

1
2
3
4
5
6
7

DR. ROBIN S WAPLES (Orcid ID : 0000-0003-3362-7590)

Article type : Resource Article

Pseudoreplication in genomics-scale datasets

Robin S. Waples^{*}, Ryan K. Waples[†], and Eric J. Ward^{*}

^{*}NOAA Fisheries, Northwest Fisheries Science Center
2725 Montlake Blvd. East, Seattle, WA 98112

[†]Department of Biology, Section for Computational and RNA Biology,
University of Copenhagen, Copenhagen, Denmark
Current address: Department of Biostatistics
University of Washington, Seattle WA

Running title: Precision in genomics datasets

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/1755-0998.13482](https://doi.org/10.1111/1755-0998.13482)

This article is protected by copyright. All rights reserved

Corresponding author:

Robin S. Waples

robin.waples@noaa.gov

+1 206 860-3254

July 2021

8 **Abstract**

9 In genomics-scale datasets, loci are closely packed within chromosomes and hence provide
10 correlated information. Averaging across loci as if they were independent creates
11 pseudoreplication, which reduces the effective degrees of freedom (df') compared to the nominal
12 degrees of freedom, df . This issue has been known for some time, but consequences have not
13 been systematically quantified across the entire genome. Here we measured pseudoreplication
14 (quantified by the ratio df'/df) for a common metric of genetic differentiation (F_{ST}) and a
15 common measure of linkage disequilibrium between pairs of loci (r^2). Based on data simulated
16 using models (*SLiM* and *msprime*) that allow efficient forward-in-time and coalescent
17 simulations while precisely controlling population pedigrees, we estimated df' and df'/df by
18 measuring the rate of decline in the variance of mean F_{ST} and mean r^2 as more loci were used.
19 For both indices, df' increases with N_e and genome size, as expected. However, even for large N_e
20 and large genomes, df' for mean r^2 plateaus after a few thousand loci, and a variance components
21 analysis indicates that the limiting factor is uncertainty associated with sampling individuals
22 rather than genes. Pseudoreplication is less extreme for F_{ST} , but $df'/df \leq 0.01$ can occur in
23 datasets using tens of thousands of loci. Commonly-used block-jackknife methods consistently
24 overestimated $\text{var}(F_{ST})$, producing very conservative confidence intervals. Predicting df' based
25 on our modeling results as a function of N_e , L , S , and genome size provides a robust way to
26 quantify precision associated with genomics-scale datasets.

27

28 Keywords: degrees of freedom; linkage disequilibrium; F_{ST} ; N_e ; genome size; jackknife
29 variance; simulations

30 1 | INTRODUCTION

31

32 It is now relatively easy to generate data at tens or hundreds of thousands of single-
33 nucleotide polymorphism (SNP) markers for non-model species, even in the absence of a
34 reference genome (da Fonseca et al. 2016; Van Wyngaarden et al. 2017; Aguirre et al. 2019;
35 Choquet et al. 2019; Minias et al. 2019). This opens up vast new opportunities for researchers
36 but also creates a host of analytical challenges, including ascertainment bias (Rosenblum and
37 Novembre 2007; Albrechtsen et al. 2010), phylogenetic inference (Leaché et al. 2015),
38 genotyping errors and missing data (Gautier et al. 2013; Graffelman et al. 2015; Huang and
39 Knowles 2014), and effects of selection (Foll and Gaggiotti 2008; Wolf and Ellegren 2017).

40 One topic that has not received sufficient attention with respect to genomics data for non-
41 model species is pseudoreplication—lack of independence among datapoints that reduces the
42 total information content. Statistical inference in biology is challenging because biological
43 systems are complex, variable, and subject to measurement and sampling errors. Replication is
44 generally necessary to ensure that apparently-interesting results are not due to small sample sizes
45 and chance. But true replication (which produces multiple independent datapoints) is difficult to
46 achieve, and flaws of experimental design and/or statistical inference that lead to overly-
47 optimistic estimates of statistical significance have been found to be widespread in ecology and
48 evolutionary biology (Hurlbert 1984; Ramage et al. 2013; Aarts et al. 2014; Colegrave and
49 Ruxton 2018; Lin et al. 2019).

50 In the present context, we are interested in use of genetic data to draw inferences about
51 key genetic parameters in real populations of non-model species. For decades, most studies were
52 limited to a few dozen genetic markers, and it was routinely assumed that each marker was
53 independent—or if not, that departures from independence were small enough to be safely
54 ignored. This assumption is no longer tenable in the age of genomics. Most species have at
55 most a few dozen chromosomes (Table 1), so in contemporary datasets a typical chromosome
56 contains thousands of markers—a situation which guarantees that pseudoreplication occurs.

57 Here we focus on two widely-used genetic metrics: F_{ST} , a measure of differentiation

58 among populations; and r^2 , a measure of linkage disequilibrium (LD) at pairs of loci. Inter-locus
59 sampling variances for both of these metrics are large, so a high degree of replication is required
60 to obtain reliable estimates of the mean. Let $\hat{F}_{ST(L)}$ be the estimate of mean F_{ST} based on data for
61 L diallelic loci, and let $\hat{r}^2_{(L)}$ be the two-locus analogue for LD based on $n = L(L-1)/2$ pairs of L
62 loci, where the ‘^’ indicates an estimate. In statistical theory, if one has k independent
63 observations of a random variable x with standard deviation σ , the standard error of the estimate
64 of $\text{mean}(x)$ is $\sigma_{\bar{x}} = \sigma/\sqrt{k}$. The genetic metrics of interest here also have the property that their
65 variances are proportional to $1/k$, assuming the data points are independent (Lewontin and
66 Krakauer 1973; Hill 1981).

67 In population genetic studies, two general sources of replication are available: one can
68 sample multiple individuals from each population, and one can sample multiple loci from each
69 individual. Often the number of individuals that can be sampled is limited by population size or
70 logistical constraints, in which case the only feasible method to increase precision is by
71 increasing the number of loci. Easy access to genomics-scale datasets has made it possible to
72 vastly increase the number of loci sampled per individual—but to what extent do these genes
73 provide independent information about population-level parameters of interest?

74 We encounter two kinds of pseudoreplication we need to be concerned with. First, if loci
75 do not freely recombine, information they provide is correlated, so (for example) $\text{var}(\hat{F}_{ST(L)})$
76 does not decline as fast with addition of more loci as would be the case if the loci were
77 independent. The total information content of a genetic dataset is constrained by the amount of
78 recombination. This constraint means, for example, that even with many millions of SNPs
79 available in human genetics, it is not possible to confidently resolve distant familial relationships
80 with genetic data alone (Thompson 2013). Pseudoreplication due to lack of independent
81 assortment applies both to \hat{F}_{ST} and \hat{r}^2 . Analyses of LD also generate a second kind of
82 pseudoreplication, caused by overlapping pairs of loci. The approximately $L^2/2$ different pairs of
83 loci are not independent because each locus occurs in $L-1$ pairwise comparisons.

84 Problems related to this lack of independence have not been quantified in any systematic
85 way for non-model species. For such species, a typical experimental design involves orders of
86 magnitude more SNP loci than individuals, and if any detailed genomic mapping information is
87 available, loci typically can only be placed on short genomic scaffolds rather than full
88 chromosomes. In this study we investigate how much precision is reduced by pseudoreplication,

89 compared to what it would be if the assumption of independence were completely satisfied.
90 Assuming one has data for L diallelic (SNP) loci, the number of datapoints (and the nominal
91 degrees of freedom, df , associated with the overall estimate) is L for F_{ST} and $n \approx L^2/2$ for the LD
92 analyses. Pseudoreplication causes the actual (effective) df' (L' or n') to be less than the total
93 number of datapoints (Cox 1984; Giesbrecht 2006). The key question thus can be framed as
94 follows: How much smaller is the effective df' than the nominal df , and how does the df'/df ratio
95 depend on aspects of experimental design (number of individuals (S) and loci sampled) and
96 uncontrollable parameters of the population(s) of interest (genome size or number of
97 chromosomes, C , and effective population size, N_e)? The ability to estimate L' or n' would allow
98 an unbiased evaluation of precision and facilitate placing accurate confidence bounds on
99 estimates of key population genetic parameters.

100 For species (including humans) with reference genomes and linkage map, these
101 dependencies have been addressed to some extent, particularly with respect to multiple testing
102 (Nyholt 2004; Pe'er et al., 2008; Galwey 2009). A popular approach is to use a weighted, block-
103 jackknife that breaks up the genome into blocks of contiguous loci and leaves each block out in
104 jackknife fashion (Busing et al. 1999). LD pruning (Purcell et al. 2007)—excluding loci to
105 reduce LD—can also be used to generate sets of loci that act more independently and thus reduce
106 prereplication. These approaches, however, require detailed mapping information, and even so
107 they only deal with correlations of the original variables. To evaluate the second type of
108 pseudoreplication noted above, it is necessary to consider second-order correlations—that is, the
109 degree to which r^2 (locus 1 \times locus 2) is correlated with r^2 (locus 1 \times locus 3) and r^2 (locus 2 \times
110 locus 3). In theory n' could be calculated this way, but it would require one to specify the
111 relevant covariance matrix, and with $\approx L^2/2$ pairwise correlations of loci the covariance matrix of
112 these correlations has $\approx L^4/8$ elements. For a dataset with 1.5 million SNPs, calculating n'
113 therefore would require one to specify over 6.3×10^{23} elements in the covariance matrix—more
114 than Avogadro's number! Not surprisingly, we are not aware of any attempts to do this.

115 The objectives of this paper are to develop model-based approximations of df' , and to
116 provide general guidance – given known or measurable covariates (C, L, S, N_e). Our approach to
117 quantifying the degree of pseudoreplication involves simulating a large number of replicate
118 populations, and for each replicate we calculate mean values for both of the genetic indices (\hat{r}^2 ,
119 \hat{F}_{ST}). Observed variances of the multilocus metrics across replicates allow us to calculate the df'

120 associated with samples of individuals and gene loci. This process was repeated for several
121 evolutionary scenarios involving different combinations of C and N_e .

