

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

Running title: Genetic differentiation from pools

Keywords: F_{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

1

Abstract

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

The advent of high throughput sequencing and genotyping technologies enables the comparison of patterns of polymorphisms at a very large number of markers. While the characterization of genetic structure from individual sequencing data remains expensive for many non-model species, it has been shown that sequencing pools of individual DNAs (Pool-seq) represents an attractive and cost-effective alternative. However, analyzing sequence read counts from a DNA pool instead of individual genotypes raises statistical challenges in deriving correct estimates of genetic differentiation. In this article, we provide a method-of-moments estimator of F_{ST} for Pool-seq data, based on an analysis-of-variance framework. We show, by means of simulations, that this new estimator is unbiased, and outperforms previously proposed estimators. We evaluate the robustness of our estimator to model misspecification, such as sequencing errors and uneven contributions of individual DNAs to the pools. Finally, by reanalyzing published Pool-seq data of different ecotypes of the prickly sculpin *Cottus asper*, we show how the use of an unbiased F_{ST} estimator may question the interpretation of population structure inferred from previous analyses.

21

INTRODUCTION

22 It has long been recognized that the subdivision of species into subpopu-
23 lations, social groups and families fosters genetic differentiation (Wahlund
24 1928; Wright 1931). Characterizing genetic differentiation as a means to infer
25 unknown population structure is therefore fundamental to population genet-
26 ics, and finds applications in multiple domains, including conservation biol-
27 ogy, invasion biology, association mapping and forensics, among many others.
28 In the late 1940s and early 1950s, Malécot (1948) and Wright (1951) intro-
29 duced F -statistics to partition genetic variation within and between groups
30 of individuals (Holsinger and Weir 2009; Bhatia et al. 2013). Since then, the
31 estimation of F -statistics has become standard practice (see, e.g., Weir 1996;
32 Weir and Hill 2002; Weir 2012), and the most commonly used estimators of
33 F_{ST} have been developed in an analysis-of-variance framework (Cockerham
34 1969, 1973; Weir and Cockerham 1984), which can be recast in terms of prob-
35 abilities of identity of pairs of homologous genes (Cockerham and Weir 1987;
36 Rousset 2007; Weir and Goudet 2017).

37 Assuming that molecular markers are neutral, estimates of F_{ST} are typ-
38 ically used to quantify genetic structure in natural populations, which is
39 then interpreted as the result of demographic history (Holsinger and Weir
40 2009): large F_{ST} values are expected for small populations among which
41 dispersal is limited (Wright 1951), or between populations that have long
42 diverged in isolation from each other (Reynolds et al. 1983); when dispersal
43 is spatially restricted, a positive relationship between F_{ST} and the geograph-
44 ical distance for pairs of populations generally holds (Slatkin 1993; Rousset
45 1997). It has also been proposed to characterize the heterogeneity of F_{ST}

46 estimates across markers for identifying loci that are targeted by selection
47 (Cavalli-Sforza 1966; Lewontin and Krakauer 1973; Beaumont and Nichols
48 1996; Vitalis et al. 2001; Akey et al. 2002; Beaumont 2005; Weir et al. 2005;
49 Lotterhos and Whitlock 2014, 2015; Whitlock and Lotterhos 2015).

50 Next-generation sequencing (NGS) technologies provide unprecedented
51 amounts of polymorphism data in both model and non-model species (Elle-
52 gren 2014). Although the sequencing strategy initially involved individually
53 tagged samples in humans (The International HapMap Consortium 2005),
54 whole-genome sequencing of pools of individuals (Pool-seq) is being increas-
55 ingly popular for population genomic studies (Schlötterer et al. 2014). Be-
56 cause it consists in sequencing libraries of pooled DNA samples and does
57 not require individual tagging of sequences, Pool-seq provides genome-wide
58 polymorphism data at considerably lower cost than sequencing of individuals
59 (Schlötterer et al. 2014). However, non-equimolar amounts of DNA from all
60 individuals in a pool and stochastic variation in the amplification efficiency
61 of individual DNAs have raised concerns with respect to the accuracy of the
62 so-obtained allele frequency estimates, particularly at low sequencing depth
63 and with small pool sizes (Cutler and Jensen 2010; Ellegren 2014; Anderson
64 et al. 2014). Nonetheless, it has been shown that, at equal sequencing effort,
65 Pool-seq provides similar, if not more accurate, allele frequency estimates
66 than individual-based analyses (Futschik and Schlötterer 2010; Gautier et al.
67 2013). The problem is different for diversity and differentiation parameters,
68 which depend on second moments of allele frequencies or, equivalently, on
69 pairwise measures of genetic identity. With Pool-seq data, however, it is
70 impossible to distinguish pairs of reads that are identical because they were

71 sequenced from a single gene, from pairs of reads that are identical because
72 they were sequenced from two distinct genes that are identical in state (IIS)
73 (Ferretti et al. 2013).

74 Appropriate estimators of diversity and differentiation parameters must
75 therefore be sought, to account for both the sampling of individual genes
76 from the pool and the sampling of reads from these genes. There has been
77 several attempts to define estimators for the parameter F_{ST} for Pool-seq data
78 (Kofler et al. 2011; Ferretti et al. 2013), from ratios of heterozygosities (or
79 from probabilities of genetic identity between pairs of reads) within and be-
80 tween pools. In the following, we will argue that these estimators are biased
81 (i.e., they do not converge towards the expected value of the parameter),
82 and that some of them have undesired statistical properties (i.e., the bias
83 depends upon sample size and coverage). Here, following Cockerham (1969),
84 Cockerham (1973), Weir and Cockerham (1984), Weir (1996), Weir and Hill
85 (2002) and Rousset (2007), we define a method-of-moments estimator of the
86 parameter F_{ST} using an analysis-of-variance framework. We then evaluate
87 the accuracy and the precision of this estimator, based on the analysis of sim-
88 ulated datasets, and compare it to estimates defined in the software package
89 PoPoolation2 (Kofler et al. 2011), and in Ferretti et al. (2013). Furthermore,
90 we test the robustness of our estimators to model misspecifications (including
91 unequal contributions of individuals in pools, and sequencing errors). Finally,
92 we reanalyze the prickly sculpin (*Cottus asper*) Pool-seq data (published by
93 Dennenmoser et al. 2017), and show how the use of biased F_{ST} estimators in
94 previous analyses may challenge the interpretation of population structure.

95 Note that throughout this article, we use the term “gene” to designate a

96 segregating genetic unit (in the sense of the “Mendelian gene” from Orgogozo
97 et al. 2016). We further use the term “read” in a narrow sense, as a sequenced
98 copy of a gene. For the sake of simplicity, we will use the term “Ind-seq” to
99 refer to analyses based on individual data in which we further assume that
100 individual genotypes are called without error.

101

MODEL

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_{ST} is best defined as:

$$F_{\text{ST}} \equiv \frac{Q_1 - Q_2}{1 - Q_2} \quad (1)$$

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_{ST} for Pool-seq data, by decomposing
108 the total variance of gene 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_{ST} . Consider a sample of n_d subpopulations, each
113 of which is made of n_i genes ($i = 1, \dots, n_d$) sequenced in pools (hence n_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 ($j = 1, \dots, n_i$) 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 ($r = 1, \dots, c_{ij}$) 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.: $X_{ij:k} \equiv \sum_r X_{ijr:k} / c_{ij}$,
121 $X_{i:k} \equiv \sum_j \sum_r X_{ijr:k} / \sum_j c_{ij}$ and $X_{\dots:k} \equiv \sum_i \sum_j \sum_r X_{ijr:k} / \sum_i \sum_j c_{ij}$. The
122 analysis of variance is based on the computation of sums of squares, as fol-

123 lows:

$$\begin{aligned}
 \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - X_{\dots:k})^2 &= \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - X_{ij\cdot:k})^2 \\
 &+ \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ij\cdot:k} - X_{i\cdot\cdot:k})^2 \\
 &+ \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{i\cdot\cdot:k} - X_{\dots:k})^2 \\
 &\equiv SSR_{:k} + SSI_{:k} + SSP_{:k} \tag{2}
 \end{aligned}$$

124 As is shown in the Supplemental File S1, the expected sums of squares depend
 125 on the expectation of the allele frequency π_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 $C_1 \equiv \sum_i \sum_j c_{ij} = \sum_j c_{ij}$ 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;

137 see Equation A15 in Supplemental File S1). For reads between genes from
 138 different pools, we have:

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

139 where $C_2 \equiv \sum_i \left(\sum_j c_{ij} \right)^2$ and $D_2^* \equiv [\sum_i C_{1i} (C_{1i} + n_i - 1) / n_i] / C_1$ (see
 140 Equation A16 in Supplemental File S1). Rearranging Equations 4–5, and
 141 summing over alleles, we get:

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

142 and:

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

143 where $n_c \equiv (C_1 - C_2/C_1) / (D_2 - D_2^*)$. Let $MSI \equiv SSI / (C_1 - D_2)$ and
 144 $MSP \equiv SSP / (D_2 - D_2^*)$. Then:

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

145 which yields the method-of-moments estimator:

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

146 where

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

147 and:

$$MSP = \frac{1}{D_2 - D_2^*} \sum_k \sum_i^{n_d} C_{1i} (\hat{\pi}_{i:k} - \hat{\pi}_k)^2 \quad (11)$$

148 (see Equations A25 and A26 in Supplemental File S1). In Equations 10
 149 and 11, $\hat{\pi}_{i:k} \equiv X_{i\cdots:k}$ is the average frequency of reads of type k within the i th
 150 pool, and $\hat{\pi}_k \equiv X_{\cdots:k}$ is the average frequency of reads of type k in the full sam-
 151 ple. Note that from the definition of $X_{\cdots:k}$, $\hat{\pi}_k \equiv \sum_i \sum_j \sum_r X_{ijr:k} / \sum_i \sum_j c_{ij} =$
 152 $\sum_i C_{1i} \hat{\pi}_{i:k} / \sum_i C_{1i}$ is the weighted average of the sample frequencies with
 153 weights equal to the pool coverage. This is equivalent to the weighted
 154 analysis-of-variance in Cockerham (1973) (see also Weir and Cockerham 1984;
 155 Weir 1996; Weir and Hill 2002; Rousset 2007; Weir and Goudet 2017). Fi-
 156 nally, the full expression of $\hat{F}_{ST}^{\text{pool}}$ in terms of sample frequencies reads:

$$\hat{F}_{ST}^{\text{pool}} = \frac{\sum_k [(C_1 - D_2) \sum_i^{n_d} C_{1i} (\hat{\pi}_{i:k} - \hat{\pi}_k)^2 - (D_2 - D_2^*) \sum_i^{n_d} C_{1i} \hat{\pi}_{i:k} (1 - \hat{\pi}_{i:k})]}{\sum_k [(C_1 - D_2) \sum_i^{n_d} C_{1i} (\hat{\pi}_{i:k} - \hat{\pi}_k)^2 + (n_c - 1) (D_2 - D_2^*) \sum_i^{n_d} C_{1i} \hat{\pi}_{i:k} (1 - \hat{\pi}_{i:k})]} \quad (12)$$

157 If we take the limit case where each gene is sequenced exactly once, we
 158 recover the Ind-seq model: assuming $c_{ij} = 1$ for all (i, j) , then $C_1 = \sum_i^{n_d} n_i$,
 159 $C_2 = \sum_i^{n_d} n_i^2$, $D_2 = n_d$ and $D_2^* = 1$. Therefore, $n_c = (C_1 - C_2/C_1) / (n_d - 1)$,
 160 and Equation 9 reduces exactly to the estimator of F_{ST} for haploids: see Weir
 161 (1996), p. 182, and Rousset (2007), p. 977.

162 As in Reynolds et al. (1983), Weir and Cockerham (1984), Weir (1996)
 163 and Rousset (2007), a multilocus estimate is derived as the sum of locus
 164 specific numerators over the sum of locus-specific denominators:

$$\hat{F}_{ST} = \frac{\sum_l MSP_l - MSI_l}{\sum_l MSP_l + (n_c - 1) MSI_l} \quad (13)$$

165 where MSI and MSP are subscripted with l to denote the l th locus. For
166 Ind-seq data, Bhatia et al. (2013) refer to this multilocus estimate as a “ratio
167 of averages” by opposition to an “average of ratios”, which would consist in av-
168 eraging single-locus F_{ST} over loci. This approach is justified in the Appendix
169 of Weir and Cockerham (1984) and in Bhatia et al. (2013), who analyzed
170 both estimates by means of coalescent simulations. Note that Equation 13
171 assumes that the pool size is equal across loci. Also note that the construc-
172 tion of the estimator in Equation 13 is different from Weir and Cockerham’s
173 (1984). These authors defined their multilocus estimator as a ratio of sums
174 of components of variance (a , b and c in their notation) over loci, which give
175 the same weight to all loci, whatever the number of sampled genes at each
176 locus. Equation 13 follows GENEPOP’s rationale (Rousset 2008), which gives
177 instead more weight to loci that are more intensively covered.

178

MATERIALS AND METHODS

179 **Simulation study**

180 *Generating individual genotypes:* we first generated individual genotypes us-
181 ing `ms` (Hudson 2002), assuming an island model of population structure
182 (Wright 1931). For each simulated scenario, we considered 8 demes, each
183 made of $N = 5,000$ haploid individuals. The migration rate (m) was fixed
184 to achieve the desired value of F_{ST} (0.05 or 0.2), using Equation 6 in Rousset
185 (1996) leading, e.g., to $M \equiv 2Nm = 16.569$ for $F_{ST} = 0.05$ and $M = 3.489$ for
186 $F_{ST} = 0.20$. The mutation rate was set at $\mu = 10^{-6}$, giving $\theta \equiv 2N\mu = 0.01$.
187 We considered either fixed, or variable sample sizes across demes. In the lat-
188 ter case, the haploid sample size n was drawn independently for each deme
189 from a Gaussian distribution with mean 100 and standard deviation 30; this
190 number was rounded up to the nearest integer, with min. 20 and max. 300
191 haploids per deme. We generated a very large number of sequences for each
192 scenario, and sampled independent single nucleotide polymorphisms (SNPs)
193 from sequences with a single segregating site. Each scenario was replicated
194 50 times (500 times for Figures 3 and S2).

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

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

199 where $y_{i:k}$ is the number of genes of type k in the i th pool, n_i is the total
200 number of genes in pool i (haploid pool size), and δ_i is the simulated total
201 coverage for pool i . In the following, we either consider a fixed coverage,
202 with $\delta_i = \Delta$ for all pools and loci, or a varying coverage across pools and
203 loci, with $\delta_i \sim \text{Pois}(\Delta)$.

204 *Sequencing error:* we simulated sequencing errors occurring at rate $\mu_e =$
205 0.001, which is typical of Illumina sequencers (Glenn 2011; Ross et al. 2013).
206 We assumed that each sequencing error modifies the allelic type of a read to
207 one of three other possible states with equal probability (there are therefore
208 four allelic types in total, corresponding to four nucleotides). Note that
209 only biallelic markers are retained in the final datasets. Also note that,
210 since we initiated this procedure with polymorphic markers only, we neglect
211 sequencing errors that would create spurious SNPs from monomorphic sites.
212 However, such SNPs should be rare in real datasets, since markers with a
213 low minimum read count (MRC) are generally filtered out.

214 *Experimental error:* non-equimolar amounts of DNA from all individuals in
215 a pool and stochastic variation in the amplification efficiency of individual
216 DNAs are sources of experimental errors in pool sequencing. To simulate
217 experimental errors, we used the model derived by Gautier et al. (2013). In
218 this model, it is assumed that the contribution $\eta_{ij} = c_{ij}/C_{1i}$ of each gene j
219 to the total coverage of the i th pool (C_{1i}) follows a Dirichlet distribution:

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

220 where the parameter ρ controls the dispersion of gene contributions around
221 the value $\eta_{ij} = 1/n_i$, expected if all genes contributed equally to the pool of
222 reads. For convenience, we define the experimental error ϵ as the coefficient
223 of variation of η_{ij} , i.e.: $\epsilon \equiv \sqrt{\mathbb{V}(\eta_{ij})}/\mathbb{E}(\eta_{ij}) = \sqrt{(n_i - 1)/(\rho + 1)}$ (see Gautier
224 et al. 2013). When ϵ tends toward 0 (or equivalently when ρ tends to infinity),
225 all individuals contribute equally to the pool, and there is no experimental
226 error. We tested the robustness of our estimates to values of ϵ comprised
227 between 0.05 and 0.5. The case $\epsilon = 0.5$ could correspond, for example, to a
228 situation where (for $n_i = 10$) 5 individuals contribute $2.8\times$ more reads than
229 the other 5 individuals.

