# Measuring Genetic Differentiation from Pool-seq Data

Valentin Hivert, Raphael Leblois, Eric Petit, Mathieu Gautier, Renaud Vitalis

HAL Id: hal-01936905

https://hal.science/hal-01936905

Submitted on 26 May 2020

# Measuring genetic differentiation from Pool-seq data

**Valentin Hivert**[*,†]**, Raphaël Leblois**[*,†]**, Eric J. Petit**[‡]**, Mathieu Gautier**[*,†,§]**, and Renaud Vitalis**[*,†,§]

[*]CBGP, INRA, CIRAD, IRD, Montpellier SupAgro, Univ Montpellier,

Montpellier, France

[†]Institut de Biologie Computationnelle, Univ Montpellier, Montpellier, France

[‡]ESE, Ecology and Ecosystem Health, INRA, Agrocampus Ouest, Rennes, France

[§]These authors are joint senior authors on this work

1

**Running title**: Genetic differentiation from pools

**Keywords**: $F_{\mathrm{ST}}$, genetic differentiation, pool sequencing, population genomics

**Corresponding author**: Renaud Vitalis

Centre de Biologie pour la Gestion des Populations

Campus International de Baillarguet, CS 30 016

34988 Montferrier-sur-Lez cedex

France

Tel : +33 (0)4 99 62 33 42

Fax : +33 (0)4 99 62 33 45

E-mail: `renaud.vitalis@inra.fr`

2

## Abstract

1

2     The advent of high throughput sequencing and genotyping tech-

3     nologies enables the comparison of patterns of polymorphisms at a

4     very large number of markers. While the characterization of genetic

5     structure from individual sequencing data remains expensive for many

6     non-model species, it has been shown that sequencing pools of indi-

7     vidual DNAs (Pool-seq) represents an attractive and cost-effective al-

8     ternative. However, analyzing sequence read counts from a DNA pool

9     instead of individual genotypes raises statistical challenges in deriving

10    correct estimates of genetic differentiation. In this article, we pro-

11    vide a method-of-moments estimator of $F_{ST}$ for Pool-seq data, based

12    on an analysis-of-variance framework. We show, by means of simula-

13    tions, that this new estimator is unbiased, and outperforms previously

14    proposed estimators. We evaluate the robustness of our estimator to

15    model misspecification, such as sequencing errors and uneven contri-

16    butions of individual DNAs to the pools. Finally, by reanalyzing pub-

17    lished Pool-seq data of different ecotypes of the prickly sculpin *Cottus*

18    *asper*, we show how the use of an unbiased $F_{ST}$ estimator may ques-

19    tion the interpretation of population structure inferred from previous

20    analyses.

3

# INTRODUCTION

It has long been recognized that the subdivision of species into subpopulations, social groups and families fosters genetic differentiation (Wahlund 1928; Wright 1931). Characterizing genetic differentiation as a means to infer unknown population structure is therefore fundamental to population genetics, and finds applications in multiple domains, including conservation biology, invasion biology, association mapping and forensics, among many others. In the late 1940s and early 1950s, Malécot (1948) and Wright (1951) introduced $F$-statistics to partition genetic variation within and between groups of individuals (Holsinger and Weir 2009; Bhatia et al. 2013). Since then, the estimation of $F$-statistics has become standard practice (see, e.g., Weir 1996; Weir and Hill 2002; Weir 2012), and the most commonly used estimators of $F_{\mathrm{ST}}$ have been developed in an analysis-of-variance framework (Cockerham 1969, 1973; Weir and Cockerham 1984), which can be recast in terms of probabilities of identity of pairs of homologous genes (Cockerham and Weir 1987; Rousset 2007; Weir and Goudet 2017).

Assuming that molecular markers are neutral, estimates of $F_{\mathrm{ST}}$ are typically used to quantify genetic structure in natural populations, which is then interpreted as the result of demographic history (Holsinger and Weir 2009): large $F_{\mathrm{ST}}$ values are expected for small populations among which dispersal is limited (Wright 1951), or between populations that have long diverged in isolation from each other (Reynolds et al. 1983); when dispersal is spatially restricted, a positive relationship between $F_{\mathrm{ST}}$ and the geographical distance for pairs of populations generally holds (Slatkin 1993; Rousset 1997). It has also been proposed to characterize the heterogeneity of $F_{\mathrm{ST}}$

estimates across markers for identifying loci that are targeted by selection (Cavalli-Sforza 1966; Lewontin and Krakauer 1973; Beaumont and Nichols 1996; Vitalis et al. 2001; Akey et al. 2002; Beaumont 2005; Weir et al. 2005; Lotterhos and Whitlock 2014, 2015; Whitlock and Lotterhos 2015).

Next-generation sequencing (NGS) technologies provide unprecedented amounts of polymorphism data in both model and non-model species (Ellegren 2014). Although the sequencing strategy initially involved individually tagged samples in humans (The International HapMap Consortium 2005), whole-genome sequencing of pools of individuals (Pool-seq) is being increasingly used for population genomic studies (Schlötterer et al. 2014). Because it consists in sequencing libraries of pooled DNA samples and does not require individual tagging of sequences, Pool-seq provides genome-wide polymorphism data at considerably lower cost than sequencing of individuals (Schlötterer et al. 2014). However, non-equimolar amounts of DNA from all individuals in a pool and stochastic variation in the amplification efficiency of individual DNAs have raised concerns with respect to the accuracy of the so-obtained allele frequency estimates, particularly at low sequencing depth and with small pool sizes (Cutler and Jensen 2010; Ellegren 2014; Anderson et al. 2014). Nonetheless, it has been shown that, at equal sequencing effort, Pool-seq provides similar, if not more accurate, allele frequency estimates than individual-based analyses (Futschik and Schlötterer 2010; Gautier et al. 2013). The problem is different for diversity and differentiation parameters, which depend on second moments of allele frequencies or, equivalently, on pairwise measures of genetic identity: with Pool-seq data, it is indeed impossible to distinguish pairs of reads that are identical because they were

5

sequenced from a single gene, from pairs of reads that are identical because

they were sequenced from two distinct genes that are identical in state (IIS)

(Ferretti et al. 2013).

Appropriate estimators of diversity and differentiation parameters must

therefore be sought, to account for both the sampling of individual genes

from the pool and the sampling of reads from these genes. There has been

several attempts to define estimators for the parameter $F_{\mathrm{ST}}$ for Pool-seq data

(Kofler et al. 2011; Ferretti et al. 2013), from ratios of heterozygosities (or

from probabilities of genetic identity between pairs of reads) within and be-

tween pools. In the following, we will argue that these estimators are biased

(i.e., they do not converge towards the expected value of the parameter),

and that some of them have undesired statistical properties (i.e., the bias

depends upon sample size and coverage). Here, following Cockerham (1969),

Cockerham (1973), Weir and Cockerham (1984), Weir (1996), Weir and Hill

(2002) and Rousset (2007), we define a method-of-moments estimator of the

parameter $F_{\mathrm{ST}}$ using an analysis-of-variance framework. We then evaluate

the accuracy and the precision of this estimator, based on the analysis of sim-

ulated datasets, and compare it to estimates defined in the software package

PoPoolation2 (Kofler et al. 2011), and in Ferretti et al. (2013). Furthermore,

we test the robustness of our estimators to model misspecifications (including

unequal contributions of individuals in pools, and sequencing errors). Finally,

we reanalyze the prickly sculpin (*Cottus asper*) Pool-seq data (published by

Dennenmoser et al. 2017), and show how the use of biased $F_{\mathrm{ST}}$ estimators in

previous analyses may challenge the interpretation of population structure.

Note that throughout this article, we use the term "gene" to designate a

6

segregating genetic unit (in the sense of the "Mendelian gene" from Orgogozo

et al. 2016). We further use the term "read" in a narrow sense, as a sequenced

copy of a gene. For the sake of simplicity, we will use the term "Ind-seq" to

refer to analyses based on individual data, for which we further assume that

individual genotypes are called without error.

7

102 $F$-statistics may be described as intra-class correlations for the probability of

103 identity in state (IIS) of pairs of genes (Cockerham and Weir 1987; Rousset

104 1996, 2007), and $F_{\mathrm{ST}}$ is best defined as:

$$F_{\mathrm{ST}} \equiv \frac{Q_1 - Q_2}{1 - Q_2} \tag{1}$$

105 where $Q_1$ is the IIS probability for genes sampled within subpopulations, and

106 $Q_2$ is the IIS probability for genes sampled between subpopulations. In the

107 following, we develop an estimator of $F_{\mathrm{ST}}$ for Pool-seq data, by decomposing

108 the total variance of read frequencies in an analysis-of-variance framework.

109 A complete derivation of the model is provided in the Supplemental File S1.

110     For the sake of clarity, the notation used throughout this article is given in

111 Table 1. We first derive our model for a single locus, and eventually provide

112 a multilocus estimator of $F_{\mathrm{ST}}$. Consider a sample of $n_{\mathrm{d}}$ subpopulations, each

113 of which is made of $n_{\mathrm{i}}$ genes ($i = 1, \ldots, n_{\mathrm{d}}$) sequenced in pools (hence $n_{\mathrm{i}}$ is

114 the haploid sample size of the $i$th pool). We define $c_{ij}$ as the number of reads

115 sequenced from gene $j$ ($j = 1, \ldots, n_{\mathrm{i}}$) in subpopulation $i$ at the locus consid-

116 ered. Note that $c_{ij}$ is a latent variable, that cannot be directly observed from

117 the data. Let $X_{ijr:k}$ be an indicator variable for read $r$ ($r = 1, \ldots, c_{ij}$) from

118 gene $j$ in subpopulation $i$, such that $X_{ijr:k} = 1$ if the $r$th read from the $j$th

119 gene in the $i$th deme is of type $k$, and $X_{ijr:k} = 0$ otherwise. In the following,

120 we use standard dot notations for sample averages, i.e.: $X_{ij\cdot:k} \equiv \sum_r X_{ijr:k}/c_{ij}$,

121 $X_{i\cdot\cdot:k} \equiv \sum_j \sum_r X_{ijr:k} / \sum_j c_{ij}$ and $X_{\cdots:k} \equiv \sum_i \sum_j \sum_r X_{ijr:k} / \sum_i \sum_j c_{ij}$. The

122 analysis of variance is based on the computation of sums of squares, as fol-

123 lows:

$$\sum_{i}^{n_{\mathrm{d}}}\sum_{j}^{n_i}\sum_{r}^{c_{ij}}\left(X_{ijr:k}-X_{\cdots:k}\right)^2 = \sum_{i}^{n_{\mathrm{d}}}\sum_{j}^{n_i}\sum_{r}^{c_{ij}}\left(X_{ijr:k}-X_{ij\cdot:k}\right)^2$$

$$+ \sum_{i}^{n_{\mathrm{d}}}\sum_{j}^{n_i}\sum_{r}^{c_{ij}}\left(X_{ij\cdot:k}-X_{i\cdot\cdot:k}\right)^2$$

$$+ \sum_{i}^{n_{\mathrm{d}}}\sum_{j}^{n_i}\sum_{r}^{c_{ij}}\left(X_{i\cdot\cdot:k}-X_{\cdots:k}\right)^2$$

$$\equiv SSR_{:k}+SSI_{:k}+SSP_{:k} \tag{2}$$

124 As is shown in the Supplemental File S1, the expected sums of squares depend

125 on the expectation of the allele frequency $\pi_k$ over all replicate populations

126 sharing the same evolutionary history, as well as on the IIS probability $Q_{1:k}$

127 that two genes in the same pool are both of type $k$, and the IIS probability

128 $Q_{2:k}$ that two genes from different pools are both of type $k$. Taking expecta-

129 tions (see the detailed computations in the Supplemental File S1), one has:

$$\mathbb{E}(SSR_{:k}) = 0 \tag{3}$$

130 for reads within individual genes, since we assume that there is no sequencing

131 error, i.e. all the reads sequenced from a single gene are identical and $X_{ijr:k}=$

132 $X_{ij\cdot:k}$ for all $r$. For reads between genes within pools, we get:

$$\mathbb{E}(SSI_{:k}) = (C_1-D_2)(\pi_k-Q_{1:k}) \tag{4}$$

133 where $C_1 \equiv \sum_i \sum_j c_{ij} = \sum_i C_{1i}$ is the total number of reads in the full sample

134 (total coverage), $C_{1i}$ is the coverage of the $i$th pool and $D_2 \equiv \sum_i (C_{1i}+n_i-1)/n_i$.

