2018). *R: A Language and Environment for Statistical Computing*. Vienna,
891 Austria. Retrieved from <https://www.r-project.org/>
892

893 Rellstab, C., Zoller, S., Tedder, A., Gugerli, F., & Fischer, M. C. (2013). Validation of SNP Allele
894 Frequencies Determined by Pooled Next-Generation Sequencing in Natural Populations
895 of a Non-Model Plant Species. *PLoS ONE*, 8(11), e80422. doi:
896 10.1371/journal.pone.0080422
897

898 Richards, F. J. (1959). A flexible growth function for empirical use. *Journal of Experimental*
899 *Botany*, 10(2), 290–301. doi: 10.1093/jxb/10.2.290
900

901 Riginos, C., Crandall, E.D., Liggins, L., Bongaerts, P., & Trembl, E.A. (2016). Navigating the
902 currents of seascape genomics: how spatial analyses can augment population genomic
903 studies. *Current Zoology*, Volume 62, Issue 6, December 2016, Pages 581–601,
904 <https://doi.org/10.1093/cz/zow067>
905

906 Rincent, R., Charcosset, A., & Moreau, L. (2017). Predicting genomic selection efficiency to
907 optimize calibration set and to assess prediction accuracy in highly structured

908 populations. *TAG. Theoretical and applied genetics. Theoretische und angewandte*
909 *Genetik*, 130(11), 2231–2247. <https://doi.org/10.1007/s00122-017-2956-7>
910

911 Robin, C., Battlay, P., & Fournier-Level, A. (2019). What can genetic association panels tell us
912 about evolutionary processes in insects? *Current Opinion in Insect Science*, 31, 99–105.
913 doi: 10.1016/j.cois.2018.12.004
914

915 Rosenberg, N. A., Huang, L., Jewett, E. M., Szpiech, Z. A., Jankovic, I., & Boehnke, M. (2010).
916 Genome-wide association studies in diverse populations. *Nature Reviews Genetics*,
917 11(5), 356–366. doi: 10.1038/nrg2760
918

919 Santos, A. S., & Gaiotto, F. A. (2020). Knowledge status and sampling strategies to maximize
920 cost-benefit ratio of studies in landscape genomics of wild plants. *Scientific Reports*,
921 10(1), 1–9. doi: 10.1038/s41598-020-60788-8
922

923 Schlötterer, C., Tobler, R., Kofler, R., & Nolte, V. (2014). Sequencing pools of individuals -
924 mining genome-wide polymorphism data without big funding. *Nature Reviews*, 15, 749–
925 765.
926

927 Selmoni, O., Vajana, E., Guillaume, A., Rochat, E., & Joost, S. (2020). Sampling strategy
928 optimization to increase statistical power in landscape genomics: A simulation-based
929 approach. *Molecular Ecology Resources*, 20(1), 154–169. doi: 10.1111/1755-0998.13095
930

931 Sexton, J. P., McIntyre, P. J., Angert, A. L., & Rice, K. J. (2009). Evolution and Ecology of
932 Species Range Limits. *Annual Review of Ecology, Evolution, and Systematics*, 40(1),
933 415–436. doi: 10.1146/annurev.ecolsys.110308.120317
934

935 Tange, O. (2011). *GNU Parallel - The Command-Line Power Tool*. The USENIX Magazine.
936

937 Visscher, P.M., Brown, M.A., McCarthy, M.I., & Yang, J. (2012). Five years of GWAS discovery.
938 *Am. J. Hum. Genet.* 90. pp 7-24
939

940 Visscher, P.M., Wray, N.R., Zhang, Q., Sklar, P., McCarthy, M.I., Brown, M.A., & Yang, J.
941 (2017). 10 Years of GWAS Discovery: Biology, Function, and Translation. *The American*
942 *Journal of Human Genetics. Volume 101, Issue 1.* pp 5-22, ISSN 0002-9297.
943 <https://doi.org/10.1016/j.ajhg.2017.06.005>.
944

945 Walsh, B., & Lynch, M. (2018). *Evolution and selection of quantitative traits* (1st ed.). Oxford
946 University Press. Retrieved from
947 http://nitro.biosci.arizona.edu/zbook/NewVolume_2/newvol2.html#2B
948

949 Wang, M., & Xu, S. (2019). Statistical power in genome-wide association studies and
950 quantitative trait locus mapping. *Heredity*, 123(3), 287–306. doi: 10.1038/s41437-019-
951 0205-3
952