230 **Other estimators**

231 For the sake of clarity, a summary of the notation of the F_{ST} estimators used
232 throughout this article is given in Table 2.

233 PP2_d : this estimator of F_{ST} is implemented by default in the software
234 package POPOOLATION2 (Kofler et al. 2011). It is based on a definition of
235 the parameter F_{ST} as the overall reduction in average heterozygosity relative
236 to the total combined population (see, e.g., Nei and Chesser 1983):

$$\text{PP2}_d \equiv \frac{\hat{H}_T - \hat{H}_S}{\hat{H}_T} \quad (16)$$

237 where \hat{H}_S is the average heterozygosity within subpopulations, and \hat{H}_T is the
238 average heterozygosity in the total population (obtained by pooling together
239 all subpopulation to form a single virtual unit). In POPOOLATION2, \hat{H}_S is

240 the unweighted average of within-subpopulation heterozygosities:

$$\hat{H}_S = \frac{1}{n_d} \sum_i^{n_d} \left(\frac{n_i}{n_i - 1} \right) \left(\frac{C_{1i}}{C_{1i} - 1} \right) \left(1 - \sum_k \hat{\pi}_{i:k}^2 \right) \quad (17)$$

241 (using the notation from Table 1). Note that in POPOOLATION2, PP2_d is
242 restricted to the case of two subpopulations only ($n_d = 2$). The two ratios in
243 the right-hand side of Equation 17 are presumably borrowed from Nei (1978)
244 to provide an unbiased estimate, although we found no formal justification
245 for the expression in Equation 17 for Pool-seq data. The total heterozygosity
246 is computed as (using the notation from Table 1):

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

247 PP2_a : this is the alternative estimator of F_{ST} provided in the software
248 package POPOOLATION2. It is based on an interpretation by Kofler et al.
249 (2011) of Karlsson et al.'s (2007) estimator of F_{ST} , as:

$$PP2_a \equiv \frac{\hat{Q}_1^r - \hat{Q}_2^r}{1 - \hat{Q}_2^r} \quad (19)$$

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

258 they were sequenced from a single gene, from pairs of reads that are identical
259 because they were sequenced from distinct, yet IIS genes.

260 FRP_{13} : this estimator of F_{ST} was developed by Ferretti et al. (2013) (see
261 their Equations 3 and 10–13). Ferretti et al. (2013) use the same definition of
262 F_{ST} as in Equation 16 above, although they estimate heterozygosities within
263 and between pools as “average pairwise nucleotide diversities”, which, from
264 their definitions, are formally equivalent to IIS probabilities. In particular,
265 they estimate the average heterozygosity within pools as (using the notation
266 from Table 1):

$$\hat{H}_S = \frac{1}{n_d} \sum_i^{n_d} \left(\frac{n_i}{n_i - 1} \right) (1 - \hat{Q}_{1i}^r) \quad (20)$$

267 and the total heterozygosity among the n_d populations as:

$$\hat{H}_T = \frac{1}{n_d^2} \left[\sum_i^{n_d} \left(\frac{n_i}{n_i - 1} \right) (1 - \hat{Q}_{1i}^r) + \sum_{i \neq i'}^{n_d} (1 - \hat{Q}_{2ii'}^r) \right] \quad (21)$$

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

269 For the comparison of Ind-seq and Pool-seq datasets, we computed F_{ST} on
270 subsamples of 5,000 loci. These subsamples were defined so that only those
271 loci that were polymorphic in all coverage conditions were retained, and the
272 same loci were used for the analysis of the corresponding Ind-seq data. For
273 the latter, we used either the Nei and Chesser’s (1983) estimator based on a
274 ratio of heterozygosity (see Equation 16 above), hereafter denoted by NC_{83} , or
275 the analysis-of-variance estimator developed by Weir and Cockerham (1984),
276 hereafter denoted by WC_{84} .

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

278 ware environment for statistical computing, version 3.3.1 (R Core Team
279 2017). All these functions were carefully checked against available software
280 packages, to ensure that they provided strictly identical results.

281 **Application example: *Cottus asper***

282 Dennenmoser et al. (2017) investigated the genomic basis of adaption to
283 osmotic conditions in the prickly sculpin (*Cottus asper*), an abundant eury-
284 haline fish in northwestern North America. To do so, they sequenced the
285 whole-genome of pools of individuals from two estuarine populations (CR,
286 Capilano River Estuary; FE, Fraser River Estuary) and two freshwater pop-
287 ulations (PI, Pitt Lake and HZ, Hatzic Lake) in southern British Columbia
288 (Canada). We downloaded the four corresponding BAM files from the Dryad
289 Digital Repository (doi: 10.5061/dryad.2qg01) and combined them into a sin-
290 gular mpileup file using `SAMtools` version 0.1.19 (Li et al. 2009) with default
291 options, except the maximum depth per BAM that was set to 5,000 reads.
292 The resulting file was further processed using a custom `awk` script, to call
293 SNPs and compute read counts, after discarding bases with a Base Align-
294 ment Quality (BAQ) score lower than 25. A position was then considered
295 as a SNP if: (*i*) only two different nucleotides with a read count > 1 were
296 observed (nucleotides with ≤ 1 read being considered as a sequencing error);
297 (*ii*) the coverage was comprised between 10 and 300 in each of the four align-
298 ment files; (*iii*) the minor allele frequency, as computed from read counts,
299 was ≥ 0.01 in the four populations. The final data set consisted of 608,879
300 SNPs.

301 Our aim here was to compare the population structure inferred from pair-
302 wise estimates of F_{ST} , using the estimator \hat{F}_{ST}^{pool} on the one hand, and PP2_d

303 on the other hand. Then, to conclude on which of the two estimators per-
304 forms better, we compared the population structure inferred from \hat{F}_{ST}^{pool} and
305 $PP2_d$ to that inferred from the Bayesian hierarchical model implemented in
306 the software package BAYPASS (Gautier 2015). BAYPASS allows indeed the
307 robust estimation of the scaled covariance matrix of allele frequencies across
308 populations for Pool-seq data, which is known to be informative about pop-
309 ulation history (Pickrell and Pritchard 2012). The elements of the estimated
310 matrix can be interpreted as pairwise and population-specific estimates of
311 differentiation (Coop et al. 2010), and therefore provide a comprehensive
312 description of population structure that makes full use of the available data.

313 **Data availability**

314 The authors state that all data necessary for confirming the conclusions
315 presented in this article are fully represented within the article, figures,
316 and tables. Supplemental Tables S1–S3 and Figures S1–S4 are available at
317 FigShare, along with a complete derivation of the model in the Supplemental
318 File S1 at FigShare.

319

RESULTS

320 **Comparing Ind-seq and Pool-seq estimates of F_{ST}**

321 Single-locus estimates \hat{F}_{ST}^{pool} are highly correlated with the classical estimates
322 WC_{84} (Weir and Cockerham 1984) computed on the individual data that were
323 used to generate the pools in our simulations (see Figure 1). The variance of
324 \hat{F}_{ST}^{pool} across independent replicates decreases as the coverage increases. The
325 correlation between \hat{F}_{ST}^{pool} and WC_{84} is stronger for multilocus estimates (see
326 Figure S1A).

327 **Comparing Pool-seq estimators of F_{ST}**

328 We found that our estimator \hat{F}_{ST}^{pool} has extremely low bias ($< 0.5\%$ over
329 all scenarios tested: see Tables 3 and S1-S3). In other words, the average
330 estimates across multiple loci and replicates closely equals the expected value
331 of the F_{ST} parameter, as given by Equation 6 in Rousset (1996), which is
332 based on the computation of IIS probabilities in an island model of population
333 structure. In all the situations examined, the bias did neither depend on the
334 sample size (i.e., the size of each pool) nor on the coverage (see Figure 2).
335 Only the variance of the estimator across independent replicates decreases as
336 the sample size increases and/or as the coverage increases. At high coverage,
337 the mean and root mean squared error (RMSE) of \hat{F}_{ST}^{pool} over independent
338 replicates are virtually indistinguishable from that of the WC_{84} estimator
339 (see Table S1).

340 Figure 3 shows the RMSE of F_{ST} estimates for a wide range of pool sizes
341 and coverage. The RMSE decreases as the pool size and/or the coverage
342 increases. The F_{ST} estimates are more precise and accurate when differen-

343 tiation is low. Figure 3 provides some clues to evaluate the pool size and
344 the coverage that is necessary to achieve the same RMSE than for Ind-seq
345 data. Consider, for example, the case of samples of $n = 20$ haploids. For
346 $F_{ST} \leq 0.05$ (in the conditions of our simulations), the RMSE of F_{ST} estimates
347 based on Pool-seq data tends to the RMSE of F_{ST} estimates based on Ind-seq
348 data either by sequencing pools of ca. 200 haploids at 20X, or by sequencing
349 pools of 20 haploids at ca. 200X. However, the same precision and accuracy
350 are achieved by sequencing ca. 50 haploids at ca. 50X.

351 Conversely, we found that PP2_d (the default estimator of F_{ST} imple-
352 mented in the software package POPOOLATION2) is biased when compared
353 to the expected value of the parameter. We observed that the bias depends
354 on both the sample size, and the coverage (see Figure 2). We note that, as the
355 coverage and the sample size increase, PP2_d converges to the estimator NC₈₃
356 (Nei and Chesser 1983) computed from individual data (see Figure S1B).
357 This argument was used by Kofler et al. (2011) to validate the approach,
358 even though the estimates PP2_d depart from the true value of the parameter
359 (Figure S1B–C).

360 The second of the two estimators of F_{ST} implemented in POPOOLATION2,
361 that we refer to as PP2_a, is also biased (see Figure 2). We note that the bias
362 decreases as the sample size increases. However, the bias does not depend
363 on the coverage (only the variance over independent replicates does). The
364 estimator developed by Ferretti et al. (2013), that we refer to as FRP₁₃, is
365 also biased (see Figure 2). However, the bias does neither depend on the pool
366 size, nor on the coverage (only the variance over independent replicates does).
367 FRP₁₃ converges to the estimator NC₈₃, computed from individual data (see

368 Figure 2). At high coverage, the mean and RMSE over independent replicates
369 are virtually indistinguishable from that of the NC_{83} estimator.

370 Last, we stress out that our estimator \hat{F}_{ST}^{pool} provides estimates for multiple
371 populations, and is therefore not restricted to pairwise analyses, contrary to
372 POPOOLATION2's estimators. We show that, even at low sample size and low
373 coverage, Pool-seq estimates of differentiation are virtually indistinguishable
374 from classical estimates for Ind-seq data (see Table 3).

375 **Robustness to unbalanced pool sizes and variable sequencing cov-** 376 **erage**

377 We evaluated the accuracy and the precision of the estimator \hat{F}_{ST}^{pool} when sam-
378 ple sizes differ across pools, and when the coverage varies across pools and loci
379 (see Figure 4). We found that, at low coverage, unequal sampling or variable
380 coverage causes a negligible departure from the median of WC_{84} estimates
381 computed on individual data, which vanishes as the coverage increases. At
382 100X coverage, the distribution of \hat{F}_{ST}^{pool} estimates is almost indistinguishable
383 from that of WC_{84} (see Figure 4 and Tables S2–S3).

384 **Robustness to sequencing and experimental errors**

385 Figure 5 shows that sequencing errors cause a negligible negative bias for
386 \hat{F}_{ST}^{pool} estimates. Filtering (using a minimum read count of 4) improves es-
387 timation slightly, but only at high coverage (Figure 6B). It must be noted,
388 though, that filtering increases the bias in the absence of sequencing error,
389 especially at low coverage (Figure 6A). With experimental error, i.e., when
390 individuals do not contribute evenly to the final set of reads, we observed a
391 positive bias for \hat{F}_{ST}^{pool} estimates (Figure 5). We note that the bias decreases

392 as the size of the pools increases. Figure S2 shows the RMSE of F_{ST} esti-
393 mates for a wider range of pool sizes, coverage and experimental error rate.
394 For $\epsilon \geq 0.25$, increasing the coverage cannot improve the quality of the in-
395 ference, if the pool size is too small. When Pool-seq experiments are prone
396 to large experimental error rates, increasing the size of pools is the only way
397 to improve the estimation of F_{ST} . Filtering (using a minimum read count of
398 4) does not improve estimation (Figure 6C).

399 **Application example**

400 The reanalysis of the prickly sculpin data revealed larger pairwise estimates of
401 multilocus F_{ST} using PP2_d estimator, as compared to \hat{F}_{ST}^{pool} (see Figure 7A).
402 Furthermore, we found that \hat{F}_{ST}^{pool} estimates are smaller for within-ecotype
403 pairwise comparisons as compared to between-ecotype comparisons. There-
404 fore, the inferred relationships between samples based on pairwise \hat{F}_{ST}^{pool} esti-
405 mates show a clear-cut structure, separating the two estuarine samples from
406 the freshwater ones (see Figure 7C). We did not recover the same structure
407 using PP2_d estimates (see Figure 7B). Supportingly, the scaled covariance
408 matrix of allele frequencies across samples is consistent with the structure
409 inferred from \hat{F}_{ST}^{pool} estimates (see Figure 7D).

410

DISCUSSION

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

419 The naive approach that consists in computing F_{ST} from read counts, as
420 if they were allele counts (e.g., as in Chen et al. 2016), ignores the extra
421 variance brought by the random sampling of reads from the gene pool dur-
422 ing Pool-seq experiments. Furthermore, such computation fails to consider
423 the actual number of lineages in the pool (haploid pool size). Altogether,
424 these limits may result in severely biased estimates of differentiation when
425 the pool size is low (see Figure S3). A possible alternative is to compute F_{ST}
426 from allele counts imputed from read counts using a maximum-likelihood
427 approach conditional on the haploid size of the pools (e.g., as in Smadja
428 et al. 2012; Leblois et al. 2018), or from allele frequencies estimated using a
429 model-based method that accounts for the sampling effects and the sequenc-
430 ing error probabilities inherent to pooled NGS experiments (see Fariello et al.
431 2017). However, these latter approaches may only be accurate in situations
432 where the coverage is much larger than pool size, allowing to reduce sampling
433 variance of reads (see Figure S3).

434 Here, we therefore developed a new estimator of the parameter F_{ST} for

435 Pool-seq data, in an analysis-of-variance framework (Cockerham 1969, 1973).
436 The accuracy of this estimator is barely distinguishable from that of the
437 Weir and Cockerham's (1984) estimator for individual data. Furthermore,
438 does neither depend on the pool size nor on the coverage, and is robust
439 to unequal pool sizes and varying coverage across demes and loci. In our
440 analysis the frequency of reads within pools is a weighted average of the
441 sample frequencies with weights equal to the pool coverage. Therefore, our
442 approach follows Cockerham's (1973) one, which he referred to as a weighted
443 analysis-of-variance (see also Weir and Cockerham 1984; Weir 1996; Weir and
444 Hill 2002; Weir and Goudet 2017).

445 With unequal pool sizes, weighted and unweighted analyses differ. As dis-
446 cussed recently in Weir and Goudet (2017), the unweighted approach seems
447 appropriate when the between component exceed the within component, i.e.
448 when F_{ST} is large (Tukey 1957). It turns out that optimal weighting depends
449 upon the parameter to be estimated (Cockerham 1973) and is only efficient
450 at lower levels of differentiation (Robertson 1962). In a likelihood analysis
451 of the island model, Rousset (2007) derived asymptotically efficient weights
452 that are proportional to n_i^2 for the sum of squares of different samples (i.e.,
453 as in Robertson 1962). To the best of our knowledge, such optimal weighting
454 has never been considered in the literature. Nevertheless, if these arguments
455 are true for estimators of variance components, they do not necessarily apply
456 to estimates of intra-class correlations (Cockerham 1973).

457 **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,
462 introduced in population genetics by Cockerham (1969, 1973), estimators
463 based on the computation of probabilities of identical response within and
464 between groups have been proposed (see, e.g., Fleiss 1971; Fleiss and Cuzick
465 1979; Mak 1988; Ridout et al. 1999; Wu et al. 2012), which were originally
466 referred to as kappa-type statistics (Fleiss 1971; Landis and Koch 1977).
467 These estimators have later been endorsed in population genetics, where the
468 “probability of identical response” was then interpreted as the frequency with
469 which the genes are alike (Cockerham 1973; Cockerham and Weir 1987; Weir
470 1996; Rousset 2007; Weir and Goudet 2017).