122 We have the following general expectations:

- 123 1) As more loci are packed into a fixed number of chromosomes, the number of loci per
124 chromosome increases, minimum distance between a new locus and existing loci shrinks
125 (Figure S1), and the marginal increase in information content provided by each new locus
126 declines. Therefore, we expect that as the ratio L/C increases, pseudoreplication will
127 increase and the ratio of effective to nominal df will decrease for both metrics. In
128 addition, the rate of decay of LD with distance between loci increases with effective
129 population size (Figure S2). So, we expect that, for a given L/C ratio, n' will be smaller
130 for populations with smaller N_e .
- 131 2) It seems intuitive that the magnitude of LD increases with the number of loci, but that is
132 not actually the case. Each new locus increases opportunities for that locus to be in LD
133 with existing loci, but this is balanced by new pairwise comparisons with loci on different
134 chromosomes. As a result, the probability that any two randomly-chosen loci will be in
135 any particular LD association is independent of the number of loci, but it does depend on
136 genome size and N_e (Waples et al. 2016). Therefore, we expect that the first type of
137 pseudoreplication for LD analyses will be inversely correlated with both N_e and C . The
138 second type of pseudoreplication arising from LD analyses—that caused by overlapping
139 pairs of the same loci—has received little study. However, we note that, of the $\sim L^2/2$
140 pairwise comparisons of L loci, only $L/2$ are completely independent (non-overlapping),
141 so the proportion of independent comparisons is $1/L$. Hence, we expect that this type of
142 pseudoreplication will increase with the number of loci.

143 We hope to find one or both of the following:

- 144 a. Patterns in empirical variances across simulated datasets that allow us to provide
145 general guidance for users interested in predicting df' , based on measurable or
146 estimable covariates (C, L, S, N_e).
- 147 b. Existing jackknife methods prove to be reliable at estimating precision for a given
148 dataset.

149

150 2 | METHODS

151

152 Below we provide an overview of methods used; for more details, see Supporting Information.

153 Table 2 summarizes notation.

154

155 **2.1 | Conceptual Framework**

156 Our focus is on actual populations, each of which has a single, realized population
157 pedigree (Wakeley et al. 2012; Ralph 2019). To mimic sampling of individuals and genes from
158 real populations, we combine efficient coalescent and forward simulation programs (*SLiM*,
159 Messer 2013; and *msprime*, Kelleher et al. 2016) that allow control over multigenerational
160 pedigrees in replicate populations (Haller et al. 2019). Our experimental design explicitly
161 models two major sources of uncertainty in estimating population-level parameters: sampling of
162 individuals and sampling of genes. Both processes are important when evaluating uncertainty
163 around parameter estimates. If replication only occurs across genes, estimates will converge on
164 values determined by the pedigree of the observed individuals (Waples and Faulkner 2009;
165 Wakeley et al. 2012; King et al. 2018). In general, however, one wants to draw inference about
166 the entire population, and uncertainty related to sampling individuals from the population cannot
167 be eliminated by intensive replication across genes.

168

169 **2.2 | Modeled Scenarios and Simulations**

170 For 16 different evolutionary scenarios (combinations of effective size
171 [$N_e=50,200,800,3200$] and number of chromosomes [$C=1,4,16,64$]), we simulated four separate
172 ancestral populations of $N=N_e$ diploid individuals for $10N_e$ generations under Wright-Fisher
173 (WF) reproduction and ensured that each ancestral population fully coalesced. The simulations
174 modeled the process of recombination and mutation in realistic-sized genomes with distinct
175 chromosomes. Genomes were simulated with varying numbers of chromosomes [1, 4, 16, 64],
176 each of which was 50 Mb (5×10^7 bp) long, with a recombination rate of 1×10^{-8} per bp per
177 generation. Each ancestral population then split into four daughter populations of size N_e , which
178 evolved in isolation for enough generations that expected $F_{ST} \approx 0.05-0.1$, common values for
179 natural populations of many species. From this point, the models differed slightly for analyses of
180 F_{ST} and LD.

181 In the LD analyses, for each of the $4 \times 4 = 16$ population pedigrees, we created eight non-
182 overlapping sets of loci by adding mutations to gene trees present in the population pedigree
183 (Figure 1). Gene trees reflect the evolutionary history of genes and haplotypes back in time. For
184 each scenario, the $4 \times 4 \times 8$ hierarchical design produced 128 replicate populations with up to 75K
185 diallelic loci (up to 100K loci for $N_e = 3200$). From each replicate population, we took four
186 random subsamples of $S = 25, 50, 100$ diploid individuals (only $S \leq 50$ for $N_e = 50$), and data for each
187 subsample were analyzed for $L = 100 - 75,000$ loci. Except for $N_e = 3200$, we also took exhaustive
188 samples of the population ($S = N_e$). This experimental design allowed us to calculate two separate
189 variance components: V_1 (same individuals, different loci) and V_2 (same loci, different but
190 potentially overlapping samples of individuals) (Table 3). Data from cells on the diagonal
191 (different individuals and different loci) were used to calculate df' for the LD analyses (n'), as
192 described below.

193 In the F_{ST} analyses, for each of the four ancestral populations, the four daughter
194 populations shown in Figure 1 allowed 6 pairwise comparisons among populations, for a total of
195 24 two-population pedigrees (Figure S3). Six mutational replicates of 200K loci were created
196 for each pedigree, and eight samples for each of several different sizes were taken from each
197 mutational replicate.

198 The resulting distributions of minor allele frequency (MAF) followed the familiar U-
199 shaped pattern expected at mutation-drift equilibrium (Figure S4). In addition to the
200 *SLiM/msprime* simulations that organize the genome into chromosomes, for LD we also modeled
201 scenarios with unlinked loci and infinite N_e to provide insights into asymptotic behavior.

202

203 2.3 | Genetic Indices

204 For each replicate in each evolutionary scenario, we calculated $\hat{F}_{ST(L)}$ and $\hat{r}^2_{(L)}$ across
205 variable numbers of loci or locus pairs. LD analyses were conducted using both all pairs of loci,
206 and only pairs on different chromosomes. For each pair, the sample estimate of the squared
207 correlation coefficient (\hat{r}^2) was computed using the Pearson product-moment correlation of
208 diploid genotypes. After applying the sample-size adjustments implemented in *LDNE* (Waples
209 and Do 2008), mean \hat{r}^2 can be used to estimate N_e . Because this method assumes all loci are
210 unlinked, use of all locus pairs leads to a predictable pattern of downward bias in \hat{N}_e that is
211 inversely proportional to $\log(C)$ (Waples et al. 2016), whereas \hat{N}_e is unbiased when computed

212 only using pairs of loci on different chromosomes (Figure S5). To reduce effects of rare alleles,
213 within each sample we omitted loci with two or fewer copies of the minor allele.

214 We used two methods for estimating the standardized variance of allele frequency among
215 populations, F_{ST} . Nei's (1973) gene diversity method calculates $G_{ST} = (H_T - H_s) / H_T$, where H_s is
216 the expected within-population heterozygosity (averaged across both populations) and H_T is
217 expected total heterozygosity, based on mean allele frequencies across both populations, without
218 an adjustment for sample size. Calculated this way, and with just two alleles at each locus, G_{ST} is
219 identical to Wright's (1951) F_{ST} (Nei and Chakravarti 1977); we refer to this estimator as \hat{F}_{ST}^{Nei} .

220 Another widely-used measure of population differentiation is the coancestry coefficient,
221 θ (Cockerham 1969; Reynolds et al. 1983); the relationship between the two parameters is given
222 by (Cockerham and Weir 1987)

$$223 \quad \theta = \frac{sG_{ST}}{G_{ST} + s - 1},$$

224 (1)

225 where s is the number of subpopulations. For $s=2$ (as considered here) this reduces to
226 $\theta = 2G_{ST} / (G_{ST} + 1)$, which is close to $2G_{ST}$ if G_{ST} is small. We considered Hudson et al.'s (1992)
227 estimator $\hat{F}_{ST}^{Hudson} = 1 - H_W / H_B$, where H_W and H_B are the mean numbers of allelic differences
228 within and between populations, respectively; the formula we used was Equation 10 in Bhatia et
229 al. (2013). Under conditions modeled here (two populations, equal sample sizes), \hat{F}_{ST}^{Hudson} is
230 identical to Weir and Cockerham's (1984) $\hat{\theta}$, or nearly so (Bhatia et al. 2013).

231 Because performance of multilocus F statistics is generally better when they are
232 computed as ratios of multi-locus means rather than means of single-locus ratios (Jorde and
233 Ryman 2007; Bhatia et al. 2013), we evaluated both methods for calculating \hat{F}_{ST}^{Nei} (\hat{F}_{ST}^{Hudson}
234 already is a ratio of averages). We considered three ascertainment schemes for identifying
235 variable loci: 1) loci variable in at least one of the two populations; 2) loci variable in the
236 ancestral population; 3) loci with overall MAF ≥ 0.05 in the two samples combined.

237

238 **2.4 | Effective Degrees of Freedom and Effective Number of Loci**

239 For analysis of population differentiation, effective df was calculated from observed
240 variances of $\hat{F}_{ST(L)}$ using the relationship $\phi_F = Var(\hat{F}_{ST}) / E^2(\hat{F}_{ST}) =$ squared coefficient of
241 variation of \hat{F}_{ST} (Lewontin and Krakauer 1973). If $\hat{F}_{ST(L)}$ is based on L independent diallelic

242 loci, the following relationship should hold: $\phi_F=2/L$. Therefore, we computed df' from the
243 empirical variance of $\hat{F}_{ST(L)}$ as $L'=2/\phi_F$. L' can be interpreted as the number of independent loci
244 that would be expected to produce the value of ϕ_F observed in the data (Cox 1984; Giesbrecht
245 2006). The ratio L'/L therefore provides an index of how much nonindependence has increased
246 variance of the estimator.

247 Many commonly-used estimators of F_{ST} and related quantities (Weir and Cockerham
248 1984; Hudson et al. 1992; Patterson et al. 2012) include adjustments for sampling individuals
249 and/or populations, and this complicates comparisons with theoretical expectations of the
250 variance-to-mean ratio that depend on raw F_{ST} values. This is why, for estimating L' , we used
251 Nei's (1973) gene diversity method, which does not include a sample-size adjustment. For
252 Hudson's estimator, we measured the rate of decline in the variance of $\hat{F}_{ST(L)}^{Hudson}$ as more loci were
253 used in the analysis, and this provided an alternative way to quantify the degree of
254 pseudoreplication.