135 $D_2$ arises from the assumption that the distribution of the read counts $c_{ij}$

136 is multinomial (i.e., that all genes contribute equally to the pool of reads;

9

137 see Equation A15 in Supplemental File S1). For reads between genes from

138 different pools, we have:

$$\mathbb{E}(SSP_{:k}) = \left(C_1 - \frac{C_2}{C_1}\right)(Q_{1:k} - Q_{2:k}) + (D_2 - D_2^\star)(\pi_k - Q_{1:k}) \quad (5)$$

139 where $C_2 \equiv \sum_i C_{1i}^2$ and $D_2^\star \equiv \left[\sum_i C_{1i}(C_{1i} + n_i - 1)/n_i\right]/C_1$ (see Equa-

140 tion A16 in Supplemental File S1). Rearranging Equations 4–5, and summing

141 over alleles, we get:

$$Q_1 - Q_2 = \frac{(C_1 - D_2)\mathbb{E}(SSP) - (D_2 - D_2^\star)\mathbb{E}(SSI)}{(C_1 - D_2)(C_1 - C_2/C_1)} \quad (6)$$

142 and:

$$1 - Q_2 = \frac{(C_1 - D_2)\mathbb{E}(SSP) + (n_c - 1)(D_2 - D_2^\star)\mathbb{E}(SSI)}{(C_1 - D_2)(C_1 - C_2/C_1)} \quad (7)$$

143 where $n_c \equiv (C_1 - C_2/C_1)/(D_2 - D_2^\star)$. Let $MSI \equiv SSI/(C_1 - D_2)$ and

144 $MSP \equiv SSP/(D_2 - D_2^\star)$. Then, using the definition of $F_{\mathrm{ST}}$ from Equation 1,

145 we have:

$$F_{\mathrm{ST}} \equiv \frac{Q_1 - Q_2}{1 - Q_2} = \frac{\mathbb{E}(MSP) - \mathbb{E}(MSI)}{\mathbb{E}(MSP) + (n_c - 1)\mathbb{E}(MSI)} \quad (8)$$

146 which yields the method-of-moments estimator:

$$\hat{F}_{\mathrm{ST}}^{\mathrm{pool}} = \frac{MSP - MSI}{MSP + (n_c - 1)MSI} \quad (9)$$

147 where

$$MSI = \frac{1}{C_1 - D_2}\sum_k \sum_i^{n_d} C_{1i}\hat{\pi}_{i:k}(1 - \hat{\pi}_{i:k}) \quad (10)$$

10

148 and:

$$MSP = \frac{1}{D_2 - D_2^\star} \sum_k \sum_i^{n_{\rm d}} C_{1i} \left( \hat{\pi}_{i:k} - \hat{\pi}_k \right)^2 \qquad (11)$$

149 (see Equations A25 and A26 in Supplemental File S1). In Equations 10

150 and 11, $\hat{\pi}_{i:k} \equiv X_{i\cdot\cdot:k}$ is the average frequency of reads of type $k$ within the $i$th

151 pool, and $\hat{\pi}_k \equiv X_{\cdot\cdot\cdot:k}$ is the average frequency of reads of type $k$ in the full sam-

152 ple. Note that from the definition of $X_{\cdot\cdot\cdot:k}$, $\hat{\pi}_k \equiv \sum_i \sum_j \sum_r X_{ijr:k} / \sum_i \sum_j c_{ij} =$

153 $\sum_i C_{1i}\hat{\pi}_{i:k} / \sum_i C_{1i}$ is the weighted average of the sample frequencies with

154 weights equal to the pool coverage. This is equivalent to the weighted

155 analysis-of-variance in Cockerham (1973) (see also Weir and Cockerham 1984;

156 Weir 1996; Weir and Hill 2002; Rousset 2007; Weir and Goudet 2017). Fi-

157 nally, the full expression of $\hat{F}_{\rm ST}^{\rm pool}$ in terms of sample frequencies reads:

$$\hat{F}_{\rm ST}^{\rm pool} = \frac{\sum_k \left[ (C_1 - D_2) \sum_i^{n_{\rm d}} C_{1i} \left( \hat{\pi}_{i:k} - \hat{\pi}_k \right)^2 - (D_2 - D_2^\star) \sum_i^{n_{\rm d}} C_{1i}\hat{\pi}_{i:k} \left( 1 - \hat{\pi}_{i:k} \right) \right]}{\sum_k \left[ (C_1 - D_2) \sum_i^{n_{\rm d}} C_{1i} \left( \hat{\pi}_{i:k} - \hat{\pi}_k \right)^2 + (n_{\rm c} - 1) (D_2 - D_2^\star) \sum_i^{n_{\rm d}} C_{1i}\hat{\pi}_{i:k} \left( 1 - \hat{\pi}_{i:k} \right) \right]}$$
$$(12)$$

158 If we take the limit case where each gene is sequenced exactly once, we

159 recover the Ind-seq model: assuming $c_{ij} = 1$ for all $(i,j)$, then $C_1 = \sum_i^{n_{\rm d}} n_i$,

160 $C_2 = \sum_i^{n_{\rm d}} n_i^2$, $D_2 = n_{\rm d}$ and $D_2^\star = 1$. Therefore, $n_{\rm c} = (C_1 - C_2/C_1) / (n_{\rm d} - 1)$,

161 and Equation 9 reduces exactly to the estimator of $F_{\rm ST}$ for haploids: see Weir

162 (1996), p. 182, and Rousset (2007), p. 977.

163 As in Reynolds et al. (1983), Weir and Cockerham (1984), Weir (1996)

164 and Rousset (2007), a multilocus estimate is derived as the sum of locus

165 specific numerators over the sum of locus-specific denominators:

$$\hat{F}_{\rm ST} = \frac{\sum_l MSP_l - MSI_l}{\sum_l MSP_l + (n_{\rm c} - 1) MSI_l} \qquad (13)$$

11

166 where $MSI$ and $MSP$ are subscripted with $l$ to denote the $l$th locus. For

167 Ind-seq data, Bhatia et al. (2013) refer to this multilocus estimate as a "ratio

168 of averages" by opposition to an "average of ratios", which would consist in av-

169 eraging single-locus $F_{\mathrm{ST}}$ over loci. This approach is justified in the Appendix

170 of Weir and Cockerham (1984) and in Bhatia et al. (2013), who analyzed

171 both estimates by means of coalescent simulations. Note that Equation 13

172 assumes that the pool size is equal across loci. Also note that the construc-

173 tion of the estimator in Equation 13 is different from Weir and Cockerham's

174 (1984). These authors defined their multilocus estimator as a ratio of sums

175 of components of variance ($a$, $b$ and $c$ in their notation) over loci, which give

176 the same weight to all loci, whatever the number of sampled genes at each lo-

177 cus. Equation 13 follows GENEPOP's rationale (Rousset 2008) instead, which

178 gives instead more weight to loci that are more intensively covered.

12

## Simulation study

*Generating individual genotypes:* we first generated individual genotypes using `ms` (Hudson 2002), assuming an island model of population structure (Wright 1931). For each simulated scenario, we considered 8 demes, each made of $N = 5,000$ haploid individuals. The migration rate $(m)$ was fixed to achieve the desired value of $F_{ST}$ (0.05 or 0.2), using Equation 6 in Rousset (1996) leading, e.g., to $M \equiv 2Nm = 16.569$ for $F_{ST} = 0.05$ and $M = 3.489$ for $F_{ST} = 0.20$. The mutation rate was set at $\mu = 10^{-6}$, giving $\theta \equiv 2N\mu = 0.01$. We considered either fixed, or variable sample sizes across demes. In the latter case, the haploid sample size $n$ was drawn independently for each deme from a Gaussian distribution with mean 100 and standard deviation 30; this number was rounded up to the nearest integer, with min. 20 and max. 300 haploids per deme. We generated a very large number of sequences for each scenario, and sampled independent single nucleotide polymorphisms (SNPs) from sequences with a single segregating site. Each scenario was replicated 50 times (500 times for Figures 3 and S2).

*Pool sequencing:* for each `ms` simulated dataset, we generated Pool-seq data by drawing reads from a binomial distribution (Gautier et al. 2013). More precisely, we assume that for each SNP, the number $r_{i:k}$ of reads of allelic type $k$ in pool $i$ follows:

$$r_{i:k} \sim \text{Bin}\left(\frac{y_{i:k}}{n_i}, \delta_i\right) \tag{14}$$

13

₂₀₀ where $y_{i:k}$ is the number of genes of type $k$ in the $i$th pool, $n_i$ is the total

₂₀₁ number of genes in pool $i$ (haploid pool size), and $\delta_i$ is the simulated total

₂₀₂ coverage for pool $i$. In the following, we either consider a fixed coverage,

₂₀₃ with $\delta_i = \Delta$ for all pools and loci, or a varying coverage across pools and

₂₀₄ loci, with $\delta_i \sim \mathrm{Pois}(\Delta)$.

₂₀₅ *Sequencing error:* we simulated sequencing errors occurring at rate $\mu_{\mathrm{e}} =$

₂₀₆ $0.001$, which is typical of Illumina sequencers (Glenn 2011; Ross et al. 2013).

₂₀₇ We assumed that each sequencing error modifies the allelic type of a read to

₂₀₈ one of three other possible states with equal probability (there are therefore

₂₀₉ four allelic types in total, corresponding to four nucleotides). Note that

₂₁₀ only biallelic markers are retained in the final datasets. Also note that,

₂₁₁ since we initiated this procedure with polymorphic markers only, we neglect

₂₁₂ sequencing errors that would create spurious SNPs from monomorphic sites.

₂₁₃ However, such SNPs should be rare in real datasets, since markers with a

₂₁₄ low minimum read count (MRC) are generally filtered out.

₂₁₅ *Experimental error:* non-equimolar amounts of DNA from all individuals in

₂₁₆ a pool and stochastic variation in the amplification efficiency of individual

₂₁₇ DNAs are sources of experimental errors in pool sequencing. To simulate

₂₁₈ experimental errors, we used the model derived by Gautier et al. (2013). In

₂₁₉ this model, it is assumed that the contribution $\eta_{ij} = c_{ij}/C_{1i}$ of each gene $j$

₂₂₀ to the total coverage of the $i$th pool ($C_{1i}$) follows a Dirichlet distribution:

$$\{\eta_{ij}\}_{1 \leq j \leq n_i} \sim \mathrm{Dir}\left(\frac{\rho}{n_i}\right) \tag{15}$$

14

221 where the parameter $\rho$ controls the dispersion of gene contributions around

222 the value $\eta_{ij} = 1/n_i$, expected if all genes contributed equally to the pool of

223 reads. For convenience, we define the experimental error $\epsilon$ as the coefficient

224 of variation of $\eta_{ij}$, i.e. $\epsilon \equiv \sqrt{\mathbb{V}(\eta_{ij})}/\mathbb{E}(\eta_{ij}) = \sqrt{(n_i - 1)/(\rho + 1)}$ (see Gautier

225 et al. 2013). When $\epsilon$ tends toward 0 (or equivalently when $\rho$ tends to infinity),

226 all individuals contribute equally to the pool, and there is no experimental

227 error. We tested the robustness of our estimates to values of $\epsilon$ comprised

228 between 0.05 and 0.5. The case $\epsilon = 0.5$ could correspond, for example, to a

229 situation where (for $n_i = 10$) 5 individuals contribute $2.8\times$ more reads than

230 the other 5 individuals.