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

485 found that analysis-of-variance estimates have better precision and accuracy
486 than IIS-based estimates, particularly for low levels of differentiation (see
487 Figure S4). Interestingly, we found that IIS-based estimates of F_{ST} for Pool-
488 seq data have generally lower bias and variance if the overall estimates of IIS
489 probabilities within and between pools are computed as unweighted averages
490 of population-specific or pairwise estimates (see Equations A39 and A43), as
491 compared to weighted averages. Equation A28 further shows that our esti-
492 mator may be rewritten as a function close to $(\hat{Q}_1 - \hat{Q}_2) / (1 - \hat{Q}_2)$, except
493 that it also depends on the sums $\sum_i (\hat{Q}_{1i} - \hat{Q}_1)$ in both the numerator and
494 the denominator. This suggests that if the Q_{1i} 's differ among subpopulations,
495 then our estimator provides an estimate of an average of population-specific
496 F_{ST} (Weir and Hill 2002; Weir and Goudet 2017).

497 It follows from the derivations in the Supplemental File S1 that the es-
498 timator PP2_a (Equation 19) is biased, because the IIS probability between
499 pairs of reads within a pool (\hat{Q}_1^r) is a biased estimator of the IIS probability
500 between pairs of distinct genes in that pool (see Equation A34 in Supplemen-
501 tal File S1). This is so, because the former confounds pairs of reads that are
502 identical because they were sequenced from a single gene copy, from pairs of
503 reads that are identical because they were sequenced from distinct, yet IIS
504 genes.

505 A more justified estimator of F_{ST} has been proposed by Ferretti et al.
506 (2013), based on previous developments by Futschik and Schlötterer (2010).
507 Note that, although they defined F_{ST} as a ratio of functions of heterozygosi-
508 ties, they actually worked with IIS probabilities (see Equations 20 and 21).
509 However, although their Equation 20 is strictly identical to our Equation A34

510 in Supplemental File S1, we note that they computed the total heterozygosity
511 by integrating over pairs of genes sampled both within and between popula-
512 tions (see Equation 21), which may explain the observed bias (see Figure 2).

513 **Comparison with alternative estimators**

514 An alternative framework to Weir and Cockerham's (1984) analysis-of-variance
515 has been developed by Masatoshi Nei and coworkers to estimate F_{ST} from
516 gene diversities (Nei 1973, 1977; Nei and Chesser 1983; Nei 1986). The es-
517 timator $PP2_d$ (see Equations 16–18) implemented in the software package
518 `POPOOLATION2` (Kofler et al. 2011) follows this logic. However, it has long
519 been recognized that both frameworks are fundamentally different in that the
520 analysis-of-variance approach considers both statistical and genetic (or evo-
521 lutionary) sampling, whereas Nei and coworkers' approach do not (Weir and
522 Cockerham 1984; Excoffier 2007; Holsinger and Weir 2009). Furthermore,
523 the expectation of Nei and coworkers' estimators depend upon the number
524 of sampled populations, with a larger bias for lower numbers of sampled pop-
525 ulations (Goudet 1993; Excoffier 2007; Weir and Goudet 2017). This is so,
526 because the computation of the total diversity in Equations 18 and 21 includes
527 the comparison of pairs of genes from the same subpopulation, whereas the
528 computation of IIS probabilities between subpopulations do not (see, e.g.,
529 Excoffier 2007). Therefore, we do not recommend using the estimator $PP2_d$
530 implemented in the software package `POPOOLATION2` (Kofler et al. 2011).

531 **Applications in evolutionary ecology studies**

532 Pool-seq is being increasingly used in many application domains (Schlötterer
533 et al. 2014), such as conservation genetics (see, e.g., Fuentes-Pardo 2017),

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

540 Here, we reanalyzed the Pool-seq data produced by Dennenmoser et al.
541 (2017), who investigated the adaptive genomic divergence between freshwa-
542 ter and brackish-water ecotypes of the prickly sculpin *C. asper*, an abundant
543 euryhaline fish in northwestern North America. Measuring pairwise genetic
544 differentiation between samples using \hat{F}_{ST}^{pool} , we found a clear-cut structure
545 separating the freshwater from the brackish-water ecotypes. Such genetic
546 structure supports the hypothesis that populations are locally adapted to
547 osmotic conditions in these two contrasted habitats, as discussed in Den-
548 nenmoser et al. (2017). This structure, which is at odds with that inferred
549 from $PP2_d$ estimates, is not only supported by the scaled covariance ma-
550 trix of allele frequencies, but also by previous microsatellite-based studies,
551 who showed that populations were genetically more differentiated between
552 ecotypes than within ecotypes (Dennenmoser et al. 2014, 2015).

553 **Limits of the model and perspectives**

554 We have shown that the stronger source of bias for the \hat{F}_{ST}^{pool} estimate is un-
555 equal contributions of individuals in pools. This is so, because we assume in
556 our model that the read counts are multinomially distributed, which supposes
557 that all genes contribute equally to the pool of reads (Gautier et al. 2013),
558 i.e. that there is no variation in DNA yield across individuals and that all

559 genes have equal sequencing coverage (Rode et al. 2018). Because the effect
560 of unequal contribution is expected to be stronger with small pool sizes, it
561 has been recommended to use pool-seq with at least 50 diploid individuals
562 per pool (Lynch et al. 2014; Schlötterer et al. 2014). However, this limit may
563 be overly conservative for allele frequency estimates (Rode et al. 2018), and
564 we have shown here that we can achieve very good precision and accuracy
565 of F_{ST} estimates with smaller pool sizes. Furthermore, because genotypic in-
566 formation is lost during Pool-seq experiments, we assume in our derivations
567 that pools are haploid (and therefore that F_{IS} is nil). Analyzing non-random
568 mating populations (e.g., in selfing species) is therefore problematic.

569 Finally, our model, as in Weir and Cockerham (1984), formally assumes
570 that all populations provide independent replicates of some evolutionary pro-
571 cess (Excoffier 2007; Holsinger and Weir 2009). This may be unrealistic in
572 many natural populations, which motivated Weir and Hill (2002) to derive a
573 population-specific estimator of F_{ST} for Ind-seq data (see also Vitalis et al.
574 2001). Even though the use of Weir and Hill's (2002) estimator is still scarce
575 in the literature (but see Weir et al. 2005; Vitalis 2012), Weir and Goudet
576 (2017) recently proposed a re-interpretation of population-specific estimates
577 of F_{ST} in terms of allelic matching proportions, which are strictly equivalent
578 to IIS probabilities between pairs of genes. It would therefore be straight-
579 forward to extend Weir and Goudet's (2017) estimator of population-specific
580 F_{ST} for the analysis of Pool-seq data, using the unbiased estimates of IIS
581 probabilities provided in the Supplemental File S1.

582

DATA ACCESSIBILITY

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

586

ACKNOWLEDGEMENTS

587 We thank Alexandre Dehne-Garcia for his assistance in using computer farms.
588 Analyses were performed on the genotoul bioinformatics platform Toulouse
589 Midi-Pyrénées (bioinfo.genotoul.fr) and the CBGP HPC computational
590 platform. This work is part of Valentin Hivert's Ph.D., who was supported
591 by a grant from the INRA's Plant Health and Environment (SPE) Division,
592 and by the BiodivERsA project EXOTIC (ANR-13-EBID-0001). Part of this
593 work was supported by the ANR project SWING (ANR-16-CE02-0015) of
594 the French National Research Agency, and by the CORBAM project of the
595 French region Hauts-de-France. We thank two anonymous reviewers for their
596 positive comments and suggestions.

597

Literature Cited

598 Akey, J. M., Zhang, G., Jin, L., and Shriver, M. D. (2002). Interrogating a
599 high-density SNP map for signatures of natural selection. *Genome Res.*,
600 12:1805–1814.

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

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

606 Beaumont, M. A. and Nichols, R. A. (1996). Evaluating loci for use in the
607 genetic analysis of population structure. *Proc. R. Soc. Lond. B*, 263:1619–
608 1626.

609 Bhatia, G., Patterson, N., Sankararaman, S., and Price, A. L. (2013). Esti-
610 mating and interpreting F_{ST} : the impact of rare variants. *Genome Res.*,
611 23:1514–1521.

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

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

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

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

- 619 Cockerham, C. C. and Weir, B. S. (1987). Analyses of gene frequencies. *Proc.*
620 *Natl. Acad. Sci. USA*, 84:8512–8514.
- 621 Collet, J. M., Fuentes, S., Hesketh, J., Hill, M. S., Innocenti, P., Morrow,
622 E. H., Fowler, K., and Reuter, M. (2016). Rapid evolution of the intersex-
623 ual genetic correlation for fitness in *Drosophila melanogaster*. *Evolution*,
624 70:781–795.
- 625 Coop, G., Witonsky, D., Di Rienzo, A., and Pritchard, J. K. (2010). Us-
626 ing environmental correlations to identify loci underlying local adaptation.
627 *Genetics*, 185:1411–1423.
- 628 Cutler, D. J. and Jensen, J. D. (2010). To pool, or not to pool? *Genetics*,
629 186:41–43.
- 630 Dennenmoser, S., Nolte, A. W., Vamosi, S. M., and Rogers S, M. (2015). Phy-
631 logeography of the prickly sculpin (*Cottus asper*) in north-western North
632 America reveals parallel phenotypic evolution across multiple coastal-
633 inland colonizations. *J. Biogeogr.*, 42:1626–1638.
- 634 Dennenmoser, S., Rogers, S. M., and Vamosi, S. M. (2014). Genetic pop-
635 ulation structure in prickly sculpin (*Cottus asper*) reflects isolation-by-
636 environment between two life-history ecotypes. *Biol. J. Linnean Soc.*,
637 113:943–957.
- 638 Dennenmoser, S., Vamosi, S. M., Nolte, S. W., and Rogers, S. M. (2017).
639 Adaptive genomic divergence under high gene flow between freshwater and
640 brackish-water ecotypes of prickly sculpin (*Cottus asper*) revealed by Pool-
641 Seq. *Mol. Ecol.*, 26:25–42.

- 642 Dexter, E., Bollens, S. M., Cordell, J., Soh, H. Y., Rollwagen-Bollens, G.,
643 Pfeifer, S. P., Goudet, J., and Vuilleumier, S. (2018). A genetic reconstruc-
644 tion of the invasion of the calanoid copepod *Pseudodiaptomus inopinus*
645 across the North American Pacific Coast. *Biol. Invasions*, 20:1577–1595.
- 646 Ellegren, H. (2014). Genome sequencing and population genomics in non-
647 model organisms. *Trends Ecol. Evol.*, 29:51–63.
- 648 Excoffier, L. (2007). Analysis of population subdivision. In Balding, D. J.,
649 Bishop, M., and Cannings, C., editors, *Handbook of Statistical Genetics*,
650 pages 980–1020, Chichester. John Wiley & Sons, Ltd.
- 651 Fariello, M. I. and Boitard, S., Mercier, S., Robelin, D., Faraut, T., Arnould,
652 C., Recoquillay, J., Bouchez, O., Salin, G., Dehais, P., Gourichon, D., Ler-
653 oux, S., Pitel, F., Leterrier, C., and SanCristobal, M. (2017). Accounting
654 for Linkage Disequilibrium in genome scans for selection without individual
655 genotypes : the local score approach. *Mol. Ecol.*, 26:3700–3714.
- 656 Ferretti, L., Ramos Onsins, S., and Pérez-Enciso, M. (2013). Population
657 genomics from pool sequencing. *Mol. Ecol.*, 22:5561–5576.
- 658 Fischer, M. C., Rellstab, C., Leuzinger, M., Roumet, M., Gugerli, F.,
659 Shimizu, K. K., Holderegger, R., and Widmer, A. (2017). Estimating ge-
660 nomic diversity and population differentiation – an empirical comparison
661 of microsatellite and SNP variation in *Arabidopsis halleri*. *BMC Genomics*,
662 18:69.
- 663 Fleiss, J. L. (1971). Measuring nominal scale agreement among many raters.
664 *Psychol. Bull.*, 76:378–382.

- 665 Fleiss, J. L. and Cuzick, J. (1979). The reliability of dichotomous judgements:
666 Unequal numbers of judges per subject. *Appl. Psychol. Meas.*, 3:537–542.
- 667 Fuentes-Pardo, A. P. and Ruzzente, D. E. (2017). Whole-genome sequencing
668 approaches for conservation biology: Advantages, limitations and practical
669 recommendations. *Mol. Ecol.*, 26:5369–5406.
- 670 Futschik, A. and Schlötterer, C. (2010). The next generation of molecu-
671 lar markers from massively parallel sequencing of pooled DNA samples.
672 *Genetics*, 186:207–218.
- 673 Gautier, M. (2015). Genome-wide scan for adaptive divergence and associa-
674 tion with population-specific covariates. *Genetics*, 201:1555–1579.
- 675 Gautier, M., Gharbi, K., Cezaerd, T., Galan, M., Loiseau, A., Thomson, M.,
676 Pudlo, P., Kerdelhué, C., and Estoup, A. (2013). Estimation of popula-
677 tion allele frequencies from next-generation sequencing data: pool-versus
678 individual-based genotyping. *Mol. Ecol.*, 22:3766–3779.
- 679 Glenn, T. C. (2011). Field guide to next-generation DNA sequencers. *Mol.*
680 *Ecol. Resour.*, 11:759–769.
- 681 Goudet, J. (1993). *The genetics of geographically structured populations*. PhD
682 thesis, University of Wales, Bangor.
- 683 Holsinger, K. S. and Weir, B. S. (2009). Genetics in geographically structured
684 populations: defining, estimating and interpreting F_{ST} . *Nat. Rev. Genet.*,
685 10:639–650.
- 686 Hudson, R. R. (2002). Generating samples under a Wright-Fisher neutral
687 model of genetic variation. *Bioinformatics*, 18:337–338.

- 688 Karlsson, E. K., Baranowska, I., Wade, C. M., Salmon Hillbertz, N. H. C.,
689 Zody, M. C., Anderson, N., Biagi, T. M., Patterson, N., Pielberg, G. R.,
690 Kulbokas, E. J., Comstock, K. E., Keller, E. T., Mesirov, J. P., von Euler,
691 H., Kämpe, O., Hedhammar, A., Lander, E. S., Andersson, G., Andersson,
692 L., and Lindblad-Toh, K. (2007). Efficient mapping of Mendelian traits in
693 dogs through genome-wide association. *Nat. Genet.*, 39:1321–1328.
- 694 Kofler, R., Pandey, R. V., and Schlötterer, C. (2011). PoPoolation2: identi-
695 fying differentiation between populations using sequencing of pooled DNA
696 samples (Pool-Seq). *Bioinformatics*, 27:3435–3436.
- 697 Landis, J. R. and Koch, G. G. (1977). A one-way components of variance
698 model for categorical data. *Biometrics*, 33:671–679.
- 699 Leblois, R., Gautier, M., Rohfritsch, A., Foucaud, J., Burban, C., Galan, M.,
700 Loiseau, A., Sauné, L., Branco, M., Gharbi, K., Vitalis, R., and Kerdelhué,
701 C. (2018). Deciphering the demographic history of allochronic differentia-
702 tion in the pine processionary moth *Thaumetopoea pityocampa*. *Mol. Ecol.*,
703 27:264–278.
- 704 Lewontin, R. C. and Krakauer, J. (1973). Distribution of gene frequency as
705 a test of the theory of the selective neutrality of polymorphism. *Genetics*,
706 74:175–195.
- 707 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth,
708 G., Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing
709 Subgroup (2009). The Sequence Alignment/Map format and SAMtools.
710 *Bioinformatics*, 25:2078–2079.

- 711 Lotterhos, K. E. and Whitlock, M. C. (2014). Evaluation of demographic
712 history and neutral parameterization on the performance of F_{ST} outlier
713 tests. *Mol. Ecol.*, 23:2178–2192.
- 714 Lotterhos, K. E. and Whitlock, M. C. (2015). The relative power of genome
715 scans to detect local adaptation depends on sampling design and statistical
716 method. *Mol. Ecol.*, 24:1031–1046.
- 717 Lynch, M., Bost, D., Wilson, S., Maruki, T., and Harrison, S. (2014).
718 Population-genetic inference from pooled-sequencing data. *Genome Biol.*
719 *Evol.*, 6:1210–1218.
- 720 Mak, T. K. (1988). Analysing intraclass correlation for dichotomous vari-
721 ables. *J. R. Stat. Soc. Ser. C Appl. Stat.*, 37:344–352.
- 722 Malécot, G. (1948). *Les Mathématiques de l’Hérédité*. Masson, Paris.
- 723 Nei, M. (1973). Analysis of gene diversity in subdivided populations. *Proc.*
724 *Natl. Acad. Sci. USA*, 70:3321–3323.
- 725 Nei, M. (1977). F -statistics and analysis of gene diversity in subdivided
726 populations. *Ann. Hum. Genet.*, 41:225–233.
- 727 Nei, M. (1978). Estimation of average heterozygosity and genetic distance
728 from a small number of individuals. *Genetics*, 89:583–590.
- 729 Nei, M. (1986). Definition and estimation of fixation indices. *Evolution*,
730 40:643–645.
- 731 Nei, M. and Chesser, R. K. (1983). Estimation of fixation indices and gene
732 diversities. *Ann. Hum. Genet.*, 47:253–259.