255 Hill (1981) showed that the relationship $\phi_r=2/n$ also applies to $\hat{r}^2_{(L)}$, under the
256 assumption that all n pairwise comparisons of L loci are independent. Accordingly, we estimated
257 df' for LD analyses as $n' = 2/\phi_r$. Because the effective number of pairs of loci is not a very
258 intuitive metric, in some cases (Figures 2 and 3) we presented LD results in terms of the effective
259 number of loci (L'), which is the number of loci that would produce the actual observed variance
260 of $\hat{r}^2_{(L)}$, if all the resulting locus pairs were independent. To a good approximation, if n' is the
261 effective number of locus pairs, then the effective number of loci that would produce n' is $L' \approx$
262 $\sqrt{2n'}$ [the exact value is $0.5+\sqrt{(0.25+2n')}$, but this simple approximation is much more intuitive].
263

264 **2.5 | Model Fitting**

265 We considered a wide range of covariates [N_e , C , S , L , and (for LD) n] as predictors of
266 df' . To account for potential non-linearities, we also considered transformed responses of these
267 original variables (Table S1). Because of asymptotic relationships between df' and the predictor
268 variables (Figure S6), we focused inference on fitting statistical models with asymptotes, rather
269 than using linear models. The general function we used to parameterize the models was

$$270 \quad f(x) = \frac{x}{\left[\frac{1}{p^r} + \frac{x^r}{q^r} \right]^{\frac{1}{r}}}, \quad (2)$$

271 where parameters p , q , and r control the shape of the function. When shape parameter $r=1$, this
272 function takes the form of the familiar Michaelis–Menten equation (also known in fisheries as
273 the Beverton–Holt stock recruitment model; Beverton and Holt 1957). We considered models
274 with both covariates and transformed functions of covariates; these latter models treated each
275 parameter (p , q , r) as linear functions of other predictors; e.g., $p = b\mathbf{X}$, where \mathbf{X} is a design
276 matrix of predictors and b are estimated coefficients. All models were fit in R (R Core Team
277 2020) using maximum likelihood. We used Akaike’s Information Criterion (AIC) to evaluate
278 which combinations of parameters were best supported. Additional details about model fitting
279 are in Supporting Information.

280

281 **2.6 | Confidence Intervals**

282 To evaluate accuracy of our estimates of df' , for selected scenarios we generated many
283 (>1000) new samples of individuals and loci, and for each sample we calculated $\hat{F}_{ST(L)}$ or $\hat{r}^2_{(L)}$
284 averaged over data for 500–5000 SNPs (for LD) or 5000–50000 SNPs (for F_{ST}). We calculated
285 confidence intervals (CIs) around these means using standard statistical theory (Equations S6 and
286 S7). Width of the CIs was calculated two ways, using: 1) our modeled estimates of df' ; 2) a
287 published jackknife method (Busing et al. 1999 for F_{ST} , Jones et al. 2016 for \hat{N}_e based on r^2).

288

289

290 **3 | RESULTS**

291

292 **3.1 | Linkage Disequilibrium**

293 The four different ancestral populations produced modest differences in mean r^2 when
294 averaged across all descendant populations, and within each daughter population the eight
295 mutational replicates also produced relatively small differences in mean r^2 (Figure S7). In
296 contrast, daughter populations descended from the same ancestral population varied substantially
297 in the mean magnitude of LD. This result emphasizes the importance of accounting for recent
298 population pedigrees in assessing variability of genetic indices and shows the sensitivity of LD-
299 based estimates of N_e to recent effective population size (Waples and Faulkner 2009).

300

301 **3.1.1 | Effective degrees of freedom:** Table S2 gives our estimates of both n' and L' for r^2 for

302 every scenario we simulated. We expected substantial effects of genome size and N_e on n' and
303 L' , and those expectations were borne out (Figure 2 shows results for L' ; see Figure S8 for the
304 same results in terms of n'). For a given number of loci, L' increased smoothly as the number of
305 chromosomes increased from 1 to 64, with results for $C = 64$ being largely indistinguishable
306 from those for unlinked loci (Figure 2, top). The influence of effective size was even stronger:
307 L' increased systematically as N_e increased from 50 to 3200 (Figure 2, bottom); however, for
308 large numbers of loci, L' for $N_e=3200$ was still substantially lower than the $L' \approx L$ that was found
309 for $N_e=\infty$. With 5000 loci, $\text{var}(r^2)$ decreased by a factor of 5 when N_e increased four-fold from
310 50 to 200, and the variance decreased by a factor of 68 when N_e increased 64-fold to 3200 (using
311 data from Table S2 with $S=25$ and $C=4$). In contrast, quadrupling the number of chromosomes
312 (from 1 to 4, with N_e fixed at 50) decreased $\text{var}(r^2)$ only by 2.5 \times , and a 64-fold increase in C
313 reduced $\text{var}(r^2)$ only by 10 \times .

314 Relatively speaking, sample size of individuals has less influence on L' : with $C=16$, a
315 four-fold increase in S reduced $\text{var}(r^2)$ by only 1.9 \times for $N_e=200$ and by only 6% for $N_e=800$.
316 There was one notable exception, however: we found a qualitative difference in L' between
317 scenarios in which the entire population was sampled and those in which only a subset of
318 individuals was analyzed (Figure 3). With exhaustive sampling ($S=N=N_e$), L' continues to rise,
319 albeit increasingly slowly, as larger and larger numbers of loci are used to compute $\hat{r}^2_{(L)}$. In
320 contrast, for incomplete sampling ($S < N$), L' rapidly plateaus once about 1000 loci are used, after
321 which additional loci do little to increase precision. As a consequence, for 10,000 SNPs L' is
322 only about 3% of the number of loci used (and n' is about 3 orders of magnitude smaller than the
323 number of pairwise comparisons; Table S2). Furthermore, as long as the entire population is not
324 exhaustively sampled, sample size has relatively little effect on L' (Figure 3). Samples of
325 individuals were all drawn from the $N=N_e$ individuals in the final simulated generation, and the
326 variance associated with replicate hypergeometric samples is proportional to $(N-S)/N$. *A priori*,
327 therefore, we expected that L' would increase smoothly with sample size. Two factors likely
328 explain the patterns actually observed.

329 First, larger samples produce less sampling error (Hill 1981; Waples 2006), which
330 reduces both $\text{var}(\hat{r}^2)$ and $\text{mean}(\hat{r}^2)$. Because n' depends on ϕ_r which is the ratio of $\text{var}(\hat{r}^2)$ to
331 $[\text{mean}(\hat{r}^2)]^2$, the net effect of increasing S is only a modest change in n' and L' . Second, the
332 variance components analysis (Figure 4) shows that whereas V_1 (variance among replicate sets of

333 loci assayed on the same individuals) continues to decline rapidly with increasing numbers of
334 loci, V_2 (variance among different samples of individuals assayed for the same loci) does not
335 decline much after about 500-1000 loci are used. As a consequence, except for small numbers of
336 loci, V_2 dominates overall $\text{var}(\hat{r}^2_{(L)})$ and hence this component largely determines n' . Although
337 sampling a larger fraction of the population does reduce V_2 , the resulting variance still is much
338 larger than V_1 and still dominates n' and L' . Note also that the actual variance of $\hat{r}^2_{(L)}$ is smaller
339 than the sum of V_1 and V_2 (Figure 4), which indicates that the two variance components must be
340 negatively correlated to some extent. In what follows, we focus on $S < N$, which is the most
341 common scenario in studies of natural populations.

342 Restricting LD analyses to pairs of loci on different chromosomes reduces the number of
343 pairwise comparisons by the proportion $1/C$, but this has little effect on n' , which is essentially
344 the same regardless whether same-chromosome comparisons are allowed or not (Figure S9; the
345 difference is a bit larger for small genomes, where linkage has a stronger effect). Unless
346 otherwise noted, all results presented are for all pairwise comparisons.

347 The fact that LD results for unlinked loci approximate those for simulations with 64
348 chromosomes, and that both scenarios show that L' and n' both fail to increase much after a few
349 thousand loci are used (Figures 2 and S8, top), indicate that physical linkage is not the primary
350 factor that reduces df' for analyses involving LD. The second factor relevant to two-locus
351 analyses involves overlapping pairs of the same loci. We expect that, if we have genotypes for
352 four unlinked loci (w, x, y, z) and compute all six pairwise correlations, the correlation $[w, x]$ will
353 be independent of the correlation $[y, z]$, but to what extent does the correlation $[w, x]$ provide
354 independent information to the correlation $[w, y]$, since one locus is shared?

355 To evaluate this factor, we simulated unlinked loci in many replicate populations and
356 calculated r^2 for each pair of loci in each replicate, across all N_e individuals in the population
357 (see Detailed Methods in Supporting Information). Then, across all replicates, we computed the
358 squared correlation between r^2 values that did and did not share one locus. Results show that
359 correlations between pairs of r^2 values that did not share a locus were essentially 0, regardless
360 how large or small N_e was (Figure S10). This is also the result one obtains if one compares
361 correlations between pairs of vectors of *i.i.d.* random variables, regardless whether the pairwise
362 correlations being compared share one variable or not (data not shown). However, when the data
363 being compared are generated by a process mediated by a pedigree (as occurs during

364 reproduction in finite populations), then the pairwise correlations are not independent when they
365 share one variable, and the degree of non-independence is inversely related to N_e (Figure S10).

366

367 **3.1.2 | Confidence intervals:** For datasets with more than 1000 loci, the number of pairwise
368 comparisons of loci is of order 10^6 or higher, in which case parametric CIs for \hat{N}_e that assume all
369 datapoints are independent become vanishingly small (Figure 5). Based on estimates of df'
370 obtained in this study, actual CIs for large numbers of loci are much wider, and this difference is
371 important to understand to avoid misleading conclusions about precision.

372 Our results provide a basis for users to develop robust CIs for their own datasets. The
373 best model to predict n' included several sets of covariates for each parameter in the Michaelis-
374 Menten asymptotic function (p, q, r). Parameter estimates, standard errors, and more details are
375 included in Supporting Information and Tables S3 and S5. When fit to the original data, the
376 correlation between $\log(\text{predicted } n')$ and $\log(\text{true } n')$ was 0.997 (Figure S11). When we
377 evaluated performance of CIs for \hat{N}_e using Equation S5 based on these predicted n' values, the
378 fraction of 90% CIs that contained the true N_e was close to the expected 0.9 (data not shown).
379 But a typical user will know covariate values only for the numbers of individuals and loci in their
380 samples and will have to estimate N_e (from the genetic data, or elsewhere) and perhaps C (e.g.,
381 from a related species), and finally estimate n' based on our modeling results. Accounting for
382 these additional sources of uncertainty reduced CI performance only slightly, such that overall
383 coverage was 90% for $S=100$, 88% for $S=50$, and 87% for $S=25$ (Table 4).

384 Performance of the Jones et al. (2016) jackknife method varied with sample size (Table
385 4), which is not surprising given that this method jackknives over individuals. For $S=100$ the
386 method produced conservative 90% CIs that included the true N_e >93% of the time, indicating
387 that n' was generally underestimated. Overall coverage was close to the expected 90% for $S=50$
388 but only 87% for $S=25$, indicating a tendency of the jackknife method to overestimate precision
389 for small samples. Across all scenarios, most of the jackknife CIs that did not contain the true
390 value were too high, meaning that the lower bound was larger than true N_e , and this effect was
391 much stronger for smaller samples.