231 **Other estimators**

232 For the sake of clarity, a summary of the notation of the $F_{\mathrm{ST}}$ estimators used

233 throughout this article is given in Table 2.

234 $\mathrm{PP2_d}$ : this estimator of $F_{\mathrm{ST}}$ is implemented by default in the software

235 package POPOOLATION2 (Kofler et al. 2011). It is based on a definition of

236 the parameter $F_{\mathrm{ST}}$ as the overall reduction in average heterozygosity relative

237 to the total combined population (see, e.g., Nei and Chesser 1983):

$$\mathrm{PP2_d} \equiv \frac{\hat{H}_{\mathrm{T}} - \hat{H}_{\mathrm{S}}}{\hat{H}_{\mathrm{T}}} \tag{16}$$

238 where $\hat{H}_{\mathrm{S}}$ is the average heterozygosity within subpopulations, and $\hat{H}_{\mathrm{T}}$ is the

239 average heterozygosity in the total population (obtained by pooling together

240 all subpopulation to form a single virtual unit). In POPOOLATION2, $\hat{H}_{\mathrm{S}}$ is

15

the unweighted average of within-subpopulation heterozygosities:

$$\hat{H}_{\mathrm{S}} = \frac{1}{n_{\mathrm{d}}} \sum_i^{n_{\mathrm{d}}} \left( \frac{n_i}{n_i - 1} \right) \left( \frac{C_{1i}}{C_{1i} - 1} \right) \left( 1 - \sum_k \hat{\pi}_{i:k}^2 \right) \tag{17}$$

(using the notation from Table 1). Note that in PoPoolation2, PP2$_{\mathrm{d}}$ is restricted to the case of two subpopulations only ($n_{\mathrm{d}} = 2$). The two ratios in the right-hand side of Equation 17 are presumably borrowed from Nei (1978) to provide an unbiased estimate, although we found no formal justification for the expression in Equation 17 for Pool-seq data. The total heterozygosity is computed as (using the notation from Table 1):

$$\hat{H}_{\mathrm{T}} = \left( \frac{\min_i(n_i)}{\min_i(n_i) - 1} \right) \left( \frac{\min_i(C_{1i})}{\min_i(C_{1i}) - 1} \right) \left( 1 - \sum_k \hat{\pi}_k^2 \right) \tag{18}$$

PP2$_{\mathrm{a}}$ : this is the alternative estimator of $F_{\mathrm{ST}}$ provided in the software package PoPoolation2. It is based on an interpretation by Kofler et al. (2011) of Karlsson et al.'s (2007) estimator of $F_{\mathrm{ST}}$, as:

$$\mathrm{PP2}_{\mathrm{a}} \equiv \frac{\hat{Q}_1^{\mathrm{r}} - \hat{Q}_2^{\mathrm{r}}}{1 - \hat{Q}_2^{\mathrm{r}}} \tag{19}$$

where $\hat{Q}_1^{\mathrm{r}}$ and $\hat{Q}_2^{\mathrm{r}}$ are the frequencies of identical pairs of reads within and between pools, respectively, computed by simple counting of IIS pairs. These are estimates of $Q_1^{\mathrm{r}}$, the IIS probability for two reads in the same pool (whether they are sequenced from the same gene or not) and $Q_2^{\mathrm{r}}$, the IIS probability for two reads in different pools. Note that the IIS probabiliy $Q_1^{\mathrm{r}}$ is different from $Q_1$ in Equation 1, which, from our definition, represents the IIS probability between distinct genes in the same pool. This approach therefore confounds pairs of reads within pools that are identical because

16

259 they were sequenced from a single gene, from pairs of reads that are identical

260 because they were sequenced from distinct, yet IIS genes.

261 FRP$_{13}$ : this estimator of $F_{ST}$ was developed by Ferretti et al. (2013) (see

262 their Equations 3 and 10–13). Ferretti et al. (2013) use the same definition of

263 $F_{ST}$ as in Equation 16 above, although they estimate heterozygosities within

264 and between pools as "average pairwise nucleotide diversities", which, from

265 their definitions, are formally equivalent to IIS probabilities. In particular,

266 they estimate the average heterozygosity within pools as (using the notation

267 from Table 1):

$$\hat{H}_{S} = \frac{1}{n_{d}} \sum_{i}^{n_{d}} \left( \frac{n_i}{n_i - 1} \right) \left( 1 - \hat{Q}_{1i}^{r} \right) \tag{20}$$

268 and the total heterozygosity among the $n_d$ populations as:

$$\hat{H}_{T} = \frac{1}{n_{d}^{2}} \left[ \sum_{i}^{n_{d}} \left( \frac{n_i}{n_i - 1} \right) \left( 1 - \hat{Q}_{1i}^{r} \right) + \sum_{i \neq i'}^{n_{d}} \left( 1 - \hat{Q}_{2ii'}^{r} \right) \right] \tag{21}$$

269 **Analyses of Ind-seq data:**

270 For the comparison of Ind-seq and Pool-seq datasets, we computed $F_{ST}$ on

271 subsamples of 5,000 loci. These subsamples were defined so that only those

272 loci that were polymorphic in all coverage conditions were retained, and the

273 same loci were used for the analysis of the corresponding Ind-seq data. For

274 the latter, we used either the Nei and Chesser's (1983) estimator based on a

275 ratio of heterozygosity (see Equation 16 above), hereafter denoted by NC$_{83}$, or

276 the analysis-of-variance estimator developed by Weir and Cockerham (1984),

277 hereafter denoted by WC$_{84}$.

278 All the estimators were computed using custom functions in the R soft-

17

²⁷⁹ ware environment for statistical computing, version 3.3.1 (R Core Team

²⁸⁰ 2017). All these functions were carefully checked against available software

²⁸¹ packages, to ensure that they provided strictly identical estimates.

### Application example: *Cottus asper*

²⁸³ Dennenmoser et al. (2017) investigated the genomic basis of adaption to

²⁸⁴ osmotic conditions in the prickly sculpin (*Cottus asper*), an abundant eury-

²⁸⁵ haline fish in northwestern North America. To do so, they sequenced the

²⁸⁶ whole-genome of pools of individuals from two estuarine populations (CR,

²⁸⁷ Capilano River Estuary; FE, Fraser River Estuary) and two freshwater pop-

²⁸⁸ ulations (PI, Pitt Lake and HZ, Hatzic Lake) in southern British Columbia

²⁸⁹ (Canada). We downloaded the four corresponding BAM files from the Dryad

²⁹⁰ Digital Repository (doi: 10.5061/dryad.2qg01) and combined them into a sin-

²⁹¹ gle mpileup file using `SAMtools` version 0.1.19 (Li et al. 2009) with default

²⁹² options, except the maximum depth per BAM that was set to 5,000 reads.

²⁹³ The resulting file was further processed using a custom `awk` script, to call

²⁹⁴ SNPs and compute read counts, after discarding bases with a Base Align-

²⁹⁵ ment Quality (BAQ) score lower than 25. A position was then considered

²⁹⁶ as a SNP if: ($i$) only two different nucleotides with a read count $> 1$ were

²⁹⁷ observed (nucleotides with $\leq 1$ read being considered as a sequencing error);

²⁹⁸ ($ii$) the coverage was comprised between 10 and 300 in each of the four align-

²⁹⁹ ment files; ($iii$) the minor allele frequency, as computed from read counts,

³⁰⁰ was $\geq 0.01$ in the four populations. The final data set consisted of 608,879

³⁰¹ SNPs.

³⁰²     Our aim here was to compare the population structure inferred from pair-

³⁰³ wise estimates of $F_{ST}$, using the estimator $\hat{F}_{ST}^{pool}$ on the one hand (Equa-

18

304 tion 12), and PP2$_\text{d}$ on the other hand. Then, to conclude on which of the

305 two estimators performs better, we compared the population structure in-

306 ferred from $\hat{F}_\text{ST}^\text{pool}$ and PP2$_\text{d}$ to that inferred from the Bayesian hierarchical

307 model implemented in the software package BayPass (Gautier 2015). Bay-

308 Pass allows the robust estimation of the scaled covariance matrix of allele

309 frequencies across populations for Pool-seq data, which is known to be infor-

310 mative about population history (Pickrell and Pritchard 2012). The elements

311 of the estimated matrix can be interpreted as pairwise and population-specific

312 estimates of differentiation (Coop et al. 2010), and therefore provide a com-

313 prehensive description of population structure that makes full use of the

314 available data.

## Data availability

316 The authors state that all data necessary for confirming the conclusions

317 presented in this article are fully represented within the article, figures,

318 and tables. Supplemental Tables S1–S3 and Figures S1–S4 are available at

319 FigShare, along with a complete derivation of the model in the Supplemental

320 File S1 at FigShare.

19

<sub>321</sub> RESULTS

## Comparing Ind-seq and Pool-seq estimates of $F_{\text{ST}}$

<sub>323</sub> Single-locus estimates $\hat{F}_{\text{ST}}^{\text{pool}}$ are highly correlated with the classical estimates

<sub>324</sub> WC$_{84}$ (Weir and Cockerham 1984) computed on the individual data that were

<sub>325</sub> used to generate the pools in our simulations (see Figure 1). The variance of

<sub>326</sub> $\hat{F}_{\text{ST}}^{\text{pool}}$ across independent replicates decreases as the coverage increases. The

<sub>327</sub> correlation between $\hat{F}_{\text{ST}}^{\text{pool}}$ and WC$_{84}$ is stronger for multilocus estimates (see

<sub>328</sub> Figure S1A).

## Comparing Pool-seq estimators of $F_{\text{ST}}$

<sub>330</sub> We found that our estimator $\hat{F}_{\text{ST}}^{\text{pool}}$ has extremely low bias ($< 0.5\%$ over

<sub>331</sub> all scenarios tested: see Tables 3 and S1-S3). In other words, the average

<sub>332</sub> estimates across multiple loci and replicates closely equal the expected value

<sub>333</sub> of the $F_{\text{ST}}$ parameter, as given by Equation 6 in Rousset (1996), which is

<sub>334</sub> based on the computation of IIS probabilities in an island model of population

<sub>335</sub> structure. In all the situations examined, the bias does neither depend on

<sub>336</sub> the sample size (i.e., the size of each pool) nor on the coverage (see Figure 2).

<sub>337</sub> Only the variance of the estimator across independent replicates decreases as

<sub>338</sub> the sample size increases and/or as the coverage increases. At high coverage,

<sub>339</sub> the mean and root mean squared error (RMSE) of $\hat{F}_{\text{ST}}^{\text{pool}}$ over independent

<sub>340</sub> replicates are virtually indistinguishable from that of the WC$_{84}$ estimator

<sub>341</sub> (see Table S1).

<sub>342</sub>   Figure 3 shows the RMSE of $F_{\text{ST}}$ estimates for a wide range of pool sizes

<sub>343</sub> and coverages. The RMSE decreases as the pool size and/or the coverage

<sub>344</sub> increases. The $F_{\text{ST}}$ estimates are more precise and accurate when differen-

345 tiation is low. Figure 3 provides some clues to evaluate the pool size and

346 the coverage that is necessary to achieve the same RMSE than for Ind-seq

347 data. Consider, for example, the case of samples of $n = 20$ haploids. For

348 $F_{ST} \leq 0.05$ (in the conditions of our simulations), the RMSE of $F_{ST}$ estimates

349 based on Pool-seq data tends to the RMSE of $F_{ST}$ estimates based on Ind-seq

350 data either by sequencing pools of ca. 200 haploids at 20X, or by sequencing

351 pools of 20 haploids at ca. 200X. However, the same precision and accuracy

352 are achieved by sequencing ca. 50 haploids at ca. 50X.