- 733 Nychka, D., Furrer, R., Paige, J., and Sain, S. (2017). *fields*: Tools for spatial
734 data. R package version 9.6.
- 735 Orgogozo, V., Peluffo, A. E., and Morizot, B. (2016). The “mendelian gene”
736 and the “molecular gene”: two relevant concepts of genetic units. In Or-
737 gogozo, V., editor, *Genes and Evolution*, volume 119 of *Current Topics in*
738 *Developmental Biology*, pages 1–26. Academic Press.
- 739 Pickrell, J. K. and Pritchard, J. K. (2012). Inference of population splits
740 and mixtures from genome-wide allele frequency data. *PLoS Genet.*,
741 8(11):e1002967.
- 742 R Core Team (2017). *R: A Language and Environment for Statistical Com-*
743 *puting*. R Foundation for Statistical Computing, Vienna, Austria.
- 744 Reynolds, J., Weir, B. S., and Cockerham, C. C. (1983). Estimation of the
745 coancestry coefficient: basis for a short-term genetic distance. *Genetics*,
746 105:767–779.
- 747 Ridout, M. S., Demktrio, C. G. B., and Firth, D. (1999). Estimating intra-
748 class correlation for binary data. *Biometrics*, 55:137–148.
- 749 Robertson, A. (1962). Weighting in the estimation of variance components
750 in the unbalanced single classification. *Biometrics*, 18:413–417.
- 751 Rode, N. O., Holtz, Y., Loridon, K., Santoni, S., Ronfort, J., and Gay, J.
752 (2018). How to optimize the precision of allele and haplotype frequency
753 estimates using pooled-sequencing data. *Mol. Ecol. Resour.*, 18:194–203.
- 754 Ross, M. G., Russ, C., Costello, M., Hollinger, A., Lennon, N. J., Hegarty,

- 755 R., Nusbaum, C., and Jaffe, D. B. (2013). Characterizing and measuring
756 bias in sequence data. *Genome Biol.*, 14:R51.
- 757 Rousset, F. (1996). Equilibrium values of measures of population subdivision
758 for stepwise mutation processes. *Genetics*, 142:1357–1362.
- 759 Rousset, F. (1997). Genetic differentiation and estimation of gene flow from
760 F -statistics under isolation by distance. *Genetics*, 145:1219–1228.
- 761 Rousset, F. (2007). Inferences from spatial population genetics. In Bald-
762 ing, D. J., Bishop, M., and Cannings, C., editors, *Handbook of Statistical*
763 *Genetics*, pages 945–979, Chichester. John Wiley & Sons, Ltd.
- 764 Rousset, F. (2008). genepop'007: a complete re-implementation of the
765 genepop software for Windows and Linux. *Mol. Ecol. Resour.*, 8:103–106.
- 766 Schlötterer, C., Tobler, R., Kofler, R., and Nolte, V. (2014). Sequencing
767 pools of individuals – mining genome-wide polymorphism data without
768 big funding. *Nat. Rev. Genet.*, 15:749–763.
- 769 Slatkin, M. (1993). Isolation by distance in equilibrium and non-equilibrium
770 populations. *Evolution*, 47:264–279.
- 771 Smadja, C. M., Canbäck, B., Vitalis, R., Gautier, M., Ferrari, J., Zhou, J.-J.,
772 and Butlin, R. K. (2012). Large-scale candidate gene scan reveals the role
773 of chemoreceptor genes in host plant specialization and speciation in the
774 pea aphid. *Evolution*, 66:2723–2738.
- 775 The International HapMap Consortium (2005). A haplotype map of the
776 human genome. *Nature*, 437:1299–1320.

- 777 Tukey, J. W. (1957). Variances of variance components: II. The unbalanced
778 single classification. *Ann. Math. Statist.*, 28:43–56.
- 779 Vitalis, R. (2012). DETSEL: An R-Package to detect marker loci responding
780 to selection. In Pompanon, F. and Bonin, A., editors, *Data Production and*
781 *Analysis in Population Genomics: Methods and Protocols*, volume 888 of
782 *Methods in Molecular Biology*, pages 277–293, New York. Humana Press.
- 783 Vitalis, R., Boursot, P., and Dawson, K. (2001). Interpretation of variation
784 across marker loci as evidence of selection. *Genetics*, 158:1811–1823.
- 785 Wahlund, S. (1928). Zusammensetzung von populationen und korrelation-
786 serscheinungen vom standpunkt der vererbungslehre aus betrachtet. *Hered-*
787 *itas*, 11:65–106.
- 788 Weir, B. S. (1996). *Genetic Data Analysis II*. Sinauer Associates, Inc.,
789 Sunderland, MA.
- 790 Weir, B. S. (2012). Estimating F -statistics: A historical view. *Philos. Sci.*,
791 79:637–643.
- 792 Weir, B. S., Cardon, L. R., Anderson, A. D., Nielsen, D. M., and Hill, W. G.
793 (2005). Measures of human population structure show heterogeneity among
794 genomic regions. *Genome Res.*, 15:1468–1476.
- 795 Weir, B. S. and Cockerham, C. C. (1984). Estimating F -statistics for the
796 analysis of population structure. *Evolution*, 38:1358–1370.
- 797 Weir, B. S. and Goudet, J. (2017). An unified characterization of population
798 structure and relatedness. *Genetics*, 206:2085–2103.

- 799 Weir, B. S. and Hill, W. G. (2002). Estimating F -statistics. *Annu. Rev.*
800 *Genet.*, 36:721–750.
- 801 Whitlock, M. C. and Lotterhos, K. E. (2015). Reliable detection of loci re-
802 sponsible for local adaptation: inference of a null model through trimming
803 the distribution of F_{ST} . *Am. Nat.*, 186:S24–S36.
- 804 Wright, S. (1931). Evolution in Mendelian populations. *Genetics*, 16:97–159.
- 805 Wright, S. (1951). The genetical structure of populations. *Ann. Eugen.*,
806 15:323–354.
- 807 Wu, S., Crespi, C. M., and Wong, W. K. (2012). Comparison of methods
808 for estimating the intraclass correlation coefficient for binary responses in
809 cancer prevention cluster randomized trials. *Contemp. Clin. Trials*, 33:869–
810 880.

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
$\hat{\pi}_k \equiv X_{\dots:k}$	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^{pool} (resp. \hat{Q}_2^{pool})	Unbiased estimator of the IIS probability for genes sampled within (resp. between) populations

Table 2 Definition of the F_{ST} estimators used in the text

Notation	Definition
\hat{F}_{ST}^{pool}	Equation 9
FRP ₁₃	Ferretti et al. (2013) and Equations 16,20–21
NC ₈₃	Nei and Chesser (1983)
PP2 _d	Kofler et al. (2011) and Equations 16–18
PP2 _a	Kofler et al. (2011) and Equation 19
WC ₈₄	Weir and Cockerham (1984)

Table 3 Overall F_{ST} estimates from multiple pools

F_{ST}	n	Pool-seq		Ind-seq
		Cov.	\hat{F}_{ST}^{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)	

Overall F_{ST} was estimated for various conditions of expected F_{ST} , pool size (n) and coverage (Cov.). For Pool-seq data, we computed our estimator \hat{F}_{ST}^{pool} (Equation 13). The mean (RMSE) over 50 independent replicates of the `ms` simulations are provided, for all populations ($n_d = 8$). For comparison, we computed WC₈₄ from allele count data inferred from individual genotypes (Ind-seq).

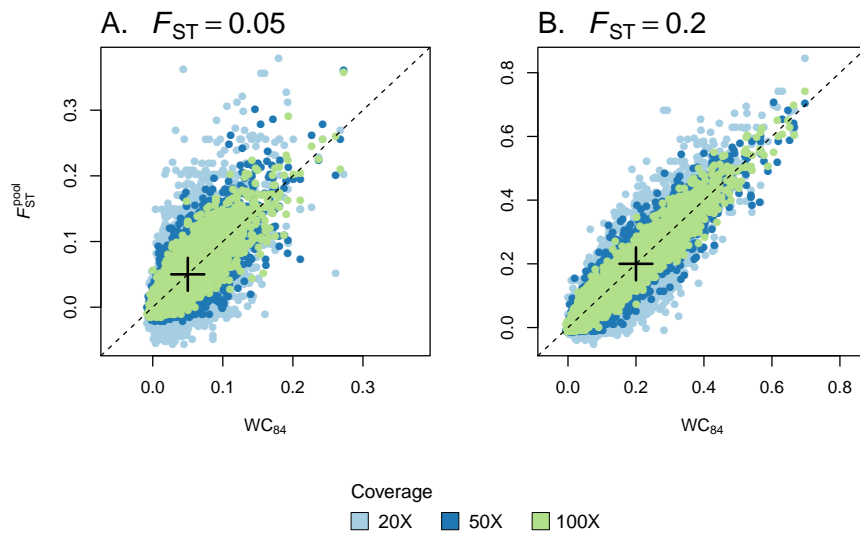


Figure 1 Single-locus estimates of F_{ST} . We compared single-locus estimates of F_{ST} based on allele count data inferred from individual genotypes (Ind-seq), using the WC_{84} estimator, to \hat{F}_{ST}^{pool} estimates from Pool-seq data. We simulated 5,000 SNPs using *ms* in an island model with $n_d = 8$ demes. We used two migration rates corresponding to $F_{ST} = 0.05$ (A) and $F_{ST} = 0.20$ (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_{ST} .

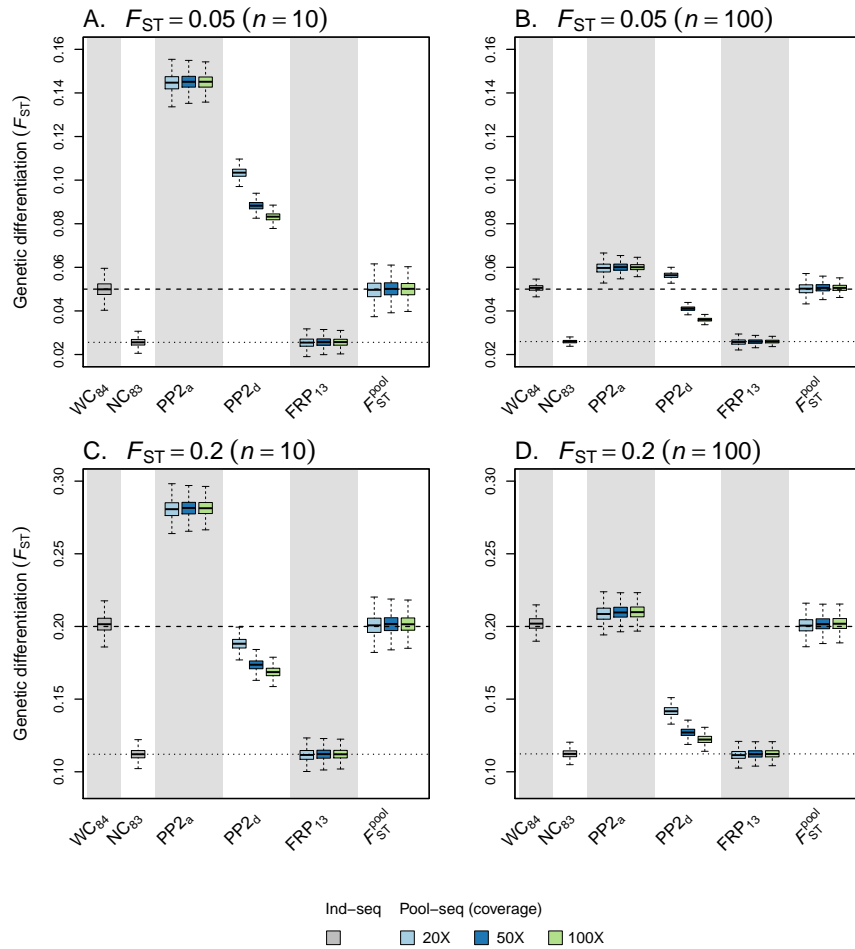


Figure 2 Precision and accuracy of pairwise estimators of F_{ST} . We considered two estimators based on allele count data inferred from individual genotypes (Ind-seq): WC_{84} and NC_{83} . For pooled data, we computed the two estimators implemented in the software package POPOOLATION2, that we refer to as $PP2_d$ and $PP2_a$, as well as the FRP_{13} estimator and our estimator \hat{F}_{ST}^{pool} (Equation 13). Each boxplot represents the distribution of multilocus F_{ST} estimates across all pairwise comparisons in an island model with $n_d = 8$ demes, and 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} and the dotted line indicates the median of the distribution of NC_{83} estimates.

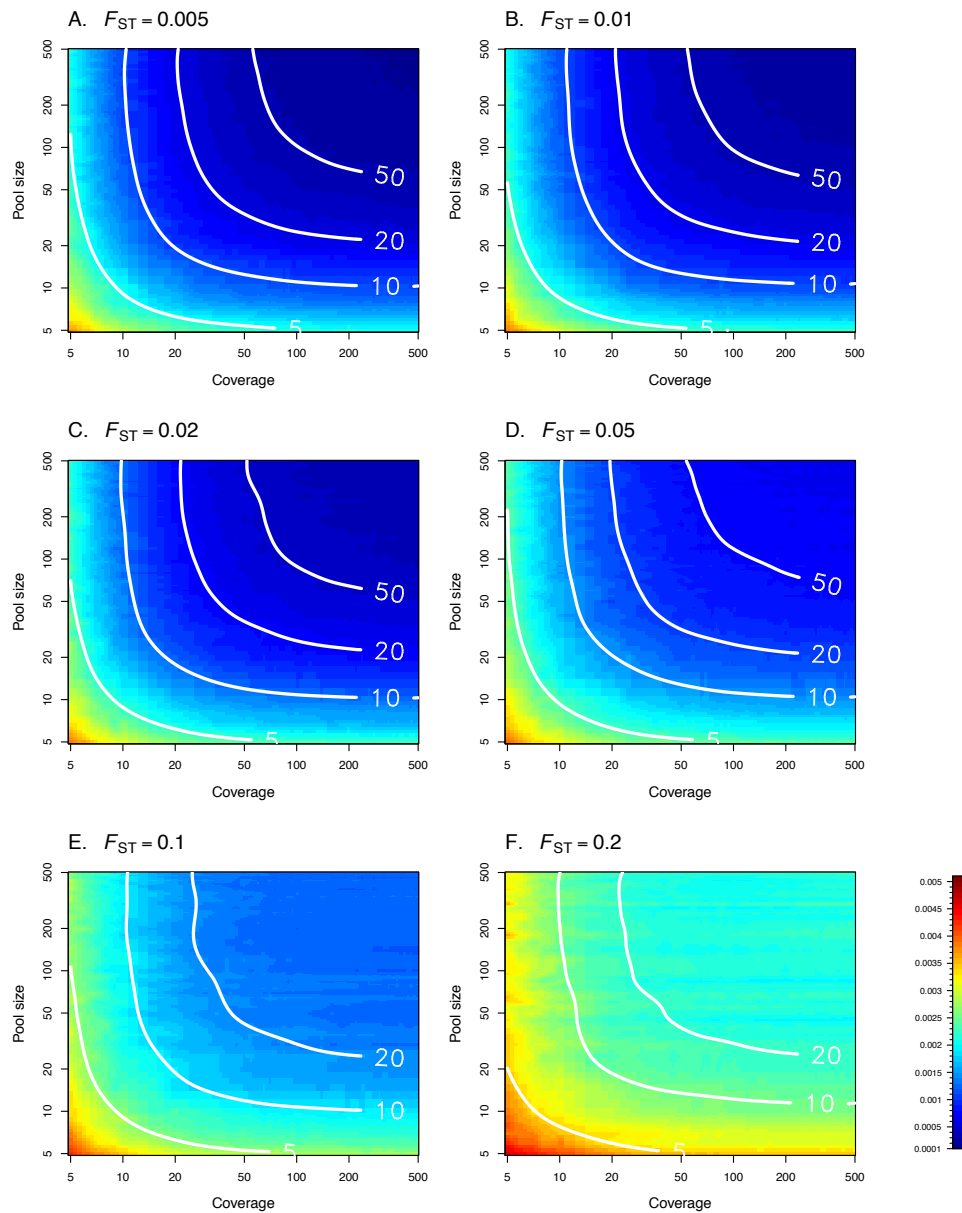


Figure 3 Root mean squared error (RMSE) of F_{ST} estimates for a wide range of pool sizes and coverage, with F_{ST} varying from 0.005 to 0.2 (A–F). Each density plot gives the RMSE of our estimator \hat{F}_{ST}^{pool} , using simple linear interpolation from a set of 44×44 pairs of pool size and coverage values. For each pool size and coverage, 500 replicates of 5,000 markers were simulated. Plain white isolines represent the RMSE of the WC_{84} estimator computed from Ind-seq data, for various sample sizes ($n = 5, 10, 20, 50$). 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).