392 Scatterplots of estimated n' vs mean r^2 for individual datasets (Figure S12) help to
393 explain performance of Jones's jackknife method. First, \hat{n}' varied considerably across replicates,
394 spanning one order of magnitude for $S=100$ and more than two orders for $S=25$. Second, for all

395 scenarios we found a strong, negative correlation between \hat{n}' and mean r^2 . Datasets with below-
396 average mean r^2 consistently were estimated to have a relatively high n' and hence relatively
397 narrow CIs. Relatively small r^2 translates into a relatively large estimate of N_e , around which the
398 CI was relatively narrow. In combination, these factors produce an excess of large \hat{N}_e estimates
399 whose lower CI bound is larger than true N_e , and this effect is exacerbated for small S . Because
400 n' is an increasing function of N_e , estimates of n' based on modeling results in this study also
401 show a negative correlation between \hat{n}' and mean r^2 , but the range of \hat{n}' for the current study was
402 about half that of the Jones jackknife method (Figure S12).

403

404 3.2 | F_{ST}

405 An initial sensitivity analysis (Figures S13-S14 and Supporting Information) showed that
406 Nei's weighted \hat{F}_{ST}^{Nei} performed better than the unweighted version, and ascertainment method 3
407 performed better than the other options. Therefore, results that follow apply to \hat{F}_{ST}^{Nei} and \hat{F}_{ST}^{Hudson}
408 with a MAF cutoff of 0.05.

409 With no pseudoreplication, variance of $\hat{F}_{ST(L)}$ should be inversely proportional to L ; the
410 actual rate of decline in $\text{var}(\hat{F}_{ST(L)})$ as the number of SNPs increases is shown in Figure 6. For a
411 scenario with $N_e=200$, $C=4$, and $S=25$, the decline is approximately log-linear up to $L \approx 2 \times 10^3$,
412 after which point addition of more loci does little to further reduce the variance. For a scenario
413 with larger N_e , S , and genome size, the rate of decline in $\text{var}(\hat{F}_{ST(L)})$ does not start to plateau until
414 the number of loci is an order of magnitude larger (Figure S15). Notably, in both scenarios $\text{var}(\hat{F}_{ST}^{Hudson})$
415 declines at the same rate as $\text{var}(\hat{F}_{ST}^{Nei})$.

416 For a given sample size, mean values of the Nei and Hudson estimators are perfectly
417 correlated with a common slope and sample-size specific intercepts (Figure S16); the slopes and
418 intercepts, however, vary with mean \hat{F}_{ST} . This linear relationship explains why the rate of
419 decline in $\text{var}(\hat{F}_{ST}^{Hudson})$ with increasing numbers of loci is nearly identical to that for $\text{var}(\hat{F}_{ST}^{Nei})$.

420

421 **3.2.1 | Effective degrees of freedom:** Table S4 gives our estimates of L' for F_{ST} for every
422 scenario we simulated. As with LD, we found that L' for F_{ST} depends strongly on both N_e and
423 genome size (Figure 7). In contrast to the two-locus results, however, we found a clear F_{ST}
424 asymptote for L' only for smaller values of N_e and C ; for larger values, L' was still increasing

425 after bringing 200K loci into the analysis. We found a moderate effect on L' of sample size of
426 individuals, which became more pronounced for large numbers of loci (Figure S17).

427 For a given number of loci, L'/L was considerably higher for F_{ST} than n'/n was for LD
428 (Figure S18). Nevertheless, in populations with small N_e and few chromosomes that are assayed
429 for large numbers of loci, the effective df associated with mean F_{ST} can be two orders of
430 magnitude or more smaller than the number of SNPs. The variance components analysis for F_{ST}
431 (Figure S19) produced a general pattern similar to that for LD (Figure 4), with one important
432 difference. For the same evolutionary scenario ($N_e=200$, $C=16$, $S=50$), whereas V_2 for LD
433 rapidly diverges from V_1 and starts to plateau after ~ 500 -1000 loci, the divergence comes much
434 later for F_{ST} , and V_2 does not begin to level off until 10^4 - 10^5 loci are included.

435 Figure 8 provides a detailed look at variation in \hat{F}_{ST}^{Hudson} across 48 samples from each of
436 two 2-population pedigrees that produced substantially different mean $\hat{F}_{ST(L)}$ values. Within
437 each pedigree, for each of eight samples of individuals, variation among six replicate samples of
438 L loci reflects variance component V_1 .

439 Our estimates of L' for $\hat{F}_{ST(L)}$ can be used to estimate realistic variances for $\hat{F}_{ST(L)}^{Hudson}$ as
440 follows: (1) Given that $L' = 2/\phi_F$ and $\phi_F = \text{var}(\hat{F}_{ST(L)}^{Nei})/[\text{mean}(\hat{F}_{ST(L)}^{Nei})]^2$, the variance of Nei's
441 multilocus estimator declines according to $\text{var}(\hat{F}_{ST(L)}^{Nei}) = \text{var}(\hat{F}_{ST(1)}^{Nei})/L'$, where $\text{var}(\hat{F}_{ST(1)}^{Nei})$ is the
442 variance among single-locus \hat{F}_{ST} values and $\text{var}(\hat{F}_{ST(1)}^{Nei}) = 2*[\text{mean}(\hat{F}_{ST(L)}^{Nei})]^2$. (2) Because $\text{var}(\hat{F}_{ST(L)}^{Hudson})$
443 declines at the same rate as $\text{var}(\hat{F}_{ST(L)}^{Nei})$, $\text{var}(\hat{F}_{ST(L)}^{Hudson}) = \text{var}(\hat{F}_{ST(1)}^{Hudson})/L'$, where $\text{var}(\hat{F}_{ST(1)}^{Hudson})$
444 is the variance of single-locus estimates and can be calculated from empirical data.

445
446 **3.2.2 | Confidence intervals:** The best model to predict L' for F_{ST} was similar in form to that
447 found for LD (Supporting Information). When fit to the original data, the correlation between
448 $\log(\text{predicted } L')$ and $\log(\text{true } L')$ was >0.99 (Figure S11). CIs based on modeled estimates of
449 $\text{var}(\hat{F}_{ST(L)}^{Nei})$ obtained in this study contained the true value 89-95% of the time (mean 91.5%), and
450 of those that did not, roughly equal numbers were too high and too low (Figure 9 and Table 5).

451 Across all scenarios [N_e, C, L, S], we found three consistent patterns in block-jackknife
452 results, and each applied equally to the Nei and Hudson estimators. (1) Block-jackknife
453 estimates of variance are positively correlated with mean $\hat{F}_{ST(L)}^{Nei}$ and mean $\hat{F}_{ST(L)}^{Hudson}$ (Figure S20);
454 (2) Virtually all block-jackknife estimates exceeded the actual variances we calculated from our

455 simulations (Figure 9 and Table 5); (3) Replicate samples generated using the same parameters
456 produced block-jackknife estimates of variance that ranged widely in magnitude (Figure S20),
457 and this variation was greater for the larger block size (one chromosome) and smaller samples of
458 individuals. Because the 5Mb blocks showed less variation among replicates, we focus on those
459 results. Upward bias in the block-jackknife estimates of variance produced overly conservative
460 confidence intervals that almost always contained the true values of $\hat{F}_{ST(L)}^{Nei}$ and $\hat{F}_{ST(L)}^{Hudson}$ (98-
461 100% coverage for 90% CIs for both Hudson and Nei estimators; Figure 9 and Tables 5 and S6).

R code that allows users to predict df' for their own data or other scenarios, for both F_{ST}
and r^2 , is available at <https://github.com/nwfsc-cb/pseudorep>.

462 .

463

464 4 | DISCUSSION

465

466 The common scenario considered in this paper involves a researcher who has collected
467 data for large numbers of genetic markers in one or a few actual populations. All real
468 populations have a single multigeneration pedigree, and a typical goal is to use genetic methods
469 to draw inferences about evolutionary process that helped shape that population pedigree. Our
470 primary interest is on quantifying uncertainty arising from two sources: sampling genes, and
471 sampling individuals. We do this by simulating many replicate populations and measuring how
472 fast variances of key genetic parameters decline as more loci are used and individuals are
473 sampled.

474 A substantial complication arises from the fact that the Wright-Fisher reproduction
475 process is inherently stochastic and has many possible realizations. As a consequence, replicate
476 WF populations have different multi-generational pedigrees and different mean values for
477 genetic indices like r^2 and F_{ST} (Cockerham and Weir 1983; Waples and Faulkner 2009), and
478 averaging across this sort of demographic variance would inflate our estimates of $\text{var}(\hat{F}_{ST})$ and
479 $\text{var}(\hat{r}^2)$ and bias results. To avoid this complication, we used WF reproduction for the forward-
480 in-time component of our simulations, but we took advantage of recently-developed methods
481 (Haller et al. 2019) to generate many replicate samples of genes and individuals drawn from the
482 same population pedigree. This eliminated the demographic component of variance so we could
483 focus on uncertainty associated with sampling of genes and individuals from a single population

484 (or pair of populations for F_{ST}). Collectively, our results demonstrate the importance of
485 accounting for differences between pedigrees of the population as a whole and pedigrees of
486 sampled individuals. The variance components analysis showed that for both r^2 and F_{ST} , the
487 primary factor limiting precision in genomics-scale datasets is uncertainty associated with
488 sampling individuals, and this uncertainty cannot be eliminated by sampling arbitrarily large
489 numbers of genes for the same individuals. This is an important consideration for experimental
490 design in genomics-scale datasets, as often it is faster, easier, and cheaper to assay large numbers
491 of genes for a small number of individuals, rather than the reverse.

492 Effects of pedigrees on statistical inference have been noted previously (e.g., Laurie and
493 Weir 2003; Wakeley et al. 2012, 2016; Ralph 2019). The standard coalescent treats every
494 independent gene as if it were produced on a different pedigree, but in real populations all genes
495 have to percolate through the single, fixed pedigree that captures the genealogical relationships
496 among individuals in the current population and their ancestors, and this pedigree-dependence
497 creates correlations among alleles at different gene loci, even across chromosomes (Bhaskar and
498 Song 2009). This effect is strongest for analyses that are sensitive to pedigrees from the most
499 recent generations, such as relatedness, admixture, population differentiation, and LD (King et
500 al. 2018; Nelson et al. 2020). Furthermore, departures from the standard coalescent model
501 increase when sample size is more than a small fraction of effective size (Bhaskar et al. 2014), a
502 condition commonly encountered in real-world applications.

503

504 **4.1 | Linkage Disequilibrium**

505 Our simulations show that, except in populations with large N_e and large genomes, once a
506 few thousand diallelic loci are used to estimate multilocus r^2 , adding more loci does little to
507 further reduce $\text{var}(\hat{r}^2)$. As a consequence, for datasets with 10^4 or more SNPs, n' can be many
508 orders of magnitude smaller than the number of pairs of loci. Put another way, except when the
509 entire population was sampled, we never estimated L' for LD to be as high as 700 effective loci,
510 even using $L = 50\text{K}$ SNPs for the largest finite N_e (3200) for the largest genome size (64
511 chromosomes) that we modeled (Table S2). This in turn means that true confidence intervals for
512 \hat{N}_e are much wider than they would be if all the pairwise comparisons were independent. The
513 modeling results find significant effects of N_e , C , S and their interactions on n' , with N_e having
514 the strongest influence.