353      Conversely, we found that $PP2_d$ (the default estimator of $F_{ST}$ imple-

354 mented in the software package PoPoolation2) is biased when compared

355 to the expected value of the parameter. We observed that the bias depends

356 on both the sample size, and the coverage (see Figure 2). We note that, as the

357 coverage and the sample size increase, $PP2_d$ converges to the estimator $NC_{83}$

358 (Nei and Chesser 1983) computed from individual data (see Figure S1B).

359 This argument was used by Kofler et al. (2011) to validate the approach,

360 even though the estimates $PP2_d$ depart from the true value of the parameter

361 (Figure S1B–C).

362      The second of the two estimators of $F_{ST}$ implemented in PoPoolation2,

363 that we refer to as $PP2_a$, is also biased (see Figure 2). We note that the bias

364 decreases as the sample size increases. However, the bias does not depend

365 on the coverage (only the variance over independent replicates does). The

366 estimator developed by Ferretti et al. (2013), that we refer to as $FRP_{13}$, is

367 also biased (see Figure 2). However, the bias does neither depend on the pool

368 size, nor on the coverage (only the variance over independent replicates does).

369 $FRP_{13}$ converges to the estimator $NC_{83}$, computed from individual data (see

21

370 Figure 2). At high coverage, the mean and RMSE over independent replicates

371 are virtually indistinguishable from that of the $NC_{83}$ estimator.

372    Last, we stress out that our estimator $\hat{F}_{ST}^{pool}$ provides estimates for multiple

373 populations, and is therefore not restricted to pairwise analyses, contrary to

374 PoPoolation2's estimators. We show that, even at low sample size and low

375 coverage, Pool-seq estimates of differentiation are virtually indistinguishable

376 from classical estimates for Ind-seq data (see Table 3).

### 377 Robustness to unbalanced pool sizes and variable sequencing cov-

### 378 erage

379 We evaluated the accuracy and the precision of the estimator $\hat{F}_{ST}^{pool}$ when sam-

380 ple sizes differ across pools, and when the coverage varies across pools and loci

381 (see Figure 4). We found that, at low coverage, unequal sampling or variable

382 coverage causes a negligible departure from the median of $WC_{84}$ estimates

383 computed on individual data, which vanishes as the coverage increases. At

384 100X coverage, the distribution of $\hat{F}_{ST}^{pool}$ estimates is almost indistinguishable

385 from that of $WC_{84}$ (see Figure 4 and Tables S2–S3).

### 386 Robustness to sequencing and experimental errors

387 Figure 5 shows that sequencing errors cause a negligible negative bias for

388 $\hat{F}_{ST}^{pool}$ estimates. Filtering (using a minimum read count of 4) improves es-

389 timation slightly, but only at high coverage (Figure 6B). It must be noted,

390 though, that filtering increases the bias in the absence of sequencing error,

391 especially at low coverage (Figure 6A). With experimental error, i.e., when

392 individuals do not contribute evenly to the final set of reads, we observed a

393 positive bias for $\hat{F}_{ST}^{pool}$ estimates (Figure 5). We note that the bias decreases

394 as the size of the pools increases. Figure S2 shows the RMSE of $F_{\mathrm{ST}}$ esti-

395 mates for a wider range of pool sizes, coverage and experimental error rate

396 ($\epsilon$). For $\epsilon \geq 0.25$, increasing the coverage cannot improve the quality of the

397 inference, if the pool size is too small. When Pool-seq experiments are prone

398 to large experimental error rates, increasing the size of pools is the only way

399 to improve the estimation of $F_{\mathrm{ST}}$. Filtering (using a minimum read count of

400 4) does not improve estimation (Figure 6C).

401 **Application example**

402 The reanalysis of the prickly sculpin data revealed larger pairwise estimates of

403 multilocus $F_{\mathrm{ST}}$ using PP2$_{\mathrm{d}}$ estimator, as compared to $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ (see Figure 7A).

404 Furthermore, we found that $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ estimates are smaller for within-ecotype

405 pairwise comparisons as compared to between-ecotype comparisons. There-

406 fore, the inferred relationships between samples based on pairwise $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ esti-

407 mates show a clear-cut structure, separating the two estuarine samples from

408 the freshwater ones (see Figure 7C). We did not recover the same structure

409 using PP2$_{\mathrm{d}}$ estimates (see Figure 7B). Supportingly, the scaled covariance

410 matrix of allele frequencies across samples is consistent with the structure

411 inferred from $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ estimates (see Figure 7D).

DISCUSSION

Whole-genome sequencing of pools of individuals is increasingly popular for population genomic research on both model and non-model species (Schlötterer et al. 2014). The development of dedicated software packages (reviewed in Schlötterer et al. 2014) has undoubtedly something to do with the breadth of research questions that have been tackled using pool-sequencing. Yet, the analysis of population structure from Pool-seq data is complicated by the double sampling process of genes from the pool and sequence reads from those genes (Ferretti et al. 2013).

The naive approach that consists in computing $F_{ST}$ from read counts, as if they were allele counts (e.g., as in Chen et al. 2016), ignores the extra variance brought by the random sampling of reads from the gene pool during Pool-seq experiments. Furthermore, such computation fails to consider the actual number of lineages in the pool (haploid pool size). Altogether, these limits may result in severely biased estimates of differentiation when the pool size is low (see Figure S3). A possible alternative is to compute $F_{ST}$ from allele counts imputed from read counts using a maximum-likelihood approach conditional on the haploid size of the pools (e.g., as in Smadja et al. 2012; Leblois et al. 2018), or from allele frequencies estimated using a model-based method that accounts for the sampling effects and the sequencing error probabilities inherent to pooled NGS experiments (see Fariello et al. 2017). However, these latter approaches may only be accurate in situations where the coverage is much larger than pool size, allowing to reduce sampling variance of reads (see Figure S3). Here, we therefore developed a new estimator of the parameter $F_{ST}$ for Pool-seq data, in an analysis-of-variance

24

437 framework (Cockerham 1969, 1973). The accuracy of this estimator is barely

438 distinguishable from that of the Weir and Cockerham's (1984) estimator for

439 individual data. Furthermore, it does neither depend on the pool size nor on

440 the coverage, and is robust to unequal pool sizes and varying coverage across

441 demes and loci.

442     In our analysis, the frequency of reads within pools is a weighted av-

443 erage of the sample frequencies, with weights equal to the pool coverage.

444 Therefore, our approach follows Cockerham's (1973) one, which he referred

445 to as a weighted analysis-of-variance (see also Weir and Cockerham 1984;

446 Weir 1996; Weir and Hill 2002; Weir and Goudet 2017). With unequal pool

447 sizes, weighted and unweighted analyses differ. As discussed recently in Weir

448 and Goudet (2017), the unweighted approach seems appropriate when the

449 between component exceeds the within component, i.e. when $F_{ST}$ is large

450 (Tukey 1957). It turns out that optimal weighting depends upon the param-

451 eter to be estimated (Cockerham 1973) and is only efficient at lower levels of

452 differentiation (Robertson 1962). In a likelihood analysis of the island model,

453 Rousset (2007) derived asymptotically efficient weights that are proportional

454 to $n_i^2$ for the sum of squares of different samples (see also Robertson 1962). To

455 the best of our knowledge, such optimal weighting has never been considered

456 in the literature.

## Analysis of variance and probabilities of identity

458 In the analysis-of-variance framework, $F_{ST}$ is defined in Equation 1 as an

459 intraclass correlation for the probability of identity in state (Cockerham and

460 Weir 1987; Rousset 1996). Extensive statistical literature is available on

461 estimators of intraclass correlations. Beside analysis-of-variance estimators,

introduced in population genetics by Cockerham (1969, 1973), estimators based on the computation of probabilities of identical response within and between groups have been proposed (see, e.g., Fleiss 1971; Fleiss and Cuzick 1979; Mak 1988; Ridout et al. 1999; Wu et al. 2012), which were originally referred to as kappa-type statistics (Fleiss 1971; Landis and Koch 1977). These estimators have later been endorsed in population genetics, where the "probability of identical response" was then interpreted as the frequency with which the genes are alike (Cockerham 1973; Cockerham and Weir 1987; Weir 1996; Rousset 2007; Weir and Goudet 2017).

This suggests that, with Pool-seq data, another strategy could consist in computing $F_{\mathrm{ST}}$ from IIS probabilities between (unobserved) pairs of genes, which requires that unbiased estimates of such quantities are derived from read count data. We have done so in the second section of the Supplemental File S1, and we provide alternative estimators of $F_{\mathrm{ST}}$ for Pool-seq data (see Equations A44 and A48 in Supplemental File S1). These estimators (denoted by $\hat{F}_{\mathrm{ST}}^{\mathrm{pool-PID}}$ and $\tilde{F}_{\mathrm{ST}}^{\mathrm{pool-PID}}$) have exactly the same form as the analysis-of-variance estimator if the pools have all the same size and if the number of reads per pool is constant (Equation A33). This echoes the derivations by Rousset (2007) for Ind-seq data, who showed that the analysis-of-variance approach (Weir and Cockerham 1984) and the simple strategy of estimating IIS probabilities by counting identical pairs of genes provide identical estimates when sample sizes are equal (see Equation A28 and also Cockerham and Weir 1987; Weir 1996; Karlsson et al. 2007). With unbalanced samples, we found that analysis-of-variance estimates have better precision and accuracy than IIS-based estimates, particularly for low levels of differ-

26

487  entiation (see Figure S4). Interestingly, we found that IIS-based estimates

488  of $F_{\text{ST}}$ for Pool-seq data have generally lower bias and variance if the over-

489  all estimates of IIS probabilities within and between pools are computed as

490  unweighted averages of population-specific or pairwise estimates (see Equa-

491  tions A39 and A43), as compared to weighted averages (Equations A46–A47).

492  Equation A28 further shows that our estimator may be rewritten as a func-

493  tion close to $\left(\hat{Q}_1 - \hat{Q}_2\right) / \left(1 - \hat{Q}_2\right)$, except that it also depends on the sum

494  $\sum_i \left(\hat{Q}_{1i} - \hat{Q}_1\right)$ in both the numerator and the denominator. This suggests

495  that if the $Q_{1i}$'s differ among subpopulations, then our estimator provides an

496  estimate of an average of population-specific $F_{\text{ST}}$ (Weir and Hill 2002; Weir

497  and Goudet 2017).

498  It follows from the derivations in the Supplemental File S1 that the es-

499  timator $\text{PP2}_{\text{a}}$ (Equation 19) is biased because the IIS probability between

500  pairs of reads within a pool $\left(\hat{Q}_1^{\text{r}}\right)$ is a biased estimator of the IIS probability

501  between pairs of distinct genes in that pool (see Equations A34–A36 in Sup-

502  plemental File S1). This is so, because the former confounds pairs of reads

503  that are identical because they were sequenced from a single gene, from pairs

504  of reads that are identical because they were sequenced from distinct, yet IIS

505  genes.

506  A more justified estimator of $F_{\text{ST}}$ has been proposed by Ferretti et al.

507  (2013), based on previous developments by Futschik and Schlötterer (2010).

508  Note that, although they defined $F_{\text{ST}}$ as a ratio of functions of heterozygosi-

509  ties, they actually worked with IIS probabilities (see Equations 20 and 21).