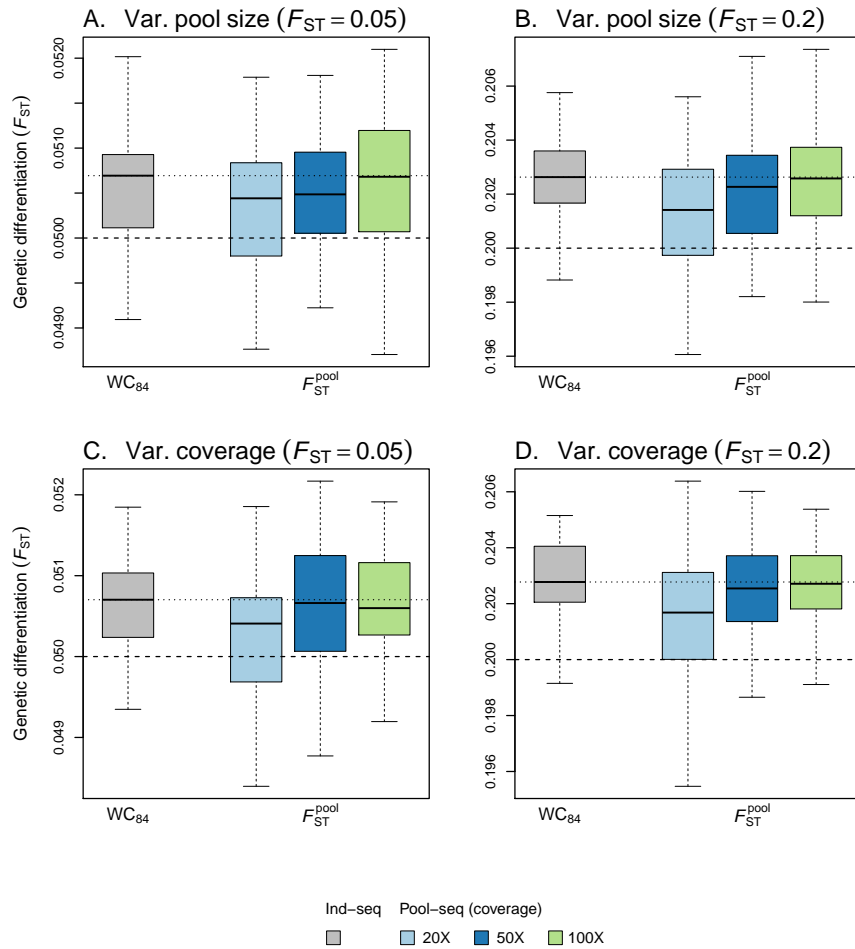


Figure 4 Precision and accuracy of F_{ST} estimates with varying pool size or varying coverage. Our estimator \hat{F}_{ST}^{pool} (Equation 13) was calculated from Pool-seq data over all loci and demes and compared to the estimator WC_{84} , computed from allele count data inferred from individual genotypes (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) or $F_{ST} = 0.20$ (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 (δ_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.

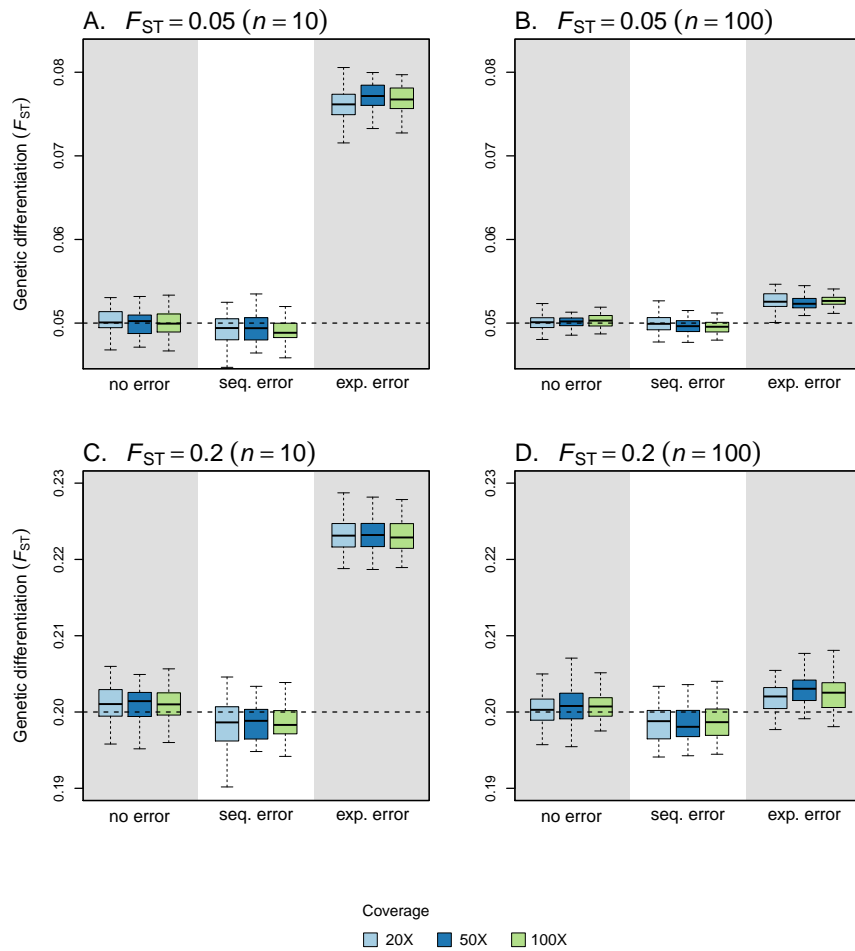


Figure 5 Precision and accuracy of F_{ST} estimates with sequencing and experimental errors. Our estimator \hat{F}_{ST}^{pool} (Equation 13) was computed from Pool-seq data over all loci and demes without error, with sequencing error (occurring at rate $\mu_e = 0.001$), and with experimental error ($\epsilon = 0.5$). 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} .

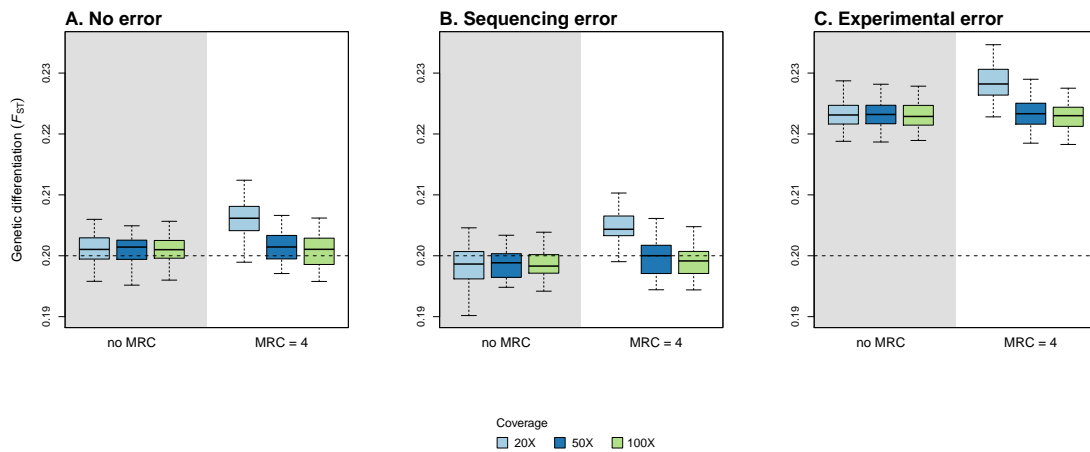


Figure 6 Precision and accuracy of F_{ST} estimates with and without filtering. Our estimator $\hat{F}_{ST}^{\text{pool}}$ (Equation 13) was computed from Pool-seq data over all loci and demes without error (A), with sequencing error (B) and with experimental error (C) (see the legend of Figure 5 for further details). For each case, we computed F_{ST} without filtering (no MRC) and with filtering (using a minimum read count $MRC = 4$). Each boxplot represents the distribution of multilocus F_{ST} estimates across 50 independent replicates of the *ms* simulations. We used a migration rate corresponding to $F_{ST} = 0.20$, 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_{ST} .

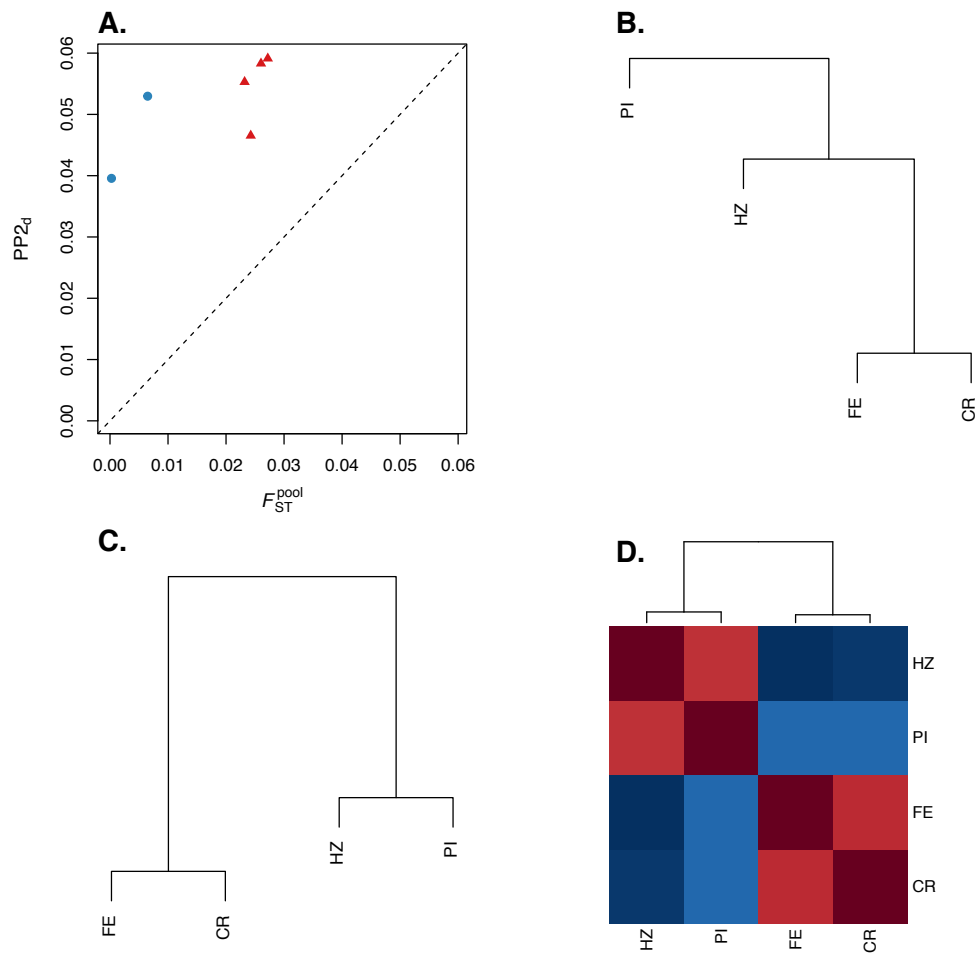


Figure 7 Analysis of the prickly sculpin (*Cottus asper*) Pool-seq data. In (A) we compare the pairwise F_{ST} estimates $PP2_d$, and \hat{F}_{ST}^{pool} (Equation 13) for all pairs of populations from the estuarine (CR and FE) and freshwater samples (PI and HZ). Within-ecotype comparisons are depicted as blue dots, and between-ecotype comparisons as red triangles. In (B–C) we show a 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.

811 SUPPLEMENTAL FILE S1: DETAILED MATHEMATICAL DERIVATIONS

812 **Analysis of variance for Pool-seq data**

813 In the following, we first derive our model for a single locus. Consider a
 814 sample of n_d subpopulations, each of which is made of n_i genes ($i = 1, \dots, n_d$)
 815 sequenced in pools (hence n_i is the haploid sample size of the i th pool). We
 816 define c_{ij} as the number of reads sequenced from gene j ($j = 1, \dots, n_i$) in
 817 subpopulation i at the locus considered. Note that c_{ij} is a latent variable,
 818 that cannot be directly observed from the data. Let $X_{ijr:k}$ be an indicator
 819 variable for read r ($r = 1, \dots, c_{ij}$) from gene j in subpopulation i , such that
 820 $X_{ijr:k} = 1$ if the r th read from the j th gene in the i th deme is of type k ,
 821 and $X_{ijr:k} = 0$ otherwise. In the following, we use standard dot notations
 822 for sample averages, i.e.: $X_{ij:k} \equiv \sum_r X_{ijr:k}/c_{ij}$, $X_{i:k} \equiv \sum_j \sum_r X_{ijr:k}/\sum_j c_{ij}$
 823 and $X_{:k} \equiv \sum_i \sum_j \sum_r X_{ijr:k}/\sum_i \sum_j c_{ij}$. The analysis of variance is based
 824 on the computation of sums of squares, as follows:

$$\begin{aligned}
 \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - X_{:k})^2 &= \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - X_{ij:k})^2 \\
 &+ \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ij:k} - X_{i:k})^2 \\
 &+ \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{i:k} - X_{:k})^2 \\
 &\equiv SSR_{:k} + SSI_{:k} + SSP_{:k} \quad (A1)
 \end{aligned}$$

825 We express the sum of squares for reads within individuals as:

$$\begin{aligned} SSR_{:k} &= \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - X_{ij:k})^2 \\ &= 0 \end{aligned} \quad (\text{A2})$$

826 since we assume that there is no sequencing error, i.e. all the reads sequenced
827 from a single gene are identical (therefore $X_{ijr:k} = X_{ij:k}$, for all r). The sum
828 of squares for genes within pools reads:

$$\begin{aligned} SSI_{:k} &= \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ij:k} - X_{i:k})^2 \\ &= \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ij:k} - \pi_k)^2 - \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{i:k} - \pi_k)^2 \\ &= \sum_i^{n_d} \sum_j^{n_i} c_{ij} (X_{ij:k} - \pi_k)^2 - \sum_i^{n_d} C_{1i} (X_{i:k} - \pi_k)^2 \end{aligned} \quad (\text{A3})$$

829 where π_k is the expectation of the frequency of allele k over independent
830 replicates of the evolutionary process, and $C_{1i} \equiv \sum_j c_{ij}$ is the total number
831 of observed reads in the i th pool. Likewise, the sum of squares for genes
832 between pools reads:

$$\begin{aligned} SSP_{:k} &= \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{i:k} - X_{...:k})^2 \\ &= \sum_i^{n_d} C_{1i} (X_{i:k} - \pi_k)^2 - C_1 (X_{...:k} - \pi_k)^2 \end{aligned} \quad (\text{A4})$$

833 where $C_1 \equiv \sum_i \sum_j c_{ij}$ is the total number of observed reads in the full sample.
834 These sums can be expressed as functions of the average frequency of reads
835 of type k for individual j : $\hat{\pi}_{ij:k} \equiv X_{ij:k}$, of the average frequency of reads of