515 The fact that n' depends heavily on N_e even for non-syntenic loci indicates that physical
516 linkage is not the major factor creating lack of independence of pairwise r^2 values; instead, most
517 pseudoreplication arises from overlapping pairs of the same loci in multiple pairwise
518 comparisons. Surprisingly (but conveniently), n' differs very little whether all pairwise
519 comparisons are used or only those on different chromosomes. Restricting comparisons to non-
520 syntenic loci reduces the number of locus pairs (n) but simultaneously increases n'/n , and the two
521 factors effectively offset each other.

522 It is somewhat ironic that the degree of physical linkage has relatively little effect on
523 pseudoreplication in analyses of LD, but this result can be understood when one considers what
524 happens as more loci are packed into a fixed number of chromosomes. Any new locus will
525 inevitably be linked with many existing loci on whatever chromosome it ends up on, but for
526 every such pairing there are $(C-1)/C$ new comparisons with existing loci on other chromosomes,
527 and these dilute any effects of linkage. Here we quantify for the first time the large amount of
528 pseudoreplication that occurs because each locus appears in many pairwise comparisons. This
529 overlapping-pairs-of-loci effect is modulated through the population pedigree and is strongest
530 when N_e is small.

531 Values of n' estimated by the Jones et al. (2016) jackknife over individuals were close to,
532 but on average slightly smaller than, overall n' we calculated from simulated data. For replicate
533 datasets simulated using the same fixed parameters, jackknife estimates of n' varied widely and
534 were negatively correlated with mean r^2 , which produced relatively narrow CIs when N_e was
535 estimated to be relatively large. Collectively, these features cause results of the jackknife
536 method to be somewhat unpredictable, being on average slightly conservative for larger sample
537 sizes and the opposite for $S=25$, with the latter producing a large excess of CIs that were higher
538 than true N_e . Overall coverage of CIs based on n' values estimated by our modeling results was
539 close to the target 90%. However, because of the strong positive correlation between \hat{n}' and N_e ,
540 our modeling results also produce CIs that are relatively narrow when \hat{N}_e is relatively high, an
541 effect that can be reduced by setting an upper limit to \hat{N}_e when estimating n' .

542

543 4.2 | F_{ST}

544 Consequences of pseudoreplication for precision of \hat{F}_{ST} are less dramatic than those for
545 r^2 , but nonetheless not trivial. L' approaches an asymptote for relatively small genomes and

546 small N_e , but for much larger numbers of loci (~10-20K) than is the case for r^2 . For relatively
547 large effective and genome sizes, L' was still increasing after 200K loci, indicating that in those
548 circumstances precision of \hat{F}_{ST} can be enhanced by very large numbers of loci.

549 Increasing C from 1 to 4 increased L' more than any comparable increases in N_e .
550 However, most higher organisms have $C > 4$ (Table 1), in which case comparable increases in
551 genome size and effective size produce roughly similar increases in precision. Relatively
552 speaking, increases in sample size have somewhat more effect on L' for \hat{F}_{ST} than they do for n'
553 for \hat{r}^2 . Although the ascertainment and estimation methods we evaluated affect mean values of
554 the index, these differences do not affect the rate at which $\text{var}(\hat{F}_{ST(L)}^{Nei})$ or $\text{var}(\hat{F}_{ST(L)}^{Hudson})$ declines
555 with addition of more loci. In related evaluations, we found that changes in the variance of
556 temporal F (Nei and Tajima 1981; Jorde and Ryman 2007) parallel those for F_{ST} (unpublished
557 data). This means that results obtained here can be applied broadly to predict realized precision
558 of temporal F statistics and related measures in genomics-scale datasets.

559 We found that widely-used block-jackknife methods consistently overestimate $\text{var}(\hat{F}_{ST})$,
560 leading to CIs that are much too conservative. In addition, for a given parameter set, the
561 estimated jackknife variance varied several-fold across replicates and was positively correlated
562 with the mean. According to Busing et al. (1999), block size should be large enough to
563 encompass all non-independence, which suggests the appropriate block size should be one
564 chromosome. Although the common block size of 5Mb might be large enough to capture
565 correlations involving sites near the center of the block, this approach arbitrarily divides a
566 continuous system of correlated loci in a way that guarantees that many tightly-linked pairs will
567 be in different blocks. Despite this apparent drawback, the undesirable attributes mentioned
568 above were less extreme for the 5Mb blocks, presumably because they produced more datapoints
569 to analyze. CIs for \hat{F}_{ST} based on L' estimated according to results of this study performed well,
570 even after accounting for uncertainty associated with estimating N_e and C .

571

572

573

574 **4.3 | Experimental Design and Practical Applications**

575 Our simulation and modeling results demonstrate that robust estimates of df' can be
576 obtained as a function of numbers of loci and individuals, genome size, and N_e . The first two

577 covariates are under control of the investigator, and the third can generally be approximated
578 reasonably well, even for non-model species. Dependence of df' on N_e introduces a
579 complication, but with even moderate amounts of genetic data one can obtain a fairly precise
580 estimate of N_e using either single-sample or two-sample (temporal) methods. Even after having
581 to estimate N_e to predict n' , LD-based confidence intervals for \hat{N}_e performed at least as well as
582 those obtained using the Jones et al. (2016) jackknife method, and with less variability among
583 replicates (Table 4; Figure S12). Why the block-jackknife method consistently overestimates
584 $\text{var}(\hat{F}_{ST})$ and produces CIs that are too wide is not clear, but it might be related to the fact that
585 blocks with arbitrary boundaries within chromosomes do not capture all dependencies among
586 loci. In any case, using our modeling results to predict L' produced robust CIs for \hat{F}_{ST} . Our
587 results should be particularly useful in planning and experimental design, as expected precision
588 for a wide range of scenarios can be evaluated quickly and easily.

589 The simulation framework we used, which combines coalescent simulations of the distant
590 past with fast and efficient Wright-Fisher forward simulations of the recent past, provides more
591 realistic results than can be achieved by either process alone (Nelson et al. 2020). Nevertheless,
592 as is inevitable our model required simplifying assumptions, so some caveats are in order. We
593 assumed closed populations and did not evaluate potential consequences of migration for $\text{var}(\hat{r}^2)$
594 or $\text{var}(\hat{F}_{ST})$. We modeled non-random associations of neutral genes and did not attempt to
595 account for correlations due to selective pressures on linked or unlinked sites, so in that respect
596 our estimates might be considered upper limits to actual df' .

597 We did not explicitly model variation in recombination rate, which is known to be
598 common both across the genome and between sexes (Ritz et al. 2017; Sardell and Kirkpatrick
599 2020), so our results reflect a generic genome-wide average. Although the genome sizes we
600 simulated (1-64 chromosomes of 50 Mb) encompassed the range of mean values reported for
601 higher organisms (Table 1), all chromosomes we modeled were the same size. Some effects of
602 unequal chromosome length can be accounted for by defining an effective number of
603 chromosomes as $C_e = 1/\sum x_i^2$, where x_i is the relative length of the i^{th} chromosome, standardized
604 such that $\sum x_i = 1$. C_e is analogous to the effective number of alleles at multiallelic loci. However,
605 C_e only deals with interactions among chromosomes and does not account for different patterns
606 of recombination within chromosomes.

607 A more general formulation that considers both intra- and inter-chromosomal effects on

608 genetic shuffling was proposed by Veller et al. (2019), who defined a metric \bar{r} , which is “the
609 probability that the alleles at two randomly chosen loci are shuffled in the production of a
610 gamete” (p. 1660). In Veller et al.’s framework, the expected value of \bar{r} can be expressed as the
611 sum of two terms: an intra-chromosomal term that is the probability that two loci are on the same
612 chromosome and shuffle their alleles; and an inter-chromosomal term that is the probability that
613 two loci are on separate chromosomes and shuffle their alleles.

614 As shown in Supporting Information, the effective number of chromosomes accounts
615 for the inter-chromosomal effect on \bar{r} , which quantifies effects of independent assortment
616 (Mendel’s Second Law), and the distribution of chromosome sizes also affects the intra-
617 chromosomal effect. For any organism with more than a few chromosomes (mean for
618 vertebrates is 25; Table 1), the inter-chromosomal effect on \bar{r} greatly exceeds that of patterns of
619 recombination within chromosomes. Therefore, use of C_e to account for effects of unequal
620 chromosome length should provide a good first-order approximation for the overall amount of
621 genetic shuffling and hence pseudoreplication.

622 Nevertheless, two additional factors contribute to the intra-chromosomal effect on \bar{r} :
623 1) the number of crossovers (COs) on each chromosome, and 2) their locations. All else being
624 equal, more COs lead to more shuffling, and COs near the center of a chromosome lead to more
625 shuffling than crossovers near the ends (Veller et al. 2019). Within chromosomes, our
626 simulations modeled the number of COs as a random Poisson variable, with locations of COs
627 randomly spaced along the chromosome. In Supporting Information, we show (Equations S14-
628 S15) how researchers who are interested in results for species with different patterns of
629 recombination within chromosomes can adjust our results to reflect the desired overall value of \bar{r}
630 and hence the appropriate level of pseudoreplication.

631 We modeled discrete generations, and our sampling design assumed that individuals were
632 sampled randomly from the entire adult population. Forward simulations used Wright-Fisher
633 dynamics with a constant number of ideal adults (N), so $N_e \approx N$. We found a qualitative difference
634 in df' for samples that included the entire population ($S=N_e$), but for real populations (typically
635 with $N_e < N$) the relevant criterion is whether all individuals have been sampled ($S=N$). Values for
636 df' reported in Table S2 for $S=N_e$ would provide a robust estimate of expected precision for the
637 special case where it is possible to assay the entire population, thus eliminating the large variance
638 component associated with sampling individuals. Finally, for many species it is most convenient

639 to sample juvenile offspring rather than adults. The variance associated with juvenile samples
640 approximates that for a very small sample of the parents (see Supporting Information).
641 Therefore, an approximate value for df' for such samples can be obtained by using the predicted
642 df' (n' or L') for small S .

643

644 **Acknowledgments**

645 We thank Eric Anderson, Nicholas Galwey, Marty Kardos, Peter Ralph, Carl Veller, and
646 John Wakeley for useful discussions and Martin Liermann for suggestions regarding model
647 fitting. Two anonymous reviewers provided useful suggestions that improved the manuscript.