510  However, although Equation 20 is strictly identical to Equation A39 in Sup-

511  plemental File S1, we note that they computed the total heterozygosity by

27

integrating over pairs of genes sampled both within and between subpopulations (compare Equation 21 with A43), which may explain the observed bias (see Figure 2).

**Comparison with alternative estimators**

An alternative framework to Weir and Cockerham's (1984) analysis-of-variance has been developed by Masatoshi Nei and coworkers to estimate $F_{ST}$ from gene diversities (Nei 1973, 1977; Nei and Chesser 1983; Nei 1986). The estimator $PP2_d$ (see Equations 16–18) implemented in the software package PoPoolation2 (Kofler et al. 2011) follows this logic. However, it has long been recognized that both frameworks are fundamentally different in that the analysis-of-variance approach considers both statistical and genetic (or evolutionary) sampling, whereas Nei and coworkers' approach do not (Weir and Cockerham 1984; Excoffier 2007; Holsinger and Weir 2009). Furthermore, the expectation of Nei and coworkers' estimators depend upon the number of sampled populations, with a larger bias for lower numbers of sampled populations (Goudet 1993; Excoffier 2007; Weir and Goudet 2017). This is so, because the computation of the total diversity in Equations 18 and 21 includes the comparison of pairs of genes from the same subpopulation, whereas the computation of IIS probabilities between subpopulations do not (see, e.g., Excoffier 2007). Therefore, we do not recommend using the estimator $PP2_d$ implemented in the software package PoPoolation2 (Kofler et al. 2011).

**Applications in evolutionary ecology studies**

Pool-seq is being increasingly used in many application domains (Schlötterer et al. 2014), such as conservation genetics (see, e.g., Fuentes-Pardo and

28

Ruzzente 2017), invasion biology (see, e.g., Dexter et al. 2018) and evolutionary biology in a broader sense (see, e.g., Collet et al. 2016). These studies use a large range of methods, which aim at characterizing fine-scaled population structure (see, e.g., Fischer et al. 2017), reconstructing past demography (see, e.g., Chen et al. 2016; Leblois et al. 2018), or identifying footprints of natural or artificial selection (see, e.g., Chen et al. 2016; Fariello et al. 2017; Leblois et al. 2018).

Here, we reanalyzed the Pool-seq data produced by Dennenmoser et al. (2017), who investigated the adaptive genomic divergence between freshwater and brackish-water ecotypes of the prickly sculpin *C. asper*, an abundant euryhaline fish in northwestern North America. Measuring pairwise genetic differentiation between samples using $\hat{F}_{\text{ST}}^{\text{pool}}$, we found a clear-cut structure separating the freshwater from the brackish-water ecotypes. Such genetic strucure supports the hypothesis that populations are locally adapted to osmotic conditions in these two contrasted habitats, as discussed in Dennenmoser et al. (2017). This structure, which is at odds with that inferred from PP2$_{\text{d}}$ estimates, is not only supported by the scaled covariance matrix of allele frequencies, but also by previous microsatellite-based studies, who showed that populations were genetically more differentiated between ecotypes than within ecotypes (Dennenmoser et al. 2014, 2015).

## Limits of the model and perspectives

We have shown that the stronger source of bias for the $\hat{F}_{\text{ST}}^{\text{pool}}$ estimate is unequal contributions of individuals in pools. This is so, because we assume in our model that the read counts are multinomially distributed, which supposes that all genes contribute equally to the pool of reads (Gautier et al. 2013),

29

i.e. that there is no variation in DNA yield across individuals and that all genes have equal sequencing coverage (Rode et al. 2018). Because the effect of unequal contribution is expected to be stronger with small pool sizes, it has been recommended to use pool-seq with at least 50 diploid individuals per pool (Lynch et al. 2014; Schlötterer et al. 2014). However, this limit may be overly conservative for allele frequency estimates (Rode et al. 2018), and we have shown here that we can achieve very good precision and accuracy of $F_{ST}$ estimates with smaller pool sizes. Furthermore, because genotypic information is lost during Pool-seq experiments, we assume in our derivations that pools are haploid (and therefore that $F_{IS}$ is nil). Analyzing non-random mating populations (e.g., in selfing species) is therefore problematic.

Finally, our model, as in Weir and Cockerham (1984), formally assumes that all populations provide independent replicates of some evolutionary process (Excoffier 2007; Holsinger and Weir 2009). This may be unrealistic in many natural populations, which motivated Weir and Hill (2002) to derive a population-specific estimator of $F_{ST}$ for Ind-seq data (see also Vitalis et al. 2001). Even though the use of Weir and Hill's (2002) estimator is still scarce in the literature (but see Weir et al. 2005; Vitalis 2012), Weir and Goudet (2017) recently proposed a re-interpretation of population-specific estimates of $F_{ST}$ in terms of allelic matching proportions, which are strictly equivalent to IIS probabilities between pairs of genes. It would therefore be straightforward to extend Weir and Goudet's (2017) estimator of population-specific $F_{ST}$ for the analysis of Pool-seq data, using the unbiased estimates of IIS probabilies provided in the Supplemental File S1.

## DATA ACCESSIBILITY

A R package, called `poolfstat`, which impletements $F_{ST}$ estimates for Pool-seq data, is available at the Comprehensive R Archive Network (CRAN): `https://cran.r-project.org/web/packages/poolfstat/index.html`.

## ACKNOWLEDGEMENTS

31

## LITERATURE CITED

Akey, J. M., Zhang, G., Jin, L., and Shriver, M. D. (2002). Interrogating a high-density SNP map for signatures of natural selection. *Genome Res.*, 12:1805–1814.

Anderson, E. C., Skaug, H. J., and Barshis, D. J. (2014). Next-generation sequencing for molecular ecology: a caveat regarding pooled samples. *Mol. Ecol.*, 23:502–512.

Beaumont, M. A. (2005). Adaptation and speciation: what can $F_{ST}$ tell us? *Trends Ecol. Evol.*, 20:435–440.

Beaumont, M. A. and Nichols, R. A. (1996). Evaluating loci for use in the genetic analysis of population structure. *Proc. R. Soc. Lond. B*, 263:1619–1626.

Bhatia, G., Patterson, N., Sankararaman, S., and Price, A. L. (2013). Estimating and interpreting $F_{ST}$: the impact of rare variants. *Genome Res.*, 23:1514–1521.

Cavalli-Sforza, L. (1966). Population structure and human evolution. *Proc. R. Soc. Lond., B, Biol. Sci.*, 164:362–379.

Chen, J., Källman, T., Ma, X.-F., Zaina, G., Morgante, M., and Lascoux, M. (2016). Identifying genetic signatures of natural selection using pooled populations sequencing in *Picea abies*. *G3*, 6:1979–1989.

Cockerham, C. C. (1969). Variance of gene frequencies. *Evolution*, 23:72–84.

Cockerham, C. C. (1973). Analyses of gene frequencies. *Genetics*, 74:679–700.

622  Cockerham, C. C. and Weir, B. S. (1987). Analyses of gene frequencies. *Proc.*

623  *Natl. Acad. Sci. USA*, 84:8512–8514.

624  Collet, J. M., Fuentes, S., Hesketh, J., Hill, M. S., Innocenti, P., Morrow,

625  E. H., Fowler, K., and Reuter, M. (2016). Rapid evolution of the intersex-

626  ual genetic correlation for fitness in *Drosophila melanogaster*. *Evolution*,

627  70:781–795.

628  Coop, G., Witonsky, D., Di Rienzo, A., and Pritchard, J. K. (2010). Us-

629  ing environmental correlations to identify loci underlying local adaptation.

630  *Genetics*, 185:1411–1423.

631  Cutler, D. J. and Jensen, J. D. (2010). To pool, or not to pool? *Genetics*,

632  186:41–43.

633  Dennenmoser, S., Nolte, A. W., Vamosi, S. M., and Rogers S, M. (2015). Phy-

634  logeography of the prickly sculpin (*Cottus asper*) in north-western North

635  America reveals parallel phenotypic evolution across multiple coastal-

636  inland colonizations. *J. Biogeogr.*, 42:1626–1638.

637  Dennenmoser, S., Rogers, S. M., and Vamosi, S. M. (2014). Genetic pop-

638  ulation structure in prickly sculpin (*Cottus asper*) reflects isolation-by-

639  environment between two life-history ecotypes. *Biol. J. Linnean Soc.*,

640  113:943–957.

641  Dennenmoser, S., Vamosi, S. M., Nolte, S. W., and Rogers, S. M. (2017).

642  Adaptive genomic divergence under high gene flow between freshwater and

643  brackish-water ecotypes of prickly sculpin (*Cottus asper*) revealed by Pool-

644  Seq. *Mol. Ecol.*, 26:25–42.

33

Dexter, E., Bollens, S. M., Cordell, J., Soh, H. Y., Rollwagen-Bollens, G., Pfeifer, S. P., Goudet, J., and Vuilleumier, S. (2018). A genetic reconstruction of the invasion of the calanoid copepod *Pseudodiaptomus inopinus* across the North American Pacific Coast. *Biol. Invasions*, 20:1577–1595.

Ellegren, H. (2014). Genome sequencing and population genomics in non-model organisms. *Trends Ecol. Evol.*, 29:51–63.

Excoffier, L. (2007). Analysis of population subdivision. In Balding, D. J., Bishop, M., and Cannings, C., editors, *Handbook of Statistical Genetics*, pages 980–1020, Chichester. John Wiley & Sons, Ltd.

Fariello, M. I., Boitard, S., Mercier, S., Robelin, D., Faraut, T., Arnould, C., Recoquillay, J., Bouchez, O., Salin, G., Dehais, P., Gourichon, D., Leroux, S., Pitel, F., Leterrier, C., and SanCristobal, M. (2017). Accounting for Linkage Disequilibrium in genome scans for selection without individual genotypes : the local score approach. *Mol. Ecol.*, 26:3700–3714.

Ferretti, L., Ramos Onsins, S., and Pérez-Enciso, M. (2013). Population genomics from pool sequencing. *Mol. Ecol.*, 22:5561–5576.

Fischer, M. C., Rellstab, C., Leuzinger, M., Roumet, M., Gugerli, F., Shimizu, K. K., Holderegger, R., and Widmer, A. (2017). Estimating genomic diversity and population differentiation – an empirical comparison of microsatellite and SNP variation in *Arabidopsis halleri*. *BMC Genomics*, 18:69.

Fleiss, J. L. (1971). Measuring nominal scale agreement among many raters. *Psychol. Bull.*, 76:378–382.

34

668  Fleiss, J. L. and Cuzick, J. (1979). The reliability of dichotomous judgements:
669      Unequal numbers of judges per subject. *Appl. Psychol. Meas.*, 3:537–542.

670  Fuentes-Pardo, A. P. and Ruzzente, D. E. (2017). Whole-genome sequencing
671      approaches for conservation biology: Advantages, limitations and practical
672      recommendations. *Mol. Ecol.*, 26:5369–5406.

673  Futschik, A. and Schlötterer, C. (2010). The next generation of molecu-
674      lar markers from massively parallel sequencing of pooled DNA samples.
675      *Genetics*, 186:207–218.

676  Gautier, M. (2015). Genome-wide scan for adaptive divergence and associa-
677      tion with population-specific covariates. *Genetics*, 201:1555–1579.

678  Gautier, M., Gharbi, K., Cezaerd, T., Galan, M., Loiseau, A., Thomson, M.,
679      Pudlo, P., Kerdelhué, C., and Estoup, A. (2013). Estimation of popula-
680      tion allele frequencies from next-generation sequencing data: pool-versus
681      individual-based genotyping. *Mol. Ecol.*, 22:3766–3779.