836 type k within the i th pool: $\hat{\pi}_{i:k} \equiv X_{i\cdots:k}$, and of the average frequency of reads
 837 of type k in the full sample: $\hat{\pi}_k \equiv X_{\cdots:k}$. Note that from the definition of
 838 $X_{\cdots:k}$, $\hat{\pi}_k \equiv \sum_i \sum_j \sum_r X_{ijr:k} / \sum_i \sum_j c_{ij} = \sum_i C_{1i} \hat{\pi}_{i:k} / \sum_i C_{1i}$ is the weighted
 839 average of the sample frequencies with weights equal to the pool coverage.
 840 Our approach is therefore equivalent to the weighted analysis-of-variance in
 841 Cockerham (1973) (see also Weir and Cockerham 1984; Weir 1996; Weir and
 842 Hill 2002; Rousset 2007; Weir and Goudet 2017). Then, developing the square
 843 in the first term in the right-hand side of Equation A3, we get:

$$\begin{aligned}
 (X_{ij:k} - \pi_k)^2 &= \left(\frac{\sum_r^{c_{ij}} (X_{ijr:k} - \pi_k)}{c_{ij}} \right)^2 \\
 &= \frac{1}{c_{ij}^2} \left(\sum_r^{c_{ij}} X_{ijr:k} - c_{ij} \pi_k \right)^2 \\
 &= \frac{1}{c_{ij}^2} \left(\sum_r^{c_{ij}} X_{ijr:k}^2 + \sum_{r \neq r'}^{c_{ij}} X_{ijr:k} X_{ijr':k} - 2c_{ij}^2 X_{ij:k} \pi_k + c_{ij}^2 \pi_k^2 \right) \\
 &= \frac{1}{c_{ij}^2} (c_{ij} X_{ij:k} + c_{ij} (c_{ij} - 1) X_{ij:k} \\
 &\quad - 2c_{ij}^2 X_{ij:k} \pi_k + c_{ij}^2 \pi_k^2) \\
 &= \hat{\pi}_{ij:k} - 2\pi_k \hat{\pi}_{ij:k} + \pi_k^2
 \end{aligned} \tag{A5}$$

844 The sums of squares also depend on the unobserved frequency of pairs of
 845 genes sampled in the i th pool that are both of type k , i.e. the probability
 846 of identity in state (IIS) for allele k , for two distinct genes in the i th pool:
 847 $\hat{Q}_{1i:k} \equiv \left(\sum_{j \neq j'} \sum_{r, r'} X_{ijr:k} X_{ij'r':k} \right) / \left(C_{1i}^2 - \sum_j c_{ij}^2 \right)$. Then, developing the

848 square in the second term in the right-hand side of Equation A3, we get:

$$\begin{aligned}
 (X_{i\cdot:k} - \pi_k)^2 &= \left(\frac{\sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - \pi_k)}{C_{1i}} \right)^2 \\
 &= \frac{1}{C_{1i}^2} \left(\sum_j^{n_i} \sum_r^{c_{ij}} X_{ijr:k} - C_{1i} \pi_k \right)^2 \\
 &= \frac{1}{C_{1i}^2} \left(\sum_j^{n_i} \sum_r^{c_{ij}} X_{ijr:k}^2 + \sum_j^{n_i} \sum_{r \neq r'}^{c_{ij}} X_{ijr:k} X_{ijr':k} \right. \\
 &\quad \left. + \sum_{j \neq j'}^{n_i} \sum_{r, r'}^{c_{ij}} X_{ijr:k} X_{ij'r':k} - 2C_{1i}^2 X_{i\cdot:k} \pi_k + C_{1i}^2 \pi_k^2 \right) \\
 &= \frac{1}{C_{1i}^2} \left(\sum_j^{n_i} c_{ij} X_{ij\cdot:k} + \sum_j^{n_i} c_{ij} (c_{ij} - 1) X_{ij\cdot:k} \right. \\
 &\quad \left. + \left(C_{1i}^2 - \sum_j^{n_i} c_{ij}^2 \right) \hat{Q}_{1i:k} - 2C_{1i}^2 X_{i\cdot:k} \pi_k + C_{1i}^2 \pi_k^2 \right) \\
 &= \frac{1}{C_{1i}^2} \left(\sum_j^{n_i} c_{ij}^2 (X_{ij\cdot:k} - X_{i\cdot:k}) + \left(C_{1i}^2 - \sum_j^{n_i} c_{ij}^2 \right) (\hat{Q}_{1i:k} - X_{i\cdot:k}) \right. \\
 &\quad \left. + C_{1i}^2 X_{i\cdot:k} - 2C_{1i}^2 X_{i\cdot:k} \pi_k + C_{1i}^2 \pi_k^2 \right) \\
 &= \hat{\pi}_{i:k} - 2\pi_k \hat{\pi}_{i:k} + \pi_k^2 + \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}^2} (\hat{\pi}_{ij:k} - \hat{\pi}_{i:k}) \\
 &\quad + \left(1 - \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}^2} \right) (\hat{Q}_{1i:k} - \hat{\pi}_{i:k}) \tag{A6}
 \end{aligned}$$

849 Last, the sums of squares depend on the unobserved frequency of pairs
 850 of genes sampled in the same pool that are both of type k , i.e. the IIS
 851 probability for allele k for two distinct genes in the same pool: $\hat{Q}_{1:k} \equiv$
 852 $\left(\sum_i \sum_{j \neq j'} \sum_{r, r'} X_{ijr:k} X_{ij'r':k} \right) / \left(C_2 - \sum_i \sum_j c_{ij}^2 \right)$, and of the unobserved
 853 frequency of pairs of genes sampled in different pools that are both of type
 854 k : $\hat{Q}_{2:k} \equiv \left(\sum_{i \neq i'} \sum_{j, j'} \sum_{r, r'} X_{ijr:k} X_{i'j'r':k} \right) / (C_1^2 - C_2)$, where $C_2 \equiv \sum_i C_{1i}^2$.

855 Developing the second term in the right-hand side of Equation A4, we get:

$$\begin{aligned}
 (X_{\dots:k} - \pi_k)^2 &= \left(\frac{\sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ijr:k} - \pi_k)}{C_1} \right)^2 \\
 &= \frac{1}{C_1^2} \left(\sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} X_{ijr:k} - C_1 \pi_k \right)^2 \\
 &= \frac{1}{C_1^2} \left(\sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} X_{ijr:k}^2 + \sum_i^{n_d} \sum_j^{n_i} \sum_{r \neq r'}^{c_{ij}} X_{ijr:k} X_{ijr':k} \right. \\
 &\quad + \sum_i^{n_d} \sum_{j \neq j'}^{n_i} \sum_{r, r'}^{c_{ij}} X_{ijr:k} X_{i'j'r':k} + \sum_{i \neq i'}^{n_d} \sum_{j, j'}^{n_i} \sum_{r, r'}^{c_{ij}} X_{ijr:k} X_{i'j'r':k} \\
 &\quad \left. - 2C_1^2 X_{\dots:k} \pi_k + C_1^2 \pi_k^2 \right) \\
 &= \frac{1}{C_1^2} \left(\sum_i^{n_d} \sum_j^{n_i} c_{ij} X_{ij:k} + \sum_i^{n_d} \sum_j^{n_i} c_{ij} (c_{ij} - 1) X_{ij:k} \right. \\
 &\quad \left. + \left(C_2 - \sum_i^{n_d} \sum_j^{n_i} c_{ij}^2 \right) \hat{Q}_{1:k} + (C_1^2 - C_2) \hat{Q}_{2:k} - 2C_1^2 X_{\dots:k} \pi_k + C_1^2 \pi_k^2 \right) \\
 &= \frac{1}{C_1^2} \left(\sum_i^{n_d} \sum_j^{n_i} c_{ij}^2 (X_{ij:k} - X_{\dots:k}) + \left(C_2 - \sum_i^{n_d} \sum_j^{n_i} c_{ij}^2 \right) (\hat{Q}_{1:k} - X_{\dots:k}) \right. \\
 &\quad \left. + (C_1^2 - C_2) (\hat{Q}_{2:k} - X_{\dots:k}) + C_1^2 X_{\dots:k} - 2C_1^2 X_{\dots:k} \pi_k + C_1^2 \pi_k^2 \right) \\
 &= \hat{\pi}_k - 2\pi_k \hat{\pi}_k + \pi_k^2 + \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1^2} (\hat{\pi}_{ij:k} - \hat{\pi}_k) \\
 &\quad + \left(\frac{C_2}{C_1^2} - \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1^2} \right) (\hat{Q}_{1:k} - \hat{\pi}_k) + \left(1 - \frac{C_2}{C_1^2} \right) (\hat{Q}_{2:k} - \hat{\pi}_k) \quad (A7)
 \end{aligned}$$

856 Hence, developing the first term in the right-hand side of Equation A3 using

857 Equation A5, we have:

$$\sum_i^{n_d} \sum_j^{n_i} c_{ij} (X_{ij:k} - \pi_k)^2 = C_1 (\hat{\pi}_k - 2\pi_k \hat{\pi}_k + \pi_k^2) \quad (A8)$$

858 Likewise, developing the second term in the right-hand side of Equation A3
 859 using Equation A6, we get:

$$\begin{aligned} \sum_i^{n_d} C_{1i} (X_{i\dots k} - \pi_k)^2 &= C_1 (\hat{\pi}_k - 2\pi_k \hat{\pi}_k + \pi_k^2) + \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} (\hat{\pi}_{ij:k} - \hat{\pi}_{i:k}) \\ &+ \sum_i^{n_d} \left(C_{1i} - \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} \right) (\hat{Q}_{1i:k} - \hat{\pi}_{i:k}) \end{aligned} \quad (\text{A9})$$

860 Last, developing the second term in the right-hand side of Equation A4 using
 861 Equation A7, we get:

$$\begin{aligned} C_1 (X_{\dots k} - \pi_k)^2 &= C_1 (\hat{\pi}_k - 2\pi_k \hat{\pi}_k + \pi_k^2) + \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1} (\hat{\pi}_{ij:k} - \hat{\pi}_k) \\ &+ \left(\frac{C_2}{C_1} - \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1} \right) (\hat{Q}_{1:k} - \hat{\pi}_k) \\ &+ \left(C_1 - \frac{C_2}{C_1} \right) (\hat{Q}_{2:k} - \hat{\pi}_k) \end{aligned} \quad (\text{A10})$$

862 Then, from Equations A3, A8 and A9:

$$\begin{aligned} SSI_{:k} &= \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} (\hat{\pi}_{i:k} - \hat{\pi}_{ij:k}) \\ &+ \sum_i^{n_d} \left(C_{1i} - \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} \right) (\hat{\pi}_{i:k} - \hat{Q}_{1i:k}) \end{aligned} \quad (\text{A11})$$

863 and from Equations A4, A9 and A10:

$$\begin{aligned} SSP_{:k} &= \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} (\hat{\pi}_{ij:k} - \hat{\pi}_{i:k}) - \sum_i^{n_d} \left(C_{1i} - \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} \right) (\hat{\pi}_{i:k} - \hat{Q}_{1i:k}) \\ &+ \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1} (\hat{\pi}_k - \hat{\pi}_{ij:k}) + \left(\frac{C_2}{C_1} - \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1} \right) (\hat{\pi}_k - \hat{Q}_{1:k}) \\ &+ \left(C_1 - \frac{C_2}{C_1} \right) (\hat{\pi}_k - \hat{Q}_{2:k}) \end{aligned} \quad (\text{A12})$$

864 Taking expectation over all possible samples from all replicate populations
 865 sharing the same evolutionary history, we get from Equation A11:

$$\begin{aligned}
 \mathbb{E}(SSI_{:k}) &= \sum_i^{n_d} \sum_j^{n_i} \mathbb{E}(\hat{\pi}_{i:k} - \hat{\pi}_{ij:k}) \mathbb{E}\left(\frac{c_{ij}^2}{C_{1i}}\right) \\
 &+ \sum_i^{n_d} \mathbb{E}(\hat{\pi}_{i:k} - \hat{Q}_{1i:k}) \mathbb{E}\left(C_{1i} - \sum_j \frac{c_{ij}^2}{C_{1i}}\right) \\
 &= (\pi_k - Q_{1:k}) \left(C_1 - \mathbb{E}\left(\sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}}\right) \right) \tag{A13}
 \end{aligned}$$

866 where $Q_{1:k}$ is the expected IIS probability that two genes in the same pool
 867 are both of type k . Likewise, from Equation A12:

$$\begin{aligned}
 \mathbb{E}(SSP_{:k}) &= \sum_i^{n_d} \sum_j^{n_i} \mathbb{E}(\hat{\pi}_{i:k} - \hat{\pi}_{ij:k}) \mathbb{E}\left(\frac{c_{ij}^2}{C_{1i}}\right) + \sum_i^{n_d} \sum_j^{n_i} \mathbb{E}(\hat{\pi}_k - \hat{\pi}_{ij:k}) \mathbb{E}\left(\frac{c_{ij}^2}{C_1}\right) \\
 &- \sum_i^{n_d} \mathbb{E}(\hat{\pi}_{i:k} - \hat{Q}_{1i:k}) \mathbb{E}\left(C_{1i} - \sum_j \frac{c_{ij}^2}{C_{1i}}\right) \\
 &+ \mathbb{E}(\hat{\pi}_k - \hat{Q}_{1:k}) \mathbb{E}\left(\frac{C_2}{C_1} - \sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1}\right) \\
 &+ \left(C_1 - \frac{C_2}{C_1}\right) \mathbb{E}(\hat{\pi}_k - \hat{Q}_{2:k}) \\
 &= (\pi_k - Q_{1:k}) \left(\frac{C_2}{C_1} - \mathbb{E}\left(\sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1}\right) \right) \\
 &- (\pi_k - Q_{1:k}) \left(C_1 - \mathbb{E}\left(\sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}}\right) \right) \\
 &+ \left(C_1 - \frac{C_2}{C_1}\right) (\pi_k - Q_{2:k}) \tag{A14}
 \end{aligned}$$

868 where $Q_{2:k}$ is the expected IIS probability that two genes from different pools
 869 are both of type k . Note that the expected sums $\mathbb{E}\left(\sum_i \sum_j c_{ij}^2\right)/C_{1i}$ and
 870 $\mathbb{E}\left(\sum_i \sum_j c_{ij}^2\right)/C_1$ in Equations A13 and A14 depend on the latent variable

871 c_{ij} , that cannot be directly observed from the data. Therefore, we must make
 872 an assumption on the distribution of the c_{ij} 's to proceed. In the following,
 873 we assume that for each pool i , c_{ij} follows a multinomial distribution with
 874 parameter C_{1i} (the number of trials, i.e. the total number of reads in the
 875 i th pool) and probabilities $(1/n_i, \dots, 1/n_i)$ for the n_i individuals in the pool.
 876 Then:

$$\begin{aligned}
 \mathbb{E} \left(\sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_{1i}} \right) &= \sum_i^{n_d} \frac{1}{C_{1i}} \sum_j^{n_i} \mathbb{E} (c_{ij}^2) \\
 &= \sum_i^{n_d} \frac{1}{C_{1i}} \sum_j^{n_i} \left(\mathbb{E} (c_{ij})^2 + \mathbb{V} (c_{ij}) \right) \\
 &= \sum_i^{n_d} \frac{1}{C_{1i}} \sum_j^{n_i} \left(\left(\frac{C_{1i}}{n_i} \right)^2 + \frac{C_{1i}}{n_i} \left(\frac{n_i - 1}{n_i} \right) \right) \\
 &= \sum_i^{n_d} \left(\frac{C_{1i}}{n_i} + \left(\frac{n_i - 1}{n_i} \right) \right) \equiv D_2 \tag{A15}
 \end{aligned}$$

877 and:

$$\begin{aligned}
 \mathbb{E} \left(\sum_i^{n_d} \sum_j^{n_i} \frac{c_{ij}^2}{C_1} \right) &= \frac{1}{C_1} \sum_i^{n_d} \sum_j^{n_i} \mathbb{E} (c_{ij}^2) \\
 &= \frac{1}{C_1} \sum_i^{n_d} C_{1i} \left[\frac{C_{1i}}{n_i} + \left(\frac{n_i - 1}{n_i} \right) \right] \equiv D_2^* \tag{A16}
 \end{aligned}$$

878 Hence, from Equations A13 and A15, we have:

$$\mathbb{E}(SSI_{:k}) = (C_1 - D_2) (\pi_k - Q_{1:k}) \tag{A17}$$

879 and from Equations A14 and A16:

$$\begin{aligned}
 \mathbb{E}(SSP_{:k}) &= \left(\frac{C_2}{C_1} - D_2^* \right) (\pi_k - Q_{1:k}) - (C_1 - D_2) (\pi_k - Q_{1:k}) \\
 &+ \left(C_1 - \frac{C_2}{C_1} \right) (\pi_k - Q_{2:k}) \\
 &= \left(C_1 - \frac{C_2}{C_1} \right) (Q_{1:k} - Q_{2:k}) \\
 &+ (D_2 - D_2^*) (\pi_k - Q_{1:k})
 \end{aligned} \tag{A18}$$