648

REFERENCES

- Aarts, E., Verhage, M., Veenvliet, J.V., Dolan, C.V. and Van Der Sluis, S., 2014. A solution to dependency: using multilevel analysis to accommodate nested data. *Nature Neuroscience*, 17, 491.
- Aguirre, N.C., Filippi, C.V., Zaina, G., Rivas, J.G., Acuña, C.V., Villalba, P.V., García, M.N., González, S., Rivarola, M., Martínez, M.C. and Puebla, A.F., 2019. Optimizing ddRADseq in non-model species: A case study in *Eucalyptus dunnii* Maiden. *Agronomy*, 9(9), p.484.
- Albrechtsen, A., Nielsen, F.C. and Nielsen, R., 2010. Ascertainment biases in SNP chips affect measures of population divergence. *Molecular Biology and Evolution*, 27(11), 2534-2547.
- Beverton, R. J. H.; Holt, S. J. (1957), On the Dynamics of Exploited Fish Populations, Fishery Investigations Series II Volume XIX, Ministry of Agriculture, Fisheries and Food.
- Bhaskar, A., Song, Y.S., 2009. Multi-locus match probability in a finite population: a fundamental difference between the Moran and Wright-Fisher models. *Bioinformatics*, 25, i187–i195.
- Bhaskar, A., Clark, A.G. and Song, Y.S., 2014. Distortion of genealogical properties when the sample is very large. *Proceedings of the National Academy of Sciences*, 111(6), pp.2385-2390.
- Bhatia, G., Patterson, N., Sankararaman, S. and Price, A.L., 2013. Estimating and interpreting F_{ST} : the impact of rare variants. *Genome research*, 23(9), pp.1514-1521.

- Busing, F.M., Meijer, E. and Van Der Leeden, R., 1999. Delete-m jackknife for unequal m. *Statistics and Computing*, 9(1), 3-8.
- Choquet, M., Smolina, I., Dhanasiri, A.K., Blanco-Bercial, L., Kopp, M., Jueterbock, A., Sundaram, A.Y. and Hoarau, G., 2019. Towards population genomics in non-model species with large genomes: a case study of the marine zooplankton *Calanus finmarchicus*. *Royal Society open science*, 6(2), p.180608.
- Cockerham, C. C. 1969. Variance of gene frequencies. *Evolution* 23:72-84.
- Cockerham, C.C. and Weir, B.S., 1983. Variance of actual inbreeding. *Theoretical population biology*, 23(1), pp.85-109.
- Cockerham, C.C. and Weir, B.S., 1987. Correlations, descent measures: drift with migration and mutation. *Proceedings of the National Academy of Sciences*, 84(23), pp.8512-8514.
- Colegrave, N. and Ruxton, G.D., 2018. Using biological insight and pragmatism when thinking about pseudoreplication. *Trends in ecology & evolution*, 33(1), pp.28-35.
- Cox, D.R., 1984. Effective degrees of freedom and the likelihood ratio test. *Biometrika*, 71(3), pp.487-493.
- da Fonseca RR, Albrechtsen A, Themudo GE, Ramos-Madrigal J, Sibbesen JA, Maretty L, Zepeda-Mendoza ML, Campos PF, Heller R, Pereira RJ. 2016. Next-generation biology: Sequencing and data analysis approaches for non-model organisms. *Marine Genomics* 30, 3-13.
- Foll, M. and Gaggiotti, O.E., 2008. A genome scan method to identify selected loci appropriate for both dominant and codominant markers: a Bayesian perspective. *Genetics* 180:977:993.
- Galwey NW. 2009. A new measure of the effective number of tests, a practical tool for comparing families of non-independent significance tests. *Genetic Epidemiology* 33(7):559-68.
- Gautier M, Gharbi K, Cezard T, Foucaud J, Kerdelhué C, Pudlo P, Cornuet JM, Estoup A. 2013. The effect of RAD allele dropout on the estimation of genetic variation within and between populations. *Mol Ecol*. 22:3165–3178.
- Giesbrecht, FG. 2006. Degrees of freedom, effective. *Encyclopedia of Statistical Sciences*. Available at <https://onlinelibrary.wiley.com/doi/abs/10.1002/0471667196.ess0539.pub2>.
- Graffelman J, Nelson S, Gogarten SM, Weir BS. 2015. Exact Inference for Hardy-Weinberg Proportions with Missing Genotypes: Single and Multiple Imputation. *G3: Genes| Genomes|*

- Genetics 5(11):2365-73.
- Haller, B.C., Galloway, J., Kelleher, J., Messer, P.W. and Ralph, P.L., 2019. Tree-sequence recording in SLiM opens new horizons for forward-time simulation of whole genomes. *Molecular ecology resources*, 19(2), pp.552-566.
- Hill, W.G. 1981. Estimation of effective population size from data on linkage disequilibrium. *Genetical Research* 38:209–216.
- Huang H, Knowles LL. Unforeseen consequences of excluding missing data from next-generation sequences: simulation study of RAD sequences. 2014. *Systematic Biology* 4:syu046.
- Hudson, R.R., Slatkin, M. and Maddison, W.P., 1992. Estimation of levels of gene flow from DNA sequence data. *Genetics*, 132(2), 583-589.
- Hurlbert, S.H., 1984. Pseudoreplication and the design of ecological field experiments. *Ecological Monographs*, 54, 187-211.
- Jones AT, Ovenden JR, Wang Y-G. 2016. Improved confidence intervals for the linkage disequilibrium method for estimating effective population size *Heredity* 117(4), 217-223.
- Jorde, P.E., and N. Ryman. 2007. Unbiased estimator for genetic drift and effective population size. *Genetics* 177: 927–935.
- Kelleher J, AM Etheridge, G McVean (2016), *Efficient Coalescent Simulation and Genealogical Analysis for Large Sample Sizes*, *PLoS Comput Biol* 12(5): e1004842. doi: 10.1371/journal.pcbi.1004842
- King, L., Wakeley, J. and Carmi, S., 2018. A non-zero variance of Tajima’s estimator for two sequences even for infinitely many unlinked loci. *Theoretical Population Biology* 122:22-29.
- Laurie, C., Weir, B.S., 2003. Dependency effects in multi-locus match probabilities. *Theor. Popul. Biol.* 63, 207–219.
- Leaché AD, Banbury BL, Felsenstein J, de Oca AN, Stamatakis A. 2015. Short tree, long tree, right tree, wrong tree: new acquisition bias corrections for inferring SNP phylogenies. *Systematic Biology*, 29:syv053.
- Lewontin, R.C. and Krakauer, J., 1973. Distribution of gene frequency as a test of the theory of the selective neutrality of polymorphisms. *Genetics*, 74(1), 175-195.
- Li, X., Zhu, C., Lin, Z., Wu, Y., Zhang, D., Bai, G., Song, W., Ma, J., Muehlbauer, G.J.,

- Scanlon, M.J. and Zhang, M., 2011. Chromosome size in diploid eukaryotic species centers on the average length with a conserved boundary. *Molecular biology and evolution*, 28(6), pp.1901-1911.
- Lin, Y., Ghazanfar, S., Wang, K.Y., Gagnon-Bartsch, J.A., Lo, K.K., Su, X., Han, Z.G., Ormerod, J.T., Speed, T.P., Yang, P. and Yang, J.Y.H., 2019. scMerge leverages factor analysis, stable expression, and pseudoreplication to merge multiple single-cell RNA-seq datasets. *Proceedings of the National Academy of Sciences*, 116(20), pp.9775-9784.
- Messer P.W. 2013. SLiM: simulating evolution with selection and linkage. *Genetics* 194:1037–1039.
- Minias, P., Dunn, P.O., Whittingham, L.A., Johnson, J.A. and Oyler-McCance, S.J., 2019. Evaluation of a Chicken 600K SNP genotyping array in non-model species of grouse. *Scientific reports*, 9(1), pp.1-10.
- Nei M. 1973. Analysis of gene diversity in subdivided populations. *Proc Nat Acad Sci, USA*. 70:3321-3323.
- Nei, M. and Chakravarti, A., 1977. Drift variances of F_{ST} and G_{ST} statistics obtained from a finite number of isolated populations. *Theoretical Population Biology* 11:307-325.
- Nei, M. and Tajima, F. 1981. Genetic drift and estimation of effective population size. 40 *Genetics* 98:625–640.
- Nelson, D., Kelleher, J., Ragsdale, A.P., Moreau, C., McVean, G. and Gravel, S., 2020. Accounting for long-range correlations in genome-wide simulations of large cohorts. *PLoS genetics*, 16(5), p.e1008619.
- Nyholt DR. 2004. A simple correction for multiple testing for single nucleotide polymorphisms in linkage disequilibrium with each other. *Am J Hum Genet* 74:765–769.
- Patterson, N. J., P. Moorjani, Y. Luo, S. Mallick, N. Rohland et al., 2012 Ancient admixture in human history. *Genetics* 192: 1065–1093.
- Pe'er, I., Yelensky, R., Altshuler, D. & Daly, M. J. Estimation of the multiple testing burden for genomewide association studies of nearly all common variants. *Genet. Epidemiol.* 32, 381–385 (2008).
- Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A., Bender, D., Maller, J., Sklar, P., De Bakker, P.I., Daly, M.J. and Sham, P.C., 2007. PLINK: a tool set for whole-genome association and population-based linkage analyses. *The American journal of human*

genetics, 81(3), pp.559-575.

- R Core Team (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>
- Ralph, P.L., 2019. An empirical approach to demographic inference with genomic data. *Theoretical population biology*, 127:91-101.
- Ramage, B.S., Sheil, D., Salim, H.M., Fletcher, C., Mustafa, N.Z.A., Luruthusamay, J.C., Harrison, R.D., Butod, E., Dzulkiply, A.D., Kassim, A.R. and Potts, M.D., 2013. Pseudoreplication in tropical forests and the resulting effects on biodiversity conservation. *Conservation Biology*, 27(2), 364-372.
- Reynolds, J., B. S. Weir, And C. C. Cockerham. 1983. Estimation of the coancestry coefficient: basis for a short-term genetic distance. *Genetics* 105:767-779.
- Ritz, K. R., M. A. Noor, and N. D. Singh. 2017. Variation in recombination rate: adaptive or not? *Trends in Genetics* 33:364–374.
- Rosenblum, E.B. and Novembre, J., 2007. Ascertainment bias in spatially structured populations: a case study in the eastern fence lizard. *Journal of Heredity*, 98(4), pp.331-336.
- Sardell JM, Kirkpatrick M. 2020. Sex differences in the recombination landscape. *The American Naturalist* 195:361–379.
- Thompson, EA. 2013. Identity by descent: variation in meiosis, across genomes, and in populations. *Genetics* 194:301-326.
- Van Wyngaarden, M., Snelgrove, P.V., DiBacco, C., Hamilton, L.C., Rodríguez-Ezpeleta, N., Jeffery, N.W., Stanley, R.R. and Bradbury, I.R., 2017. Identifying patterns of dispersal, connectivity and selection in the sea scallop, *Placopecten magellanicus*, using RAD seq-derived SNPs. *Evolutionary Applications*, 10(1), 102-117.
- Veller, C., Kleckner, N. and Nowak, M.A., 2019. A rigorous measure of genome-wide genetic shuffling that takes into account crossover positions and Mendel's second law. *Proceedings of the National Academy of Sciences*, 116(5), pp.1659-1668.
- Wakeley, J., King, L., Low, B.S., Ramachandran, S., 2012. Gene genealogies within a fixed pedigree, and the robustness of Kingman's coalescent. *Genetics* 190, 1433–1435.
- Wakeley, J., King, L. and Wilton, P.R., 2016. Effects of the population pedigree on genetic signatures of historical demographic events. *Proceedings of the National Academy of Sciences*, 113(29), pp.7994-8001.