682  Glenn, T. C. (2011). Field guide to next-generation DNA sequencers. *Mol.*
683      *Ecol. Resour.*, 11:759–769.

684  Goudet, J. (1993). *The genetics of geographically structured populations.* PhD
685      thesis, University of Wales, Bangor.

686  Holsinger, K. S. and Weir, B. S. (2009). Genetics in geographically structured
687      populations: defining, estimating and interpreting $F_{ST}$. *Nat. Rev. Genet.*,
688      10:639–650.

689  Hudson, R. R. (2002). Generating samples under a Wright-Fisher neutral
690      model of genetic variation. *Bioinformatics*, 18:337–338.

35

Karlsson, E. K., Baranowska, I., Wade, C. M., Salmon Hillbertz, N. H. C., Zody, M. C., Anderson, N., Biagi, T. M., Patterson, N., Pielberg, G. R., Kulbokas, E. J., Comstock, K. E., Keller, E. T., Mesirov, J. P., von Euler, H., Kämpe, O., Hedhammar, A., Lander, E. S., Andersson, G., Andersson, L., and Lindblad-Toh, K. (2007). Efficient mapping of Mendelian traits in dogs through genome-wide association. *Nat. Genet.*, 39:1321–1328.

Kofler, R., Pandey, R. V., and Schlötterer, C. (2011). PoPoolation2: identifying differentiation between populations using sequencing of pooled DNA samples (Pool-Seq). *Bioinformatics*, 27:3435–3436.

Landis, J. R. and Koch, G. G. (1977). A one-way components of variance model for categorical data. *Biometrics*, 33:671–679.

Leblois, R., Gautier, M., Rohfritsch, A., Foucaud, J., Burban, C., Galan, M., Loiseau, A., Sauné, L., Branco, M., Gharbi, K., Vitalis, R., and Kerdelhué, C. (2018). Deciphering the demographic history of allochronic differentiation in the pine processionary moth *Thaumetopoea pityocampa*. *Mol. Ecol.*, 27:264–278.

Lewontin, R. C. and Krakauer, J. (1973). Distribution of gene frequency as a test of the theory of the selective neutrality of polymorphism. *Genetics*, 74:175–195.

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25:2078–2079.

Lotterhos, K. E. and Whitlock, M. C. (2014). Evaluation of demographic history and neutral parameterization on the performance of $F_{\mathrm{ST}}$ outlier tests. *Mol. Ecol.*, 23:2178–2192.

Lotterhos, K. E. and Whitlock, M. C. (2015). The relative power of genome scans to detect local adaptation depends on sampling design and statistical method. *Mol. Ecol.*, 24:1031–1046.

Lynch, M., Bost, D., Wilson, S., Maruki, T., and Harrison, S. (2014). Population-genetic inference from pooled-sequencing data. *Genome Biol. Evol.*, 6:1210–1218.

Mak, T. K. (1988). Analysing intraclass correlation for dichotomous variables. *J. R. Stat. Soc. Ser. C Appl. Stat.*, 37:344–352.

Malécot, G. (1948). *Les Mathématiques de l'Hérédité*. Masson, Paris.

Nei, M. (1973). Analysis of gene diversity in subdivided populations. *Proc. Natl. Acad. Sci. USA*, 70:3321–3323.

Nei, M. (1977). *F*-statistics and analysis of gene diversity in subdivided populations. *Ann. Hum. Genet.*, 41:225–233.

Nei, M. (1978). Estimation of average heterozygosity and genetic distance from a small number of individuals. *Genetics*, 89:583–590.

Nei, M. (1986). Definition and estimation of fixation indices. *Evolution*, 40:643–645.

Nei, M. and Chesser, R. K. (1983). Estimation of fixation indices and gene diversities. *Ann. Hum. Genet.*, 47:253–259.

736    Nychka, D., Furrer, R., Paige, J., and Sain, S. (2017). fields: Tools for spatial
737        data. R package version 9.6.

738    Orgogozo, V., Peluffo, A. E., and Morizot, B. (2016). The "mendelian gene"
739        and the "molecular gene": two relevant concepts of genetic units. In Or-
740        gogozo, V., editor, *Genes and Evolution*, volume 119 of *Current Topics in*
741        *Developmental Biology*, pages 1–26. Academic Press.

742    Pickrell, J. K. and Pritchard, J. K. (2012). Inference of population splits
743        and mixtures from genome-wide allele frequency data. *PLoS Genet.*,
744        8(11):e1002967.

745    R Core Team (2017). *R: A Language and Environment for Statistical Com-*
746        *puting*. R Foundation for Statistical Computing, Vienna, Austria.

747    Reynolds, J., Weir, B. S., and Cockerham, C. C. (1983). Estimation of the
748        coancestry coefficient: basis for a short-term genetic distance. *Genetics*,
749        105:767–779.

750    Ridout, M. S., Demktrio, C. G. B., and Firth, D. (1999). Estimating intra-
751        class correlation for binary data. *Biometrics*, 55:137–148.

752    Robertson, A. (1962). Weighting in the estimation of variance components
753        in the unbalanced single classification. *Biometrics*, 18:413–417.

754    Rode, N. O., Holtz, Y., Loridon, K., Santoni, S., Ronfort, J., and Gay, J.
755        (2018). How to optimize the precision of allele and haplotype frequency
756        estimates using pooled-sequencing data. *Mol. Ecol. Resour.*, 18:194–203.

757    Ross, M. G., Russ, C., Costello, M., Hollinger, A., Lennon, N. J., Hegarty,

38

R., Nusbaum, C., and Jaffe, D. B. (2013). Characterizing and measuring bias in sequence data. *Genome Biol.*, 14:R51.

Rousset, F. (1996). Equilibrium values of measures of population subdivision for stepwise mutation processes. *Genetics*, 142:1357–1362.

Rousset, F. (1997). Genetic differentiation and estimation of gene flow from *F*-statistics under isolation by distance. *Genetics*, 145:1219–1228.

Rousset, F. (2007). Inferences from spatial population genetics. In Balding, D. J., Bishop, M., and Cannings, C., editors, *Handbook of Statistical Genetics*, pages 945–979, Chichester. John Wiley & Sons, Ltd.

Rousset, F. (2008). genepop'007: a complete re-implementation of the genepop software for Windows and Linux. *Mol. Ecol. Resour.*, 8:103–106.

Schlötterer, C., Tobler, R., Kofler, R., and Nolte, V. (2014). Sequencing pools of individuals – mining genome-wide polymorphism data without big funding. *Nat. Rev. Genet.*, 15:749–763.

Slatkin, M. (1993). Isolation by distance in equilibrium and non-equilibrium populations. *Evolution*, 47:264–279.

Smadja, C. M., Canbäck, B., Vitalis, R., Gautier, M., Ferrari, J., Zhou, J.-J., and Butlin, R. K. (2012). Large-scale candidate gene scan reveals the role of chemoreceptor genes in host plant specialization and speciation in the pea aphid. *Evolution*, 66:2723–2738.

The International HapMap Consortium (2005). A haplotype map of the human genome. *Nature*, 437:1299–1320.

39

Tukey, J. W. (1957). Variances of variance components: II. The unbalanced single classification. *Ann. Math. Statist.*, 28:43–56.

Vitalis, R. (2012). DetSel: An R-Package to detect marker loci responding to selection. In Pompanon, F. and Bonin, A., editors, *Data Production and Analysis in Population Genomics: Methods and Protocols*, volume 888 of *Methods in Molecular Biology*, pages 277–293, New York. Humana Press.

Vitalis, R., Boursot, P., and Dawson, K. (2001). Interpretation of variation across marker loci as evidence of selection. *Genetics*, 158:1811–1823.

Wahlund, S. (1928). Zusammens etzung von populationen und korrelationserscheinungen vom standpunkt der vererbungslehre aus betrachtet. *Hereditas*, 11:65–106.

Weir, B. S. (1996). *Genetic Data Analysis II.* Sinauer Associates, Inc., Sunderland, MA.

Weir, B. S. (2012). Estimating $F$-statistics: A historical view. *Philos. Sci.*, 79:637–643.

Weir, B. S., Cardon, L. R., Anderson, A. D., Nielsen, D. M., and Hill, W. G. (2005). Measures of human population structure show heterogeneity among genomic regions. *Genome Res.*, 15:1468–1476.

Weir, B. S. and Cockerham, C. C. (1984). Estimating $F$-statistics for the analysis of population structure. *Evolution*, 38:1358–1370.

Weir, B. S. and Goudet, J. (2017). An unified characterization of population structure and relatedness. *Genetics*, 206:2085–2103.

40

Weir, B. S. and Hill, W. G. (2002). Estimating $F$-statistics. *Annu. Rev. Genet.*, 36:721–750.

Whitlock, M. C. and Lotterhos, K. E. (2015). Reliable detection of loci responsible for local adaptation: inference of a null model through trimming the distribution of $F_{\text{ST}}$. *Am. Nat.*, 186:S24–S36.

Wright, S. (1931). Evolution in Mendelian populations. *Genetics*, 16:97–159.

Wright, S. (1951). The genetical structure of populations. *Ann. Eugen.*, 15:323–354.

Wu, S., Crespi, C. M., and Wong, W. K. (2012). Comparison of methods for estimating the intraclass correlation coefficient for binary responses in cancer prevention cluster randomized trials. *Contemp. Clin. Trials*, 33:869–880.

## Table 1  Summary of main notations

| Notation | Parameter definition |
| --- | --- |
| $X_{ijr:k}$ | Indicator variable: $X_{ijr:k} = 1$ if the $r$th read from the $j$th individual in the $i$th pool is of type $k$, and $X_{ijr:k} = 0$ otherwise |
| $r_{i:k} = \sum_j \sum_r X_{ijr:k}$ | Number of reads of type $k$ in the $i$th pool |
| $c_{ij}$ | Number of reads sequenced from individual $j$ in sub-population $i$ (unobserved individual coverage) |
| $C_{1i} \equiv \sum_j c_{ij}$ | Total number of reads in the $i$th pool (pool coverage) |
| $C_1 \equiv \sum_i C_{1i}$ | Total number of reads in the full sample (total coverage) |
| $C_2 \equiv \sum_i C_{1i}^2$ | Squared number of reads in the full sample |
| $n_i$ | Total number of genes the $i$th pool (haploid pool size) |
| $y_{i:k}$ | (Unobserved) number of genes of type $k$ in the $i$th pool |
| $\pi_k \equiv \mathbb{E}(X_{ijr:k})$ | Expected frequency of reads of type $k$ in the full sample |
| $\hat{\pi}_{ij:k} \equiv X_{ij\cdot:k}$ | (Unobserved) average frequency of reads of type $k$ for individual $j$ in the $i$th pool |
| $\hat{\pi}_{i:k} \equiv X_{i\cdot\cdot:k}$ | Average frequency of reads of type $k$ in the $i$th pool |
| $\hat{\pi}_k \equiv X_{\cdots:k}$ | Average frequency of reads of type $k$ in the full sample |
| $Q_1$ (resp. $Q_2$) | IIS probability for two genes sampled within (resp. between) pools |
| $Q_1^r$ (resp. $Q_2^r$) | IIS probability for two reads sampled within (resp. between) pools |
| $\hat{Q}_1^{\mathrm{pool}}$ (resp. $\hat{Q}_2^{\mathrm{pool}}$) | Unbiased estimator of the IIS probability for genes sampled within (resp. between) populations |