953 Ward, Sarah M., & Marie Jasieniuk. (2009). Review: Sampling Weedy and Invasive Plant

- 954 Populations for Genetic Diversity Analysis. *Weed Science* 57, no. 6 : 593-602.
955 <http://www.jstor.org/stable/40586878>.
956
- 957 Weiss, Matthew, Richard A. Snieszko, Daniela Puiu, Marc W. Crepeau, Kristian Stevens, Steven
958 L. Salzberg, Charles H. Langley, David B. Neale, and Amanda R. De La Torre. (2020).
959 "Genomic Basis of White Pine Blister Rust Quantitative Disease Resistance and Its
960 Relationship with Qualitative Resistance." *The Plant Journal* 104, no. 2: 365–76.
961 <https://doi.org/10.1111/tpj.14928>.
962
- 963 Weir, B. S., & Cockerham, C. C. (1984). Estimating F-Statistics for the Analysis of Population
964 Structure Author (s): B . S . Weir and C . Clark Cockerham Reviewed work (s) :
965 Published by : Society for the Study of Evolution Stable URL : [http://www.jstor.org/stable/](http://www.jstor.org/stable/2408641)
966 2408641 . *Evolution*, 38(6), 1358–1370. doi: 128.103.149.52
967
- 968 Williams, B. K., & Brown, E. D. (2019). Sampling and analysis frameworks for inference in
969 ecology. *Methods in Ecology and Evolution*, 10(11), 1832–1842. doi: 10.1111/2041-
970 210X.13279
971
- 972 Willing, E.-M., Bentzen, P., Oosterhout, C. V., Hoffmann, M., Cable, J., Breden, F., ... Dreyer, C.
973 (2010). Genome-wide single nucleotide polymorphisms reveal population history and
974 adaptive divergence in wild guppies. *Molecular Ecology*, 19(5), 968–984. doi:
975 10.1111/j.1365-294X.2010.04528.x
976
- 977 Wojcik, G. L., Graff, M., Nishimura, K. K., Tao, R., Haessler, J., Gignoux, C. R., ... Carlson, C. S.
978 (2019). Genetic analyses of diverse populations improves discovery for complex traits.
979 *Nature*, 570(7762), 514–518. doi: 10.1038/s41586-019-1310-4
980
- 981 Yang, Jinliang, Jiang, H., Yeh, C.-T., Yu, J., Jeddloh, J. A., Nettleton, D., & Schnable, P. S.
982 (2015). Extreme-phenotype genome-wide association study (XP-GWAS): a method for
983 identifying trait-associated variants by sequencing pools of individuals selected from a
984 diversity panel. *The Plant Journal*, 84(3), 587–596. doi: 10.1111/tpj.13029
985
- 986 Yoo, A. B., Jette, M. A., & Grondona, M. (2003). SLURM: Simple Linux Utility for Resource
987 Management. In D. Feitelson, L. Rudolph, & U. Schwiegelshohn (Eds.), *Job Scheduling*
988 *Strategies for Parallel Processing* (pp. 44–60). Berlin, Heidelberg: Springer Berlin
989 Heidelberg.
990
- 991 Zaitlen, N., Kraft, P., Patterson, N., Pasaniuc, B., Bhatia, G., Pollack, S., & Price, A. L. (2013).
992 Using Extended Genealogy to Estimate Components of Heritability for 23 Quantitative
993 and Dichotomous Traits. *PLoS Genetics*, 9(5). doi: 10.1371/journal.pgen.1003520
994
- 995 Zappa Nardelli, F., Belyakova, J., Pelenitsyn, A., Chung, B., Bezanson, J., & Vitek, J. (2018).
996 Julia Subtyping: A Rational Reconstruction. *Proc. ACM Program. Lang.*, 2(OOPSLA),
997 113:1–113:27. doi: 10.1145/3276483
998
- 999 Zhou, X., & Stephens, M. (2012). Genome-wide efficient mixed-model analysis for association

1000 studies. *Nature Genetics*, 44(7), 821–824. doi: 10.1038/ng.2310

1001

1002 Zhu, Y., Bergland, A. O., González, J., & Petrov, D. A. (2012). Empirical Validation of Pooled

1003 Whole Genome Population Re-Sequencing in *Drosophila melanogaster*. *PLoS ONE*,

1004 7(7), e41901. doi: 10.1371/journal.pone.0041901