- Waples RK, Larson WA, and Waples RS. 2016. Estimating contemporary effective population size in non-model species using linkage disequilibrium across thousands of loci. *Heredity* 117:233-240; doi:10.1038/hdy.2016.60
- Waples, R.S. 2006. A bias correction for estimates of effective population size based on linkage disequilibrium at unlinked gene loci. *Conservation Genetics* 7:167-184.
- Waples, R.S., and C. Do. 2008. *LDNe*: A program for estimating effective population size from data on linkage disequilibrium. *Mol. Ecol. Resources* 8:753-756.
- Waples, R.S., and J.R. Faulkner. 2009. Modeling evolutionary processes in small populations: Not as ideal as you think. *Molecular Ecology* 18:1834-1847.
- Weir, B.S. and Cockerham, C.C., 1984. Estimating F-statistics for the analysis of population structure. *Evolution*, 38(6), 1358-1370.
- Wright 1951. The genetical structure of populations. *Annals of Eugenics* 15: 323-354.
- Wolf, J.B. and Ellegren, H., 2017. Making sense of genomic islands of differentiation in light of speciation. *Nature Reviews Genetics*, 18(2), p.87.

Data Accessibility

This manuscript did not generate new empirical data. Code to conduct the simulations and related analyses is available at <https://github.com/nwfsc-cb/pseudorep>.

Author Contributions

RSW conceived the study; RSW and RKW designed the study and simulated and analyzed the data. EW conducted the analyses to model n' as a function of key covariates. RSW wrote the manuscript with input from RKW and EW, and all authors edited the manuscript.

Table 1. Summary of information on haploid ($1n$) chromosome number and mean chromosome size and genome size for major taxonomic groups (data from Li et al. 2011).

Chromosome number				Size (Mb)	
Mean	Min	Max	sd	Chromosome	Genome

Prokaryotes	2.2	2	4	0.50	2.5	5.7
Unicellular eukaryotes	16.6	2	35	9.33	1.7	23.6
Invertebrates	11.1	2	16	5.18	21.7	168.4
Vascular plants	13.4	5	20	5.66	47.7	558.0
Vertebrates	25.2	8	38	6.27	85.4	1933.1

Table 2. Notation used in this study

N	Number of individuals in each population (we assumed constant size with discrete generations)
N_e	Effective population size. Because we modeled Wright-Fisher reproduction, each generation $E(N_e) = N$; due to random demographic stochasticity, however, each replicate population has a different realized, multigeneration N_e
C	Number of chromosomes, each of length 50 Mb
S	Number of individuals sampled, drawn from the $N_e = N$ individuals in the final generation
L	Number of diallelic (SNP) loci; also the nominal degrees of freedom for F_{ST}
n	Nominal degrees of freedom = number of pairwise comparisons of loci for the LD analyses. $n = L(L-1)/2 \approx L^2/2$
$\hat{F}_{ST(L)}$	An estimate of the standardized variance of allele frequency between two populations, based on data for L diallelic loci
θ	The coancestry coefficient, which is related to F_{ST} as shown in Equation 1
\hat{F}_{ST}^{Nei}	A weighted version of Nei's (1973) estimator of F_{ST} , commonly referred to as G_{ST} . For two populations and diallelic loci, G_{ST} is identical to Wright's (1951) F_{ST} .
\hat{F}_{ST}^{Hudson}	Hudson et al.'s (1992) estimator of F_{ST} , as in Bhatia et al. (2013). For two populations and diallelic loci, \hat{F}_{ST}^{Hudson} is identical to Weir and Cockerham's (1984) $\hat{\theta}$.
$\hat{r}^2_{(L)}$	An estimate of the mean squared correlation of alleles at different loci, based on data for all pairwise comparisons of L diallelic loci

ϕ	The squared coefficient of variation of \hat{F}_{ST} or \hat{r}^2 ; $\phi = \text{var}/\text{mean}^2$
L'	Effective degrees of freedom for F_{ST} analyses; $L' = 2/\phi_F$
n'	Effective degrees of freedom for LD analyses; $n' = 2/\phi_r$
V_1	The component of $\text{var}(\hat{F}_{ST})$ or $\text{var}(\hat{r}^2)$ that reflects uncertainty associated with sampling replicate sets of genes for the same individuals
V_2	The component of $\text{var}(\hat{F}_{ST})$ or $\text{var}(\hat{r}^2)$ that reflects uncertainty associated with taking replicate samples of individuals assayed for the same genes
CI	Confidence interval

Table 3. Experimental design for sampling individuals and genes for analyses of LD. From each population, potentially-overlapping subsets of individuals (rows) were drawn from the N_e total individuals in the final generation. Subsets of L loci (columns) were non-overlapping. Mean r^2 was calculated for each cell. The variance among mean r^2 within rows was used to estimate variance component V_1 (same individuals, different loci), and the variance within columns was used to estimate variance component V_2 (same loci, different individuals). A measure of pseudoreplication, ϕ_r , was calculated across mean r^2 for the cells (in **bold**) along the diagonal (different sets of individuals and loci). This sampling design was repeated for different numbers of loci (L), sampled individuals (S), chromosome number (C), and N_e . See the text and Figure S3 for details of a modified sampling design for F_{ST} that involved pairwise comparisons of daughter populations.

		Sets of loci				
		1	2	3	4	V_1
	A	r^2_{A1}	r^2_{A2}	r^2_{A3}	r^2_{A4}	$\text{var}(r^2_{A*})$
Samples of	B	r^2_{B1}	r^2_{B2}	r^2_{B3}	r^2_{B4}	$\text{var}(r^2_{B*})$
individuals	C	r^2_{C1}	r^2_{C2}	r^2_{C3}	r^2_{C4}	$\text{var}(r^2_{C*})$

D	r^2_{D1}	r^2_{D2}	r^2_{D3}	r^2_{D4}	$\text{var}(r^2_{D*})$
V_2	$\text{var}(r^{2*}_1)$	$\text{var}(r^{2*}_2)$	$\text{var}(r^{2*}_3)$	$\text{var}(r^{2*}_4)$	

Table 4. Effects of N_e , number of chromosomes (C), and number of individuals sampled (S) on coverage of confidence intervals (CIs) around mean r^2 . Results are shown for CIs based on the Jones et al. (2016) jackknife method and effective $df(n')$ estimated using the modeling results from this study, which required estimating N_e and C . Shown are the percentages of 1024 replicate samples whose 90% CIs included the true value (“In”), were entirely above the true value (“Above”), or were entirely below (“Below”). Each cell represents results averaged over simulations with data for 500, 1000, and 5000 SNPs, and the bottom set of rows averages results across all scenarios, by sample size and method. Results shown used data for pairs of loci on different chromosomes.

N_e	C	S	Jackknife			This study		
			%Above	%Below	%In	%Above	%Below	%In
200	4	25	9.4	1.7	88.9	12.7	1.8	85.5
200	4	50	7.1	2.3	90.6	8	3.8	88.2
200	4	100	4.9	2.5	92.7	4.6	3.6	91.8
200	16	25	9.5	1.3	89.1	9.4	2.4	88.2
200	16	50	5.8	1.6	92.6	7.1	5.0	87.9
200	16	100	3.5	1.7	94.8	4.7	4.4	90.9
800	4	25	14.8	0.6	84.6	11.7	0.8	87.5
800	4	50	8.4	1.2	90.4	10.1	1.5	88.4
800	4	100	6.1	2.7	91.2	7.4	4.5	88.1
800	16	25	17.1	0.5	82.4	11.4	1.2	87.5

800	16	50	8.9	1.0	90.1	9.7	2.9	87.4
800	16	100	5	1.9	93.1	6.1	3.7	90.1

		25	12.7	1.0	86.2	11.3	1.6	87.2
	Means	50	7.5	1.5	90.9	8.7	3.3	88.0
		100	4.9	2.2	92.9	5.7	4.1	90.2

Table 5. Effects of N_e , number of chromosomes (C), and number of individuals sampled (S) on coverage of confidence intervals (CIs) around \hat{F}_{ST} . Results are shown for CIs based on block jackknife estimates of $\text{var}(\hat{F}_{ST}^{Hudson})$ and effective $df(L')$ for \hat{F}_{ST}^{Nei} estimated using the modeling results from this study, which required estimating N_e and C . Shown are the percentages of 1152 replicate samples whose 90% CIs included the true value (“In”), were entirely above the true value (“Above”), or were entirely below (“Below”). Each cell represents results averaged over simulations with data for 5000, 20000, and 50000 SNPs, and the bottom set of rows averages results across all scenarios, by sample size and method. These results used a block size of 5 Mb. Very similar results were found for block sizes of 1 chromosome and for block-jackknife estimates of $\text{var}(\hat{F}_{ST}^{Nei})$ (see Table S6).