880 Summing over alleles, we get the following expressions for the expected sums
 881 of squares for genes between individuals within pools:

$$\mathbb{E}(SSI) = \sum_k \mathbb{E}(SSI_{:k}) = (C_1 - D_2) (1 - Q_1) \tag{A19}$$

882 and for genes between individuals from different pools:

$$\begin{aligned}
 \mathbb{E}(SSP) &= \sum_k \mathbb{E}(SSP_{:k}) \\
 &= \left(C_1 - \frac{C_2}{C_1} \right) (Q_1 - Q_2) + (D_2 - D_2^*) (1 - Q_1)
 \end{aligned} \tag{A20}$$

883 Rearranging Equations A19–A20, we get:

$$Q_1 - Q_2 = \frac{(C_1 - D_2) \mathbb{E}(SSP) - (D_2 - D_2^*) \mathbb{E}(SSI)}{(C_1 - D_2) (C_1 - C_2/C_1)} \tag{A21}$$

884 and:

$$1 - Q_2 = \frac{(C_1 - D_2) \mathbb{E}(SSP) + (n_c - 1) (D_2 - D_2^*) \mathbb{E}(SSI)}{(C_1 - D_2) (C_1 - C_2/C_1)} \tag{A22}$$

885 where $n_c \equiv (C_1 - C_2/C_1) / (D_2 - D_2^*)$. Let $MSI \equiv SSI / (C_1 - D_2)$ and
 886 $MSP \equiv SSP / (D_2 - D_2^*)$. Then, rearranging Equations A21–A22, we get:

$$F_{ST} \equiv \frac{Q_1 - Q_2}{1 - Q_2} = \frac{\mathbb{E}(MSP) - \mathbb{E}(MSI)}{\mathbb{E}(MSP) + (n_c - 1) \mathbb{E}(MSI)} \quad (\text{A23})$$

887 which yields the method-of-moments estimator:

$$\hat{F}_{ST}^{\text{pool}} = \frac{MSP - MSI}{MSP + (n_c - 1) MSI} \quad (\text{A24})$$

888 Since SSI (Equation A3) and SSP (Equation A4) may be rewritten in terms
 889 of sample frequencies as:

$$\begin{aligned} SSI &= \sum_k SSI_{:k} = \sum_k \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{ij:k} - X_{i:k})^2 \\ &= \sum_k \sum_i^{n_d} C_{1i} \hat{\pi}_{i:k} (1 - \hat{\pi}_{i:k}) \end{aligned} \quad (\text{A25})$$

890 and:

$$\begin{aligned} SSP &= \sum_k SSP_k = \sum_k \sum_i^{n_d} \sum_j^{n_i} \sum_r^{c_{ij}} (X_{i:k} - X_{\dots:k})^2 \\ &= \sum_k \sum_i^{n_d} C_{1i} (\hat{\pi}_{i:k} - \hat{\pi}_k)^2 \end{aligned} \quad (\text{A26})$$

891 our estimator then takes the form:

$$\hat{F}_{ST}^{\text{pool}} = \frac{\sum_k [(C_1 - D_2) \sum_i^{n_d} C_{1i} (\hat{\pi}_{i:k} - \hat{\pi}_k)^2 - (D_2 - D_2^*) \sum_i^{n_d} C_{1i} \hat{\pi}_{i:k} (1 - \hat{\pi}_{i:k})]}{\sum_k [(C_1 - D_2) \sum_i^{n_d} C_{1i} (\hat{\pi}_{i:k} - \hat{\pi}_k)^2 + (n_c - 1) (D_2 - D_2^*) \sum_i^{n_d} C_{1i} \hat{\pi}_{i:k} (1 - \hat{\pi}_{i:k})]} \quad (\text{A27})$$

892 The estimator in Equation A24 can also be expressed as a function of the

893 frequencies of identical pairs of genes $\hat{Q}_1 = \sum_k \hat{Q}_{1:k}$ and $\hat{Q}_2 = \sum_k \hat{Q}_{2:k}$, as:

$$\hat{F}_{\text{ST}}^{\text{pool}} = \frac{\left(\hat{Q}_1 - \hat{Q}_2\right) \alpha + \left(C_1 - \sum_i \sum_j \frac{c_{ij}^2}{C_1}\right) \beta}{\left(1 - \hat{Q}_2\right) \alpha + \left(C_2/C_1 - \sum_i \sum_j \frac{c_{ij}^2}{C_1}\right) \beta} \quad (\text{A28})$$

894 where:

$$\alpha \equiv \left(C_1 - \sum_i \sum_j \frac{c_{ij}^2}{C_{1i}}\right) \left(C_1 - \frac{C_2}{C_1}\right) \quad (\text{A29})$$

895 and:

$$\beta \equiv \sum_i \left(C_{1i} - \sum_j \frac{c_{ij}^2}{C_{1i}}\right) \left(\hat{Q}_{1i} - \hat{Q}_1\right) \quad (\text{A30})$$

896 If we take the limit case where the number of sequenced reads per gene is
 897 constant, i.e. if $C_{1i} = C$, for all $i \in (1, \dots, n_d)$, then it can be shown that
 898 Equation A28 reduces exactly to Equations 28A29–28A30 in Rousset (2007),
 899 p. 977. Furthermore, if the pools have all the same size, i.e. if $n_i = n$ for all
 900 $i \in (1, \dots, n_d)$, then $\hat{F}_{\text{ST}}^{\text{pool}} = \left(\hat{Q}_1 - \hat{Q}_2\right) / \left(1 - \hat{Q}_2\right)$.

901 If the pools have all the same size and if the number of reads per pool is
 902 constant, then one can also show that Equations A25–A26 reduce to:

$$SSI = n_d(C - 1) \left(1 - \hat{Q}_1^r\right) \quad (\text{A31})$$

903 and:

$$SSP = C(n_d - 1) \left(1 - \hat{Q}_2^r\right) - (n_d - 1)(C - 1) \left(1 - \hat{Q}_1^r\right) \quad (\text{A32})$$

904 where \hat{Q}_1^r and \hat{Q}_2^r are the frequencies of identical pairs of reads within and be-
 905 tween pools, respectively, computed by simple counting of IIS pairs. These

906 are (unweighted) averages of the population-specific estimates \hat{Q}_{1i}^r (Equa-
907 tion A34) and the pairwise estimates $\hat{Q}_{2ii'}^r$ (Equation A40), respectively.
908 Then, from Equation A24, we get:

$$\hat{F}_{ST}^{\text{pool}} = 1 - \left(\frac{1 - \hat{Q}_1^r}{1 - \hat{Q}_2^r} \right) \left(\frac{n}{n-1} \right) \quad (\text{A33})$$

909 **IIS probabilities for Pool-seq data**

910 In this Appendix, we provide unbiased estimates of IIS probabilities between
 911 pairs of genes, computed from read count data. Let $r_{i:k} = \sum_j \sum_r X_{ijr:k}$ be
 912 the number of reads of type k in the i th pool. A straightforward estimate of
 913 the IIS probability between pairs of reads in the i th pool is given by:

$$\hat{Q}_{1i}^r \equiv \frac{\sum_k r_{i:k} (r_{i:k} - 1)}{C_{1i} (C_{1i} - 1)} \quad (\text{A34})$$

914 where $C_{1i} = \sum_k r_{i:k}$. As above (see Equations A15 and A16), we assume that
 915 in each pool, the conditional distribution of the read counts $r_{i:k}$, given the
 916 (unobserved) allele counts $y_{i:k}$, is binomial, i.e.: $r_{i:k} \mid y_{i:k} \sim \text{Bin}(y_{i:k}/n_i, C_{1i})$.
 917 The conditional expectation of the number of reads is therefore given by:
 918 $\mathbb{E}(r_{i:k} \mid y_{i:k}) = C_{1i} (y_{i:k}/n_i)$, and the conditional expectation of the squared
 919 number of reads by: $\mathbb{E}(r_{i:k}^2 \mid y_{i:k}) = C_{1i} (C_{1i} - 1) (y_{i:k}/n_i)^2 + C_{1i} (y_{i:k}/n_i)$.
 920 Therefore, the conditional expectation of the IIS probability between pairs
 921 of reads in the i th pool reads:

$$\mathbb{E}(\hat{Q}_{1i}^r \mid y_{i:k}) = \frac{\sum_k \mathbb{E}(r_{i:k}^2 - r_{i:k})}{C_{1i} (C_{1i} - 1)} = \sum_k \left(\frac{y_{i:k}}{n_i} \right)^2 \quad (\text{A35})$$

922 Since

$$\hat{Q}_{1i} \equiv \frac{\sum_k y_{i:k} (y_{i:k} - 1)}{n_i (n_i - 1)} \quad (\text{A36})$$

923 is an unbiased estimate of the IIS probability between pairs of distinct genes
 924 in the i th pool, Equation A35 implies that \hat{Q}_{1i}^r (Equation A34) is a biased
 925 estimate of that quantity (i.e., the IIS probability between pairs of reads
 926 within a pool is a biased estimate of the IIS probability between pairs of

927 distinct genes in that pool). This is so, because the former confounds pairs
 928 of reads that are identical because they were sequenced from a single gene
 929 copy, from pairs of reads (from distinct gene copies) that are identical because
 930 they share a common ancestor. However, inspection of Equation A35 suggests
 931 that an unbiased estimate of \hat{Q}_{1i} may be given by:

$$\hat{Q}_{1i}^{\text{pool}} \equiv 1 - \frac{n_i}{n_i - 1} (1 - \hat{Q}_{1i}^r) \quad (\text{A37})$$

932 Taking expectation of Equation A37, we get indeed:

$$\begin{aligned} \mathbb{E} \left(\hat{Q}_{1i}^{\text{pool}} \mid y_{i:k} \right) &= \frac{n_i}{n_i - 1} \mathbb{E} \left(\hat{Q}_{1i}^r \right) - \frac{1}{n_i - 1} \\ &= \frac{n_i}{n_i - 1} \sum_k \left(\frac{y_{i:k}}{n_i} \right)^2 - \frac{n_i}{n_i(n_i - 1)} \\ &= \frac{\sum_k y_{i:k}^2}{n_i(n_i - 1)} - \frac{\sum_k y_{i:k}}{n_i(n_i - 1)} \\ &= \frac{\sum_k y_{i:k}(y_{i:k} - 1)}{n_i(n_i - 1)} \equiv \hat{Q}_{1i} \end{aligned} \quad (\text{A38})$$

933 Following Weir and Goudet (2017), we define the overall IIS probability be-
 934 tween pairs of genes within pools as the unweighted average of population-
 935 specific estimates, leading to:

$$\hat{Q}_1^{\text{pool}} \equiv \frac{\sum_i \hat{Q}_{1i}^{\text{pool}}}{n_d} \quad (\text{A39})$$

936 A straightforward estimate of the IIS probability between pairs of reads
 937 taken in different pools i and i' is given by:

$$\hat{Q}_{2ii'}^r \equiv \frac{\sum_k r_{i:k} r_{i':k}}{C_{1i} C_{1i'}} \quad (\text{A40})$$

938 Since we assume that pools are conditionally independent, taking expectation

939 gives:

$$\begin{aligned}\mathbb{E}\left(\hat{Q}_{2ii'}^r \mid y_{i:k}, y_{i':k}\right) &= \frac{\sum_k \mathbb{E}(r_{i:k})\mathbb{E}(r_{i':k})}{C_{1i}C_{1i'}} \\ &= \sum_k \left(\frac{y_{i:k}y_{i':k}}{n_i n_{i'}}\right) \equiv \hat{Q}_{2ii'}\end{aligned}\quad (\text{A41})$$

940 Therefore, the IIS probability between pairs of reads sampled in different
941 pools is an unbiased estimate of the IIS probability between pairs of genes in
942 these pools, and an unbiased estimate of the IIS probability of genes sampled
943 from different pools is given by:

$$\hat{Q}_{2ii'}^{\text{pool}} \equiv \hat{Q}_{2ii'}^r \quad (\text{A42})$$

944 As above, we define the overall IIS probability between pairs of genes sampled
945 from different pools as the unweighted average of pairwise estimates, i.e.:

$$\hat{Q}_2^{\text{pool}} \equiv \frac{\sum_{i \neq i'} \hat{Q}_{2ii'}^{\text{pool}}}{n_d(n_d - 1)} \quad (\text{A43})$$

946 We can then derive an IIS-based estimator of F_{ST} , as:

$$\begin{aligned}\hat{F}_{\text{ST}}^{\text{pool-PID}} &\equiv \frac{\hat{Q}_1^{\text{pool}} - \hat{Q}_2^{\text{pool}}}{1 - \hat{Q}_2^{\text{pool}}} = 1 - \frac{1 - \hat{Q}_1^{\text{pool}}}{1 - \hat{Q}_2^{\text{pool}}} \\ &= 1 - \frac{\sum_i \left[\left(1 - \hat{Q}_{1i}^r\right) n_i / (n_i - 1) \right]}{\sum_{i \neq i'} \left(1 - \hat{Q}_{2ii'}^r\right) / (n_d - 1)}\end{aligned}\quad (\text{A44})$$

947 which, to the extent that we may take the expectation of a ratio to be the
948 ratio of expectations, is unbiased. If the pools have all the same size (i.e., if

949 $n_i = n$ for all i), then Equation A44 reduces to:

$$\hat{F}_{\text{ST}}^{\text{pool-PID}} = 1 - \left(\frac{1 - \hat{Q}_1^r}{1 - \hat{Q}_2^r} \right) \left(\frac{n}{n-1} \right) \quad (\text{A45})$$

950 where $\hat{Q}_1^r \equiv \sum_i \hat{Q}_{1i}^r / n_d$ and $\hat{Q}_2^r \equiv \sum_{i \neq i'} \hat{Q}_{2ii'}^r / [n_d(n_d - 1)]$. Note that Equa-
951 tion A45 is strictly identical to Equation A33. Therefore, if the pools have all
952 the same size and if the number of reads per pool is constant, the analysis-
953 of-variance estimator $\hat{F}_{\text{ST}}^{\text{pool}}$ is strictly equivalent to the estimator $\hat{F}_{\text{ST}}^{\text{pool-PID}}$
954 based on the computation of IIS probabilities between pairs of reads, with
955 appropriate bias correction (see Equation A37). This echoes the derivations
956 by Rousset (2007) for Ind-seq data, who showed that the analysis-of-variance
957 approach (Weir and Cockerham 1984) and the simple strategy of estimat-
958 ing IIS probabilities by counting identical pairs of genes provides identical
959 estimates when sample sizes are equal (see also Cockerham and Weir 1987;
960 Karlsson et al. 2007).

961 Alternatively, the overall IIS probability between pairs of genes within
962 pools may be defined as the weighted average of population-specific estimates,
963 with weights equal to the number of pairs of genes in each pool (see Rousset
964 2007), i.e.:

$$\tilde{Q}_1^{\text{pool}} \equiv \frac{\sum_i n_i(n_i - 1) \hat{Q}_{1i}^{\text{pool}}}{\sum_i n_i(n_i - 1)} \quad (\text{A46})$$

965 Likewise, the overall IIS probability between pairs of genes sampled from
966 different pools may be defined as the weighted average of pairwise estimates,
967 with weights equal to the number of pairs of genes sampled between pools,

968 i.e.:

$$\tilde{Q}_2^{\text{pool}} \equiv \frac{\sum_{i \neq i'} n_i n'_i \hat{Q}_{2ii'}^{\text{pool}}}{\sum_{i \neq i'} n_i n'_i} \quad (\text{A47})$$

969 We can then derive an IIS-based estimator of F_{ST} , using weighted IIS prob-
 970 abilities, as:

$$\begin{aligned} \tilde{F}_{\text{ST}}^{\text{pool-PID}} &\equiv \frac{\tilde{Q}_1^{\text{pool}} - \tilde{Q}_2^{\text{pool}}}{1 - \tilde{Q}_2^{\text{pool}}} = 1 - \frac{1 - \tilde{Q}_1^{\text{pool}}}{1 - \tilde{Q}_2^{\text{pool}}} \\ &= 1 - \frac{\sum_i \left[n_i^2 \left(1 - \hat{Q}_{1i}^r \right) \right] / \sum_i n_i (n_i - 1)}{\sum_{i \neq i'} n_i n'_i \left(1 - \hat{Q}_{2ii'}^r \right) / \sum_{i \neq i'} n_i n'_i} \end{aligned} \quad (\text{A48})$$

971 If the pools have all the same size (i.e., if $n_i = n$ for all i), then Equation A48
 972 reduces to Equation A45, and $\tilde{F}_{\text{ST}}^{\text{pool-PID}} = \hat{F}_{\text{ST}}^{\text{pool-PID}}$. With unbalanced sam-
 973 ples, simulation analyses show that $\tilde{F}_{\text{ST}}^{\text{pool-PID}}$ has larger bias and variance
 974 than $\hat{F}_{\text{ST}}^{\text{pool-PID}}$, in particular for low levels of differentiation (see Figure S4).