42

**Table 2  Definition of the $F_{\mathrm{ST}}$ estimators used in the text**

| Notation | Definition |
|---|---|
| $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ | Equation 12 |
| $\mathrm{FRP}_{13}$ | Ferretti et al. (2013) and Equations 16,20–21 |
| $\mathrm{NC}_{83}$ | Nei and Chesser (1983) |
| $\mathrm{PP2}_{\mathrm{d}}$ | Kofler et al. (2011) and Equations 16–18 |
| $\mathrm{PP2}_{\mathrm{a}}$ | Kofler et al. (2011) and Equation 19 |
| $\mathrm{WC}_{84}$ | Weir and Cockerham (1984) |

43

**Table 3   Overall $F_{\mathrm{ST}}$ estimates from multiple pools**

| $F_{\mathrm{ST}}$ | $n$ | Pool-seq | | Ind-seq |
|---|---|---|---|---|
| | | Cov. | $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ | WC$_{84}$ |
| 0.05 | 10 | 20× | 0.050 (0.002) | |
| 0.05 | 10 | 50× | 0.051 (0.002) | 0.050 (0.002) |
| 0.05 | 10 | 100× | 0.050 (0.002) | |
| | | | | |
| 0.05 | 100 | 20× | 0.050 (0.001) | |
| 0.05 | 100 | 50× | 0.050 (0.001) | 0.051 (0.001) |
| 0.05 | 100 | 100× | 0.050 (0.001) | |
| | | | | |
| 0.20 | 10 | 20× | 0.200 (0.002) | |
| 0.20 | 10 | 50× | 0.201 (0.002) | 0.201 (0.002) |
| 0.20 | 10 | 100× | 0.201 (0.002) | |
| | | | | |
| 0.20 | 100 | 20× | 0.201 (0.003) | |
| 0.20 | 100 | 50× | 0.202 (0.003) | 0.203 (0.003) |
| 0.20 | 100 | 100× | 0.203 (0.003) | |

Multilocus $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ estimates were computed for various conditions of expected $F_{\mathrm{ST}}$, pool size ($n$) and coverage (Cov.) in an island model with $n_{\mathrm{d}} = 8$ subpopulations (pools). The mean (RMSE) is over 50 independent simulated datasets, each made of 5,000 loci. For comparison, we computed multilocus WC$_{84}$ estimates from individual genotypes (Ind-seq).
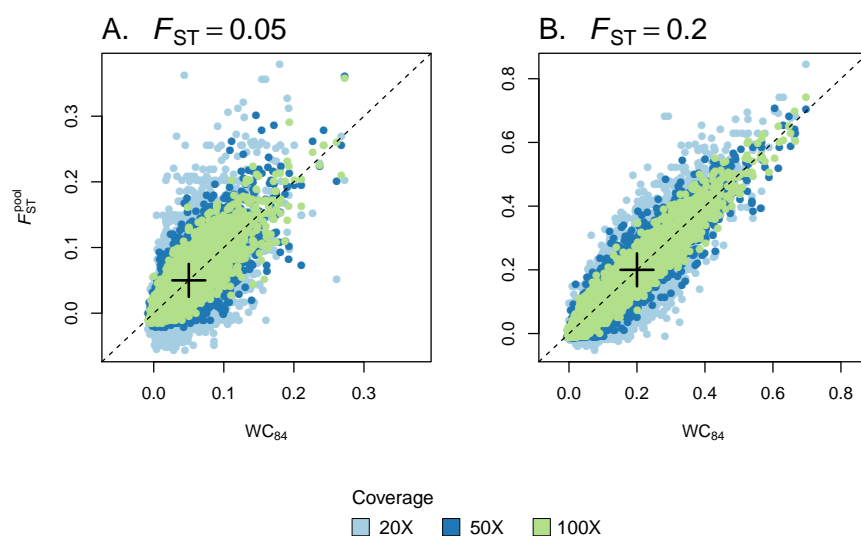
44

**Figure 1** Single-locus estimats of $F_{\mathrm{ST}}$. We compared single-locus estimates of $F_{\mathrm{ST}}$ based on allele count data inferred from individual genotypes (Ind-seq), using the $\mathrm{WC}_{84}$ estimator, to $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ estimates from Pool-seq data. We simulated 5,000 SNPs using `ms` in an island model with $n_{\mathrm{d}} = 8$ demes. We used two migration rates corresponding to $F_{\mathrm{ST}} = 0.05$ (A) and $F_{\mathrm{ST}} = 0.20$ (B). The size of each pool was fixed to 100. We show the results for different coverages (20X, 50X and 100X). In each graph, the cross indicates the simulated value of $F_{\mathrm{ST}}$.

**Figure 2** Precision and accuracy of pairwise estimators of $F_{\mathrm{ST}}$. We considered two estimators based on allele count data inferred from individual genotypes (Ind-seq): $WC_{84}$ and $NC_{83}$. For pooled data, we computed the two estimators implemented in the software package POPOOLATION2, that we refer to as $PP2_d$ and $PP2_a$, as well as the $FRP_{13}$ estimator and our estimator $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$. Each boxplot represents the distribution of multilocus $F_{\mathrm{ST}}$ estimates across all pairwise comparisons in an island model with $n_{\mathrm{d}} = 8$ demes, and across 50 independent replicates of the `ms` simulations. We used two migration rates, corresponding to $F_{\mathrm{ST}} = 0.05$ (A–B) and $F_{\mathrm{ST}} = 0.20$ (C–D). The size of each pool was either fixed to 10 (A and C) or to 100 (B and D). For Pool-seq data, we show the results for different coverages (20X, 50X and 100X). In each graph, the dashed line indicates the simulated value of $F_{\mathrm{ST}}$ and the dotted line indicates the median of the distribution of $NC_{83}$ estimates.
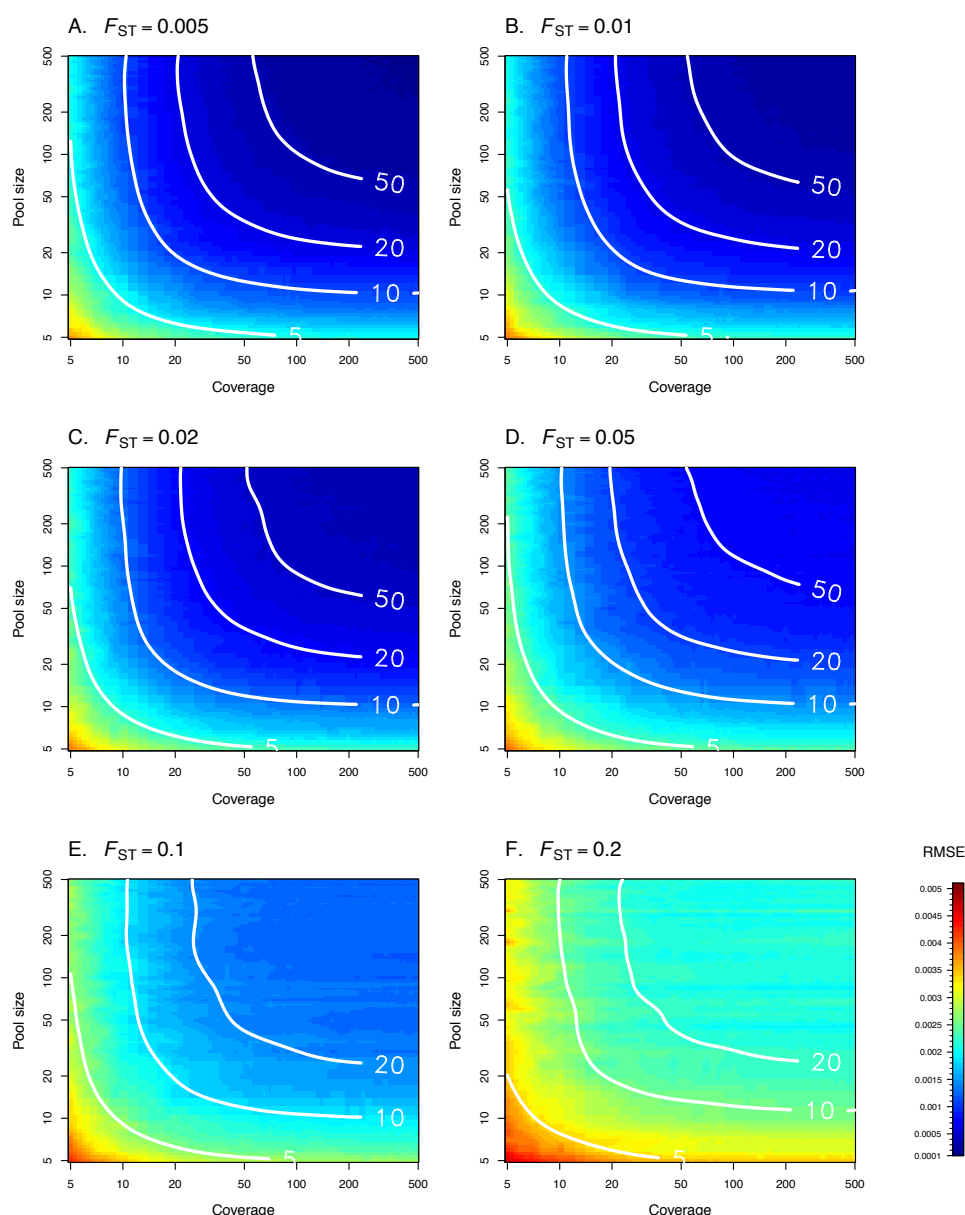
**Figure 3** Precision and accuracy of our estimator $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ as a function of pool size and coverage, for simulated $F_{\mathrm{ST}}$ values ranging from 0.005 to 0.2 (A–F). Each density plot, which represents the root mean squared error (RMSE) of the estimator $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$, was obtained using simple linear interpolation from a set of $44 \times 44$ pairs of pool size and coverage values. For each pool size and coverage, 500 replicates of 5,000 markers were simulated from an island model with $n_{\mathrm{d}} = 8$ demes. Plain white isolines represent the RMSE of the WC$_{84}$ estimator computed from Ind-seq data, for various sample sizes ($n = 5, 10, 20,$ and $50$). Each isoline was fitted using a thin plate spline regression with smoothing parameter $\lambda = 0.005$, implemented in the `fields` package for R (Nychka et al. 2017).
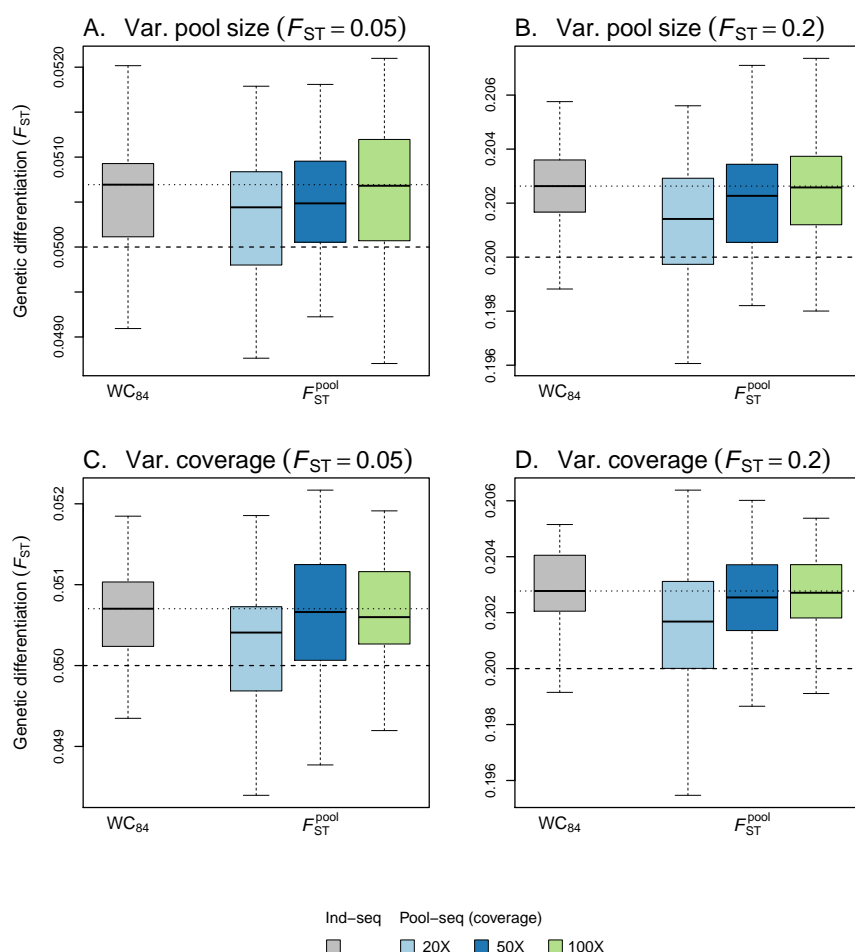
**Figure 4** Precision and accuracy of $F_{ST}$ estimates with varying pool size or varying coverage. Our estimator $\hat{F}_{ST}^{pool}$ was calculated from Pool-seq data over all demes and loci and compared to the estimator $WC_{84}$, computed from individual genotypes (Ind-seq). Each boxplot represents the distribution of multilocus $F_{ST}$ estimates across 50 independent replicates of the `ms` simulations. We used two migration rates, corresponding to $F_{ST} = 0.05$ (A and C) and $F_{ST} = 0.20$ (B and D). In A–B the pool size was variable across demes, with haploid sample size $n$ drawn independently for each deme from a Gaussian distribution with mean 100 and standard deviation 30; $n$ was rounded up to the nearest integer, with min. 20 and max. 300 haploids per deme. In C–D, the pool size was fixed ($n = 100$), and the coverage ($\delta_i$) was varying across demes and loci, with $\delta_i \sim \text{Pois}(\Delta)$ where $\Delta \in \{20, 50, 100\}$. For Pool-seq data, we show the results for different coverages (20X, 50X and 100X). In each graph, the dashed line indicates the simulated value of $F_{ST}$ and the dotted line indicates the median of the distribution of $WC_{84}$ estimates.
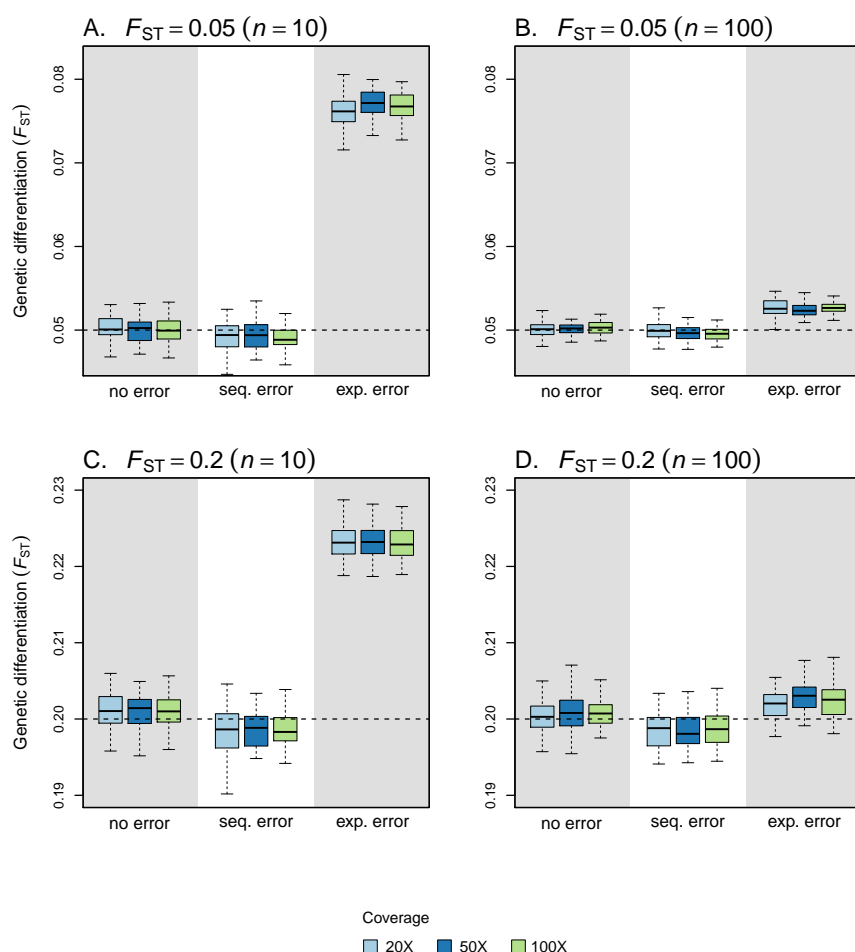
**Figure 5** Precision and accuracy of $F_{ST}$ estimates with sequencing and experimental errors. Our estimator $\hat{F}_{ST}^{pool}$ was computed from Pool-seq data over all demes and loci without error, with sequencing error (occurring at rate $\mu_e = 0.001$), and with experimental error ($\epsilon = 0.5$). Each boxplot represents the distribution of multilocus $F_{ST}$ estimates across 50 independent replicates of the `ms` simulations. We used two migration rates, corresponding to $F_{ST} = 0.05$ (A–B) or $F_{ST} = 0.20$ (C–D). The size of each pool was either fixed to 10 (A and C) or to 100 (B and D). For Pool-seq data, we show the results for different coverages (20X, 50X and 100X). In each graph, the dashed line indicates the simulated value of $F_{ST}$.
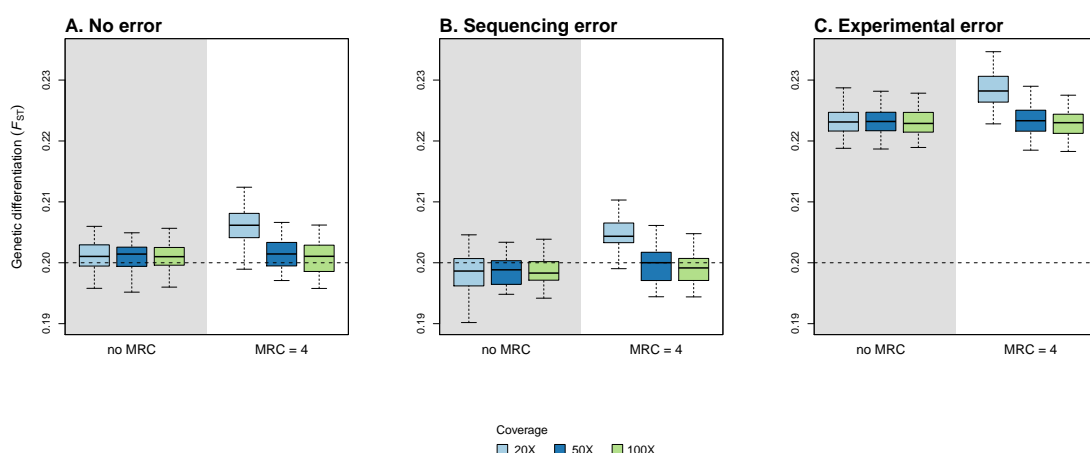
**Figure 6** Precision and accuracy of $F_{\mathrm{ST}}$ estimates with and without filtering. Our estimator $\hat{F}_{\mathrm{ST}}^{\mathrm{pool}}$ was computed from Pool-seq data over all demes and loci without error (A), with sequencing error (B) and with experimental error (C) (see the legend of Figure 5 for further details). For each case, we computed $F_{\mathrm{ST}}$ without filtering (no MRC) and with filtering (using a minimum read count MRC = 4). Each boxplot represents the distribution of multilocus $F_{\mathrm{ST}}$ estimates across 50 independent replicates of the `ms` simulations. We used a migration rate corresponding to $F_{\mathrm{ST}} = 0.20$, and pool size $n = 10$. We show the results for different coverages (20X, 50X and 100X). In each graph, the dashed line indicates the simulated value of $F_{\mathrm{ST}}$.
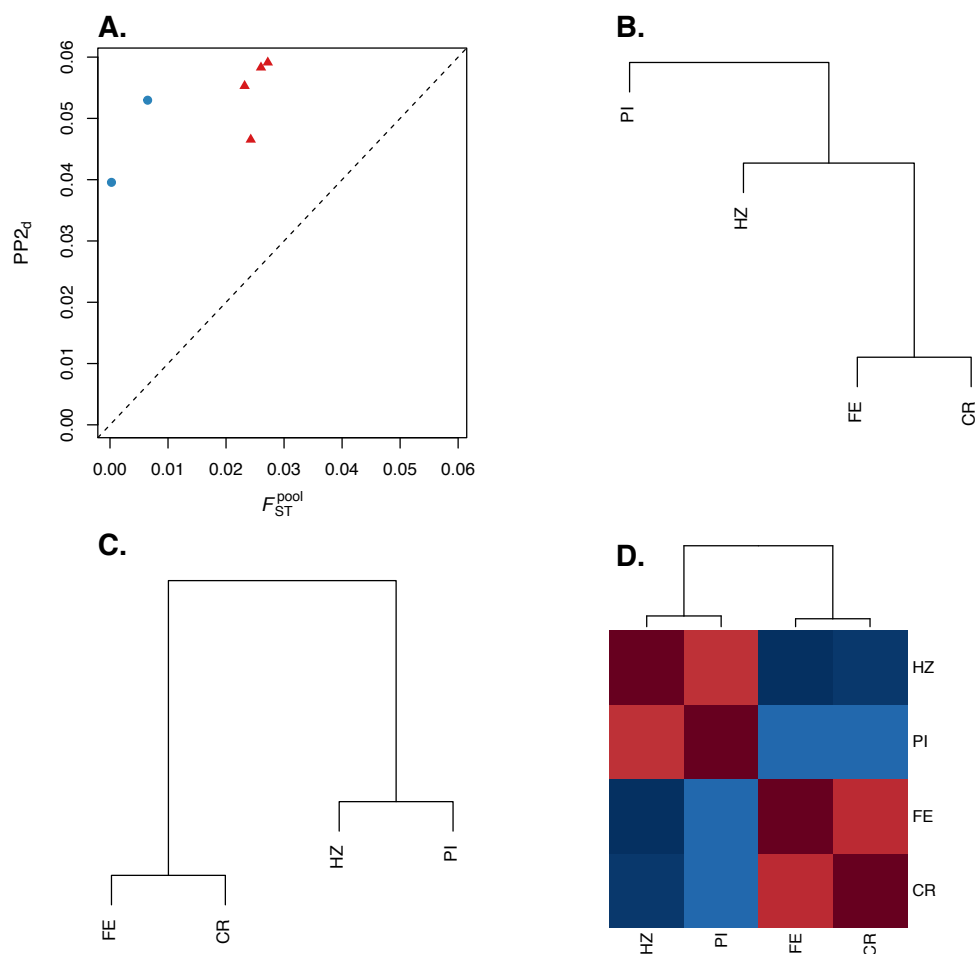
**Figure 7** Reanalysis of the prickly sculpin (*Cottus asper*) Pool-seq data. In (A) we compare the pairwise $F_{ST}$ estimates $PP2_d$, and $\hat{F}_{ST}^{pool}$ for all pairs of populations from the estuarine (CR and FE) and freshwater samples (PI and HZ). Within-ecotype comparisons are depicted as blue dots, and between-ecotype comparisons as red triangles. In (B–C) we show UPGMA hierarchical cluster analyses based on $PP2_d$ (B) and $\hat{F}_{ST}^{pool}$ (C) pairwise estimates. In (D), we show a heatmap representation of the scaled covariance matrix among the four *C. asper* populations, inferred from the Bayesian hierarchical model implemented in the software package BAYPASS.