N_e	C	S	Block jackknife			This study		
			%Above	%Below	%In	%Above	%Below	%In
200	4	25	0.7	0.6	98.7	2.4	3.9	93.7
200	4	50	0.0	0.0	100.0	2.3	2.7	94.9
200	4	100	0.0	0.0	100.0	4.0	5.4	90.6
200	16	25	0.2	0.8	99.0	2.4	3.8	93.8
200	16	50	0.2	0.2	99.6	2.9	4.0	93.1
200	16	100	0.1	0.1	99.8	3.8	4.6	91.6
800	4	25	0.5	0.3	99.2	3.0	5.6	91.3
800	4	50	0.1	0.5	99.3	4.4	6.6	89.0

800	4	100	0.1	0.1	99.8	4.7	6.4	88.9
800	16	25	0.7	1.2	98.1	3.8	4.8	91.4
800	16	50	0.5	1.2	98.4	3.7	5.7	90.6
800	16	100	0.4	0.8	98.8	5.1	5.4	89.5

		25	0.5	0.7	98.7	2.9	4.5	92.5
Means		50	0.2	0.5	99.3	3.3	4.8	91.9
		100	0.2	0.3	99.6	4.4	5.4	90.2

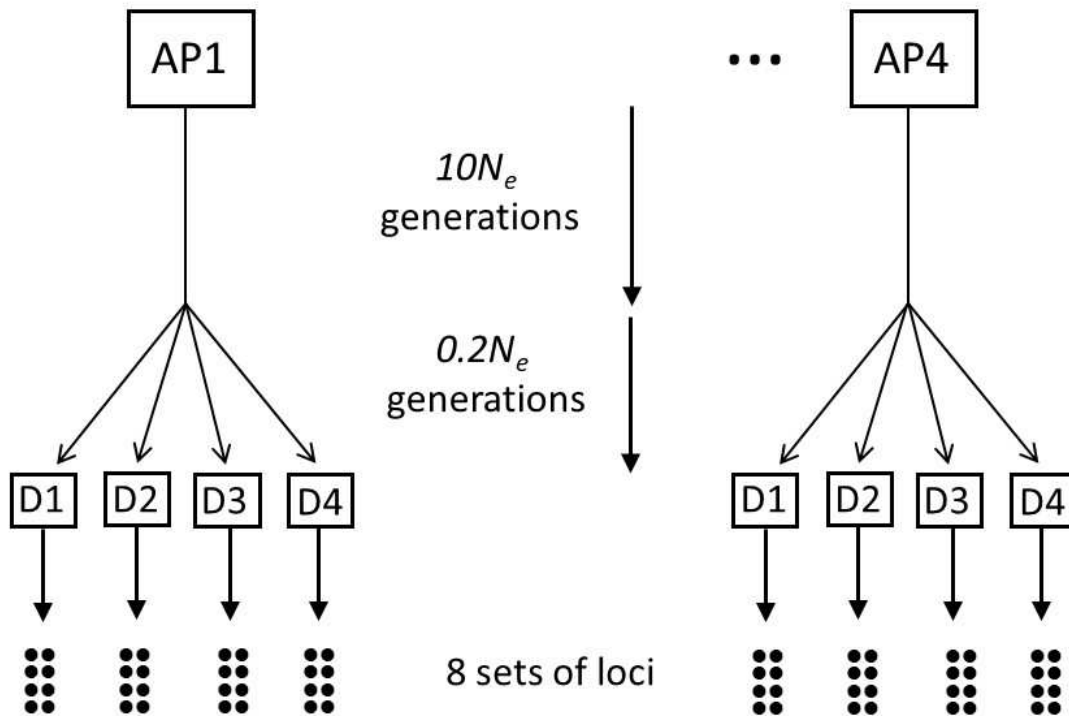


Figure 1. Experimental design for simulations. For each evolutionary scenario (combination of N_e and genome size), four ancestral populations (AP1–AP4) were simulated to ensure coalescence ($10N_e$ generations), at which point each ancestral population generated four daughter populations (D1–D4) with the same N_e . The $4 \times 4 = 16$ daughter populations then evolved independently under isolation for $t = 0.2N_e$ generations. Subsequently, the model differed slightly for the F_{ST} and LD analyses. In the latter (as depicted above), for each daughter population, eight mutational replicates (different set of loci) were generated based on the same pedigree, producing a total of 128 replicates for each evolutionary scenario. For F_{ST} , each set of four daughter populations allowed six pairwise comparisons of populations, and for each two-population pedigree six mutational replicates were generated (Figure S4).

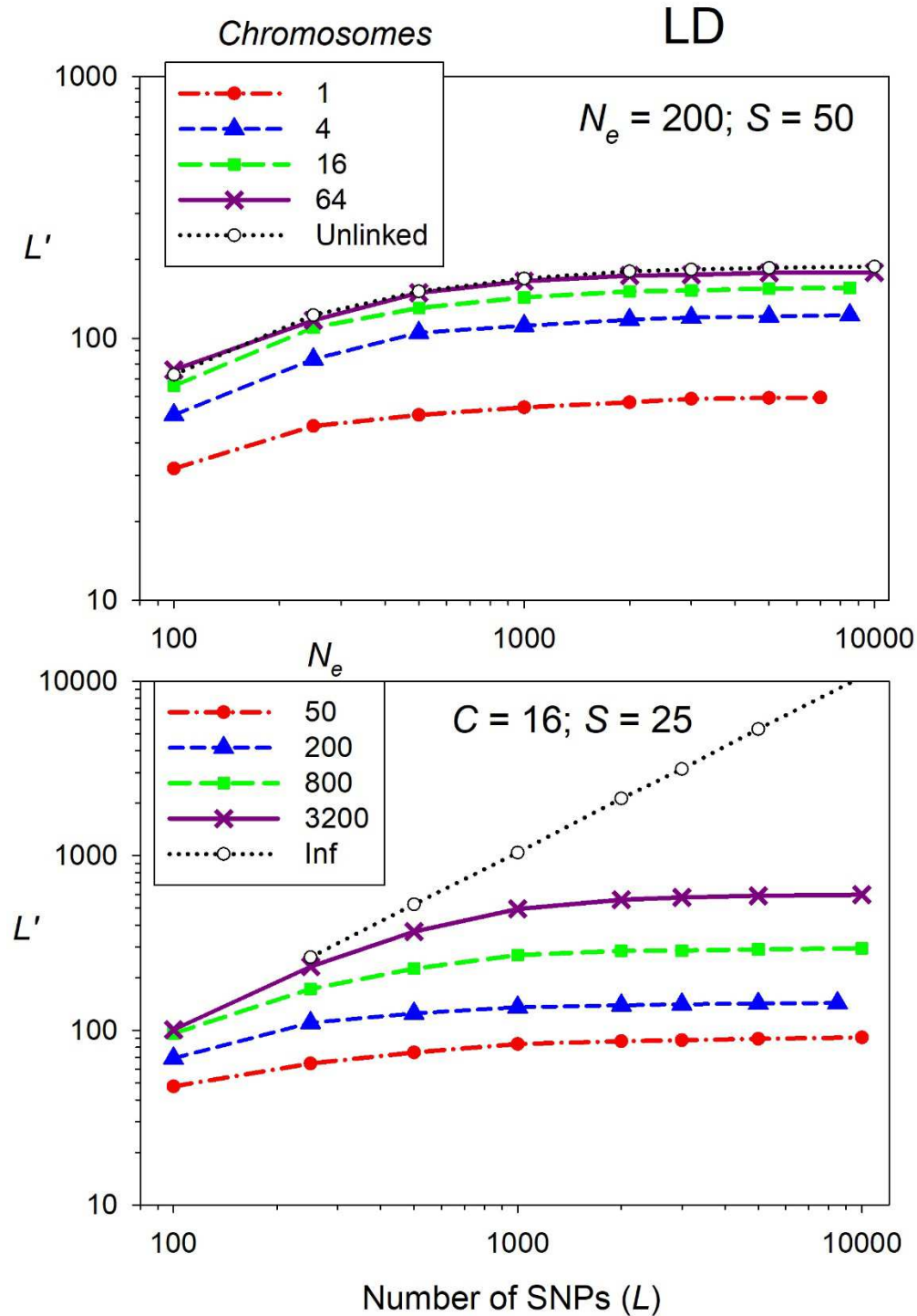


Figure 2. Effective number of loci (L') for mean r^2 as a function of the number of diallelic (SNP) loci, L . Top: Influence of number of chromosomes (C), with $N_e = 200$ and $S = 50$. Bottom: Influence of N_e , with $C = 16$ and $S = 25$. Mean r^2 was calculated across all $n = L(L-1)/2$ pairs of loci. Figure S8 (Supplementary Information) shows these same results except the Y axis

is plotted as the effective number of locus pairs (n').

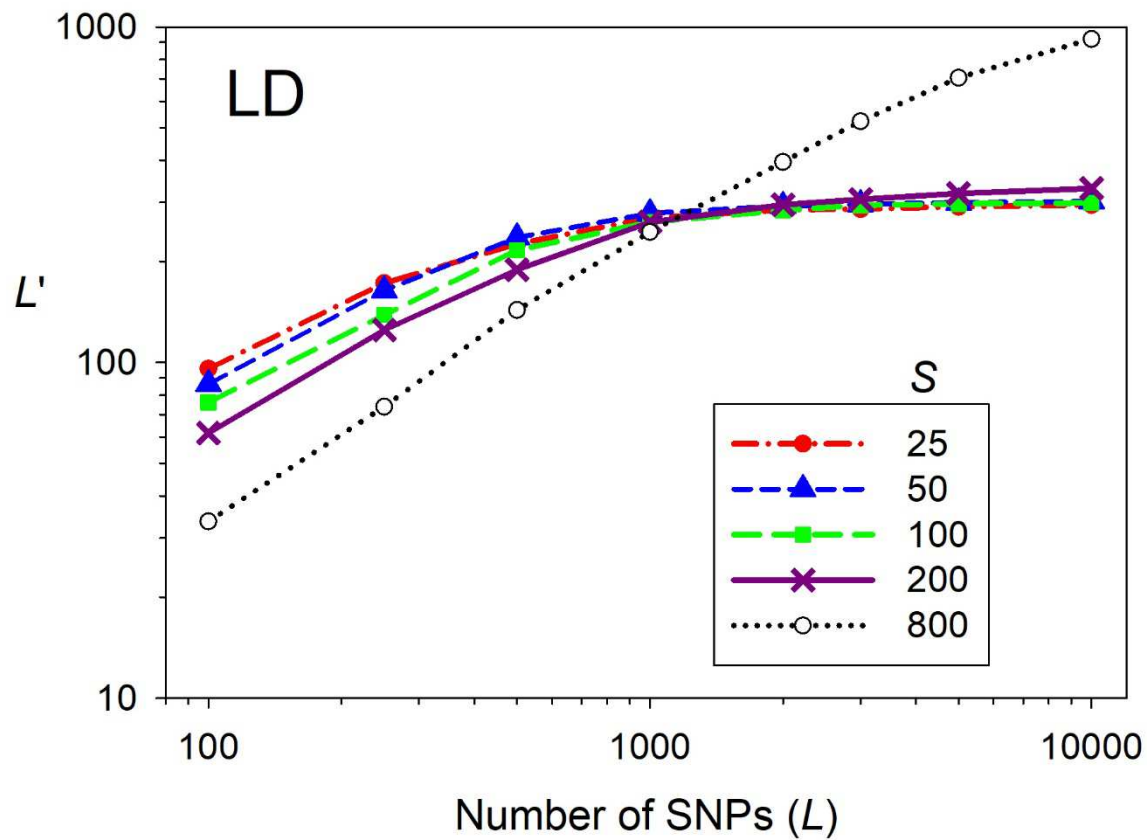


Figure 3. Effective number of loci (L') for mean r^2 as a function of the sample size of individuals ($S = 25-800$) and the number of diallelic (SNP) loci, L . Results are for $N_e = 800$, $C = 16$, and using all pairwise comparisons of loci.