Table S1 Comparison of pairwise F_{ST} estimates

F_{ST}	n	Pool-seq		Ind-seq
		Cov.	\hat{F}_{ST}^{pool}	WC ₈₄
0.05	10	20×	0.051 (0.004)	
0.05	10	50×	0.051 (0.004)	0.051 (0.003)
0.05	10	100×	0.051 (0.003)	
0.05	100	20×	0.051 (0.003)	
0.05	100	50×	0.051 (0.003)	0.051 (0.002)
0.05	100	100×	0.051 (0.002)	
0.20	10	20×	0.203 (0.007)	
0.20	10	50×	0.202 (0.006)	0.202 (0.007)
0.20	10	100×	0.201 (0.006)	
0.20	100	20×	0.201 (0.006)	
0.20	100	50×	0.201 (0.006)	0.201 (0.005)
0.20	100	100×	0.202 (0.005)	

F_{ST} was estimated for various conditions of expected F_{ST} , pool size (n) and coverage (Cov.). For Pool-seq data, we computed our estimator \hat{F}_{ST}^{pool} (Equation 13). The mean (RMSE) over 50 independent replicates of the `ms` simulations are provided for a single pair of populations. For comparison, we computed WC₈₄ from allele count data inferred from individual genotypes (Ind-seq).

Table S2 Effect of unequal sampling on pairwise F_{ST} estimates

F_{ST}	n	Pool-seq		Ind-seq
		Cov.	\hat{F}_{ST}^{pool}	WC ₈₄
0.05	$\mathcal{N}(100, 30)$	20×	0.051 (0.003)	
0.05	$\mathcal{N}(100, 30)$	50×	0.052 (0.003)	0.051 (0.002)
0.05	$\mathcal{N}(100, 30)$	100×	0.051 (0.002)	
0.20	$\mathcal{N}(100, 30)$	20×	0.202 (0.007)	
0.20	$\mathcal{N}(100, 30)$	50×	0.202 (0.006)	0.202 (0.006)
0.20	$\mathcal{N}(100, 30)$	100×	0.202 (0.006)	

Pairwise F_{ST} was estimated for various conditions of expected F_{ST} and coverage (Cov.). The pool size (n) 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. For Pool-seq data, we computed our estimator \hat{F}_{ST}^{pool} (Equation 13). The mean (RMSE) over 50 independent replicates of the `ms` simulations are provided, for a single pair of populations. For comparison, we computed WC₈₄ (Weir and Cockerham 1984) from allele count data inferred from individual genotypes (Ind-seq).

Table S3 Effect of variable coverage on pairwise F_{ST} estimates

F_{ST}	n	Pool-seq		Ind-seq
		Δ	\hat{F}_{ST}^{pool}	WC ₈₄
0.05	10	20	0.050 (0.006)	
0.05	10	50	0.050 (0.004)	0.050 (0.004)
0.05	10	100	0.050 (0.004)	
0.05	100	20	0.051 (0.003)	
0.05	100	50	0.051 (0.002)	0.051 (0.002)
0.05	100	100	0.051 (0.002)	
0.20	10	20	0.200 (0.007)	
0.20	10	50	0.200 (0.007)	0.200 (0.007)
0.20	10	100	0.200 (0.007)	
0.20	100	20	0.202 (0.006)	
0.20	100	50	0.203 (0.006)	0.203 (0.005)
0.20	100	100	0.203 (0.005)	

Pairwise F_{ST} was estimated for various conditions of expected F_{ST} and pool size (n). The coverage (δ_i) was varying across demes and loci, with $\delta_i \sim \text{Pois}(\Delta)$. For Pool-seq data, we computed our estimator \hat{F}_{ST}^{pool} (Equation 13). The mean (RMSE) over 50 independent replicates of the **ms** simulations are provided, for a single pair of populations. For comparison, we computed WC₈₄ from allele count data inferred from individual genotypes (Ind-seq).

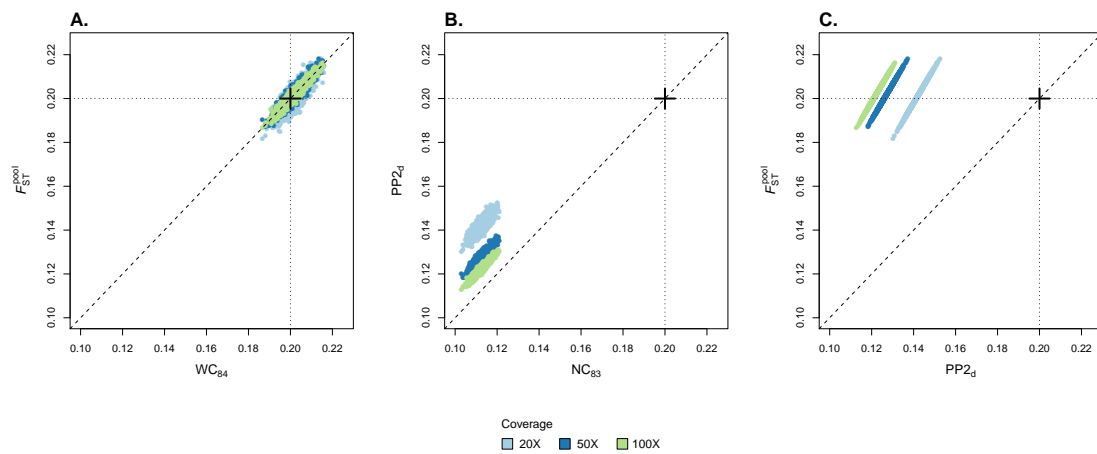


Figure S1 Pairwise estimators of F_{ST} . A. Multilocus estimates \hat{F}_{ST}^{pool} (computed using Equation 13) as a function of WC_{84} estimates computed from allele count data inferred from individual genotypes. B. Multilocus estimates $PP2_d$, as a function of NC_{83} estimates computed from allele count data inferred from individual genotypes. C. Multilocus estimates \hat{F}_{ST}^{pool} as a function of $PP2_d$ estimates. In each graph, the dots represent multilocus estimates of F_{ST} across all pairs of subpopulations from an 8-island model, and across 50 replicate *ms* simulations. We specified the migration rate corresponding to $F_{ST} = 0.20$. The size of each pool was fixed to 100. The results are shown for different coverages (20X, 50X and 100X). The cross indicates the simulated value of the parameter F_{ST} .

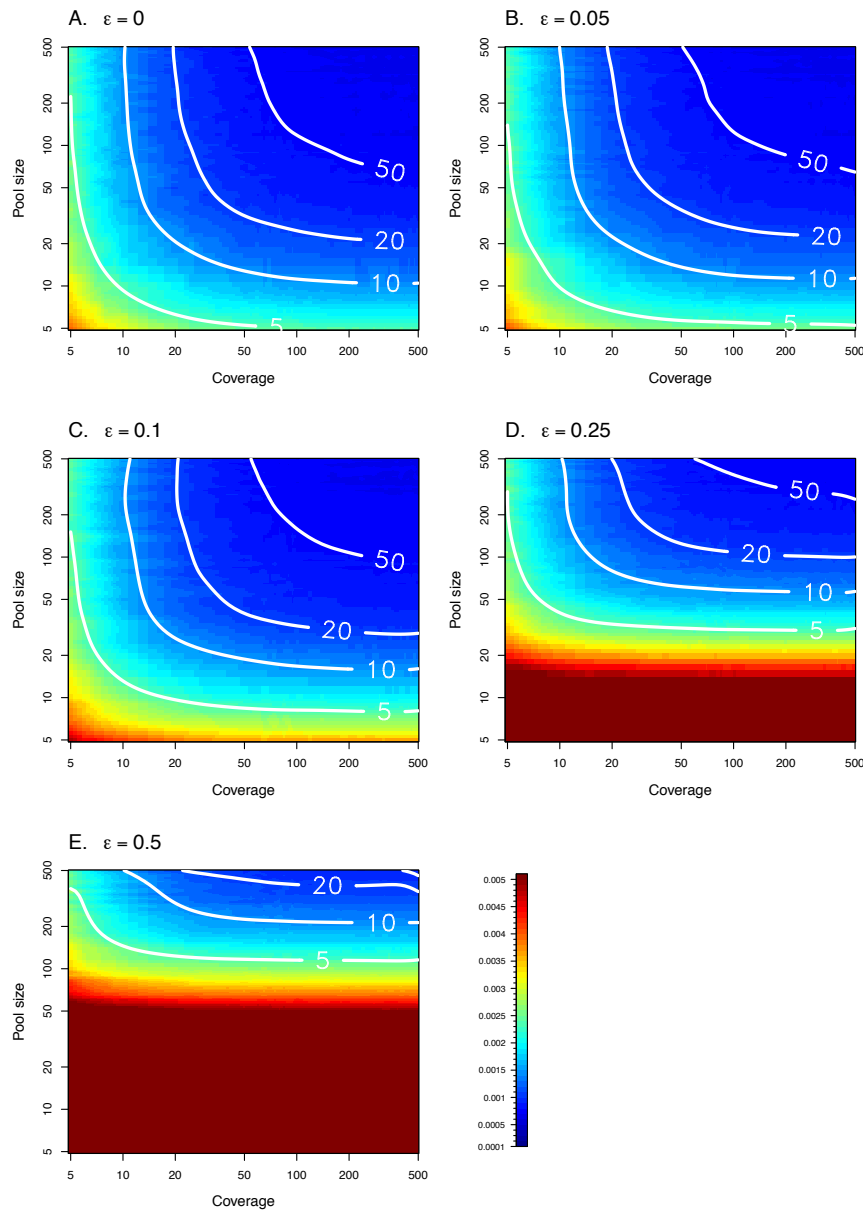


Figure S2 Root mean squared error (RMSE) of F_{ST} estimates for a wide range of pool sizes and coverage, with experimental error rate ϵ varying from 0 to 0.5 (A–E). Each density plot gives the RMSE of our estimator \hat{F}_{ST}^{pool} , using simple linear interpolation from a set of 44×44 pairs of pool size and coverage values. For each pool size and coverage, 500 replicates of 5,000 markers were simulated. Plain white isolines represent the RMSE of the WC_{84} estimator computed from Ind-seq data, for various sample sizes ($n = 5, 10, 20, 50$). 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).

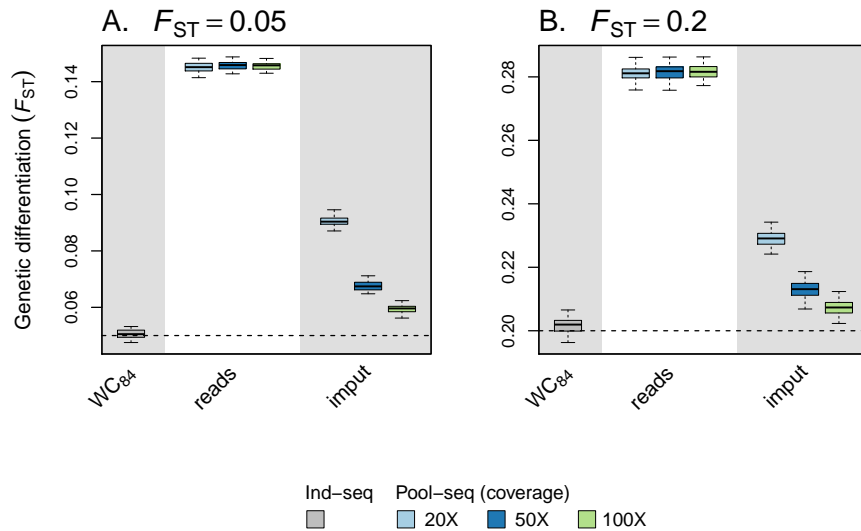


Figure S3 Global estimators of F_{ST} . We considered one estimator based on allele count data inferred from individual genotypes (Ind-seq): WC_{84} . For pooled data, we computed F_{ST} using the WC_{84} estimator: (i) directly from read counts, as if they were allele counts (“reads”); (ii) from allele counts imputed by maximum-likelihood (“imput”), as in Leblois et al. (2018). Each boxplot represents the distribution of multilocus F_{ST} estimates across all demes comparisons in an 8-island model, and across 50 independent replicates of the *ms* simulations. We used two migration rates, corresponding to $F_{ST} = 0.05$ (A) or $F_{ST} = 0.20$ (B). The size of each pool was fixed to 10. 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} .

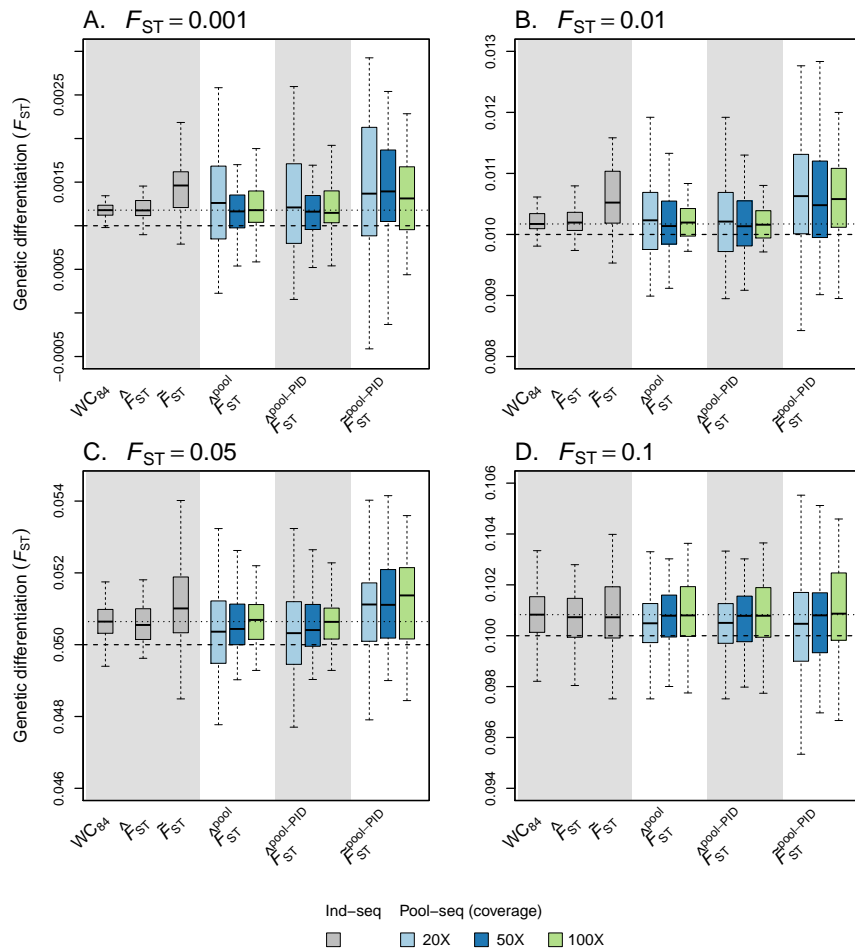


Figure S4 Precision and accuracy of alternative estimators of F_{ST} with varying pool size, for various levels of differentiation (A–D). The haploid pool size n drawn independently for each deme from a Gaussian distribution with mean 100 and standard deviation 30; n was rounded up to the nearest integer, with min. 20 and max. 300 haploids per deme. We considered three estimators based on allele count data inferred from individual genotypes (Ind-seq): WC_{84} , $\hat{F}_{ST} \equiv (\hat{Q}_1 - \hat{Q}_2) / (1 - \hat{Q}_2)$ (where \hat{Q}_1 and \hat{Q}_2 are the weighted frequencies of identical pairs of genes within and between subpopulations, respectively, with weights equal to the number of pairs of genes) and $\tilde{F}_{ST} \equiv (\tilde{Q}_1 - \tilde{Q}_2) / (1 - \tilde{Q}_2)$ (where \tilde{Q}_1 and \tilde{Q}_2 are the unweighted frequencies of identical pairs of genes within and between subpopulations, respectively). For Pool-seq data, we considered the estimators \hat{F}_{ST}^{pool} (Equation 12), $\hat{F}_{ST}^{pool-PID}$ (Equation A44) and $\tilde{F}_{ST}^{pool-PID}$ (Equation A45). Each boxplot represents the distribution of multilocus F_{ST} across 50 independent replicates of the ms simulations. 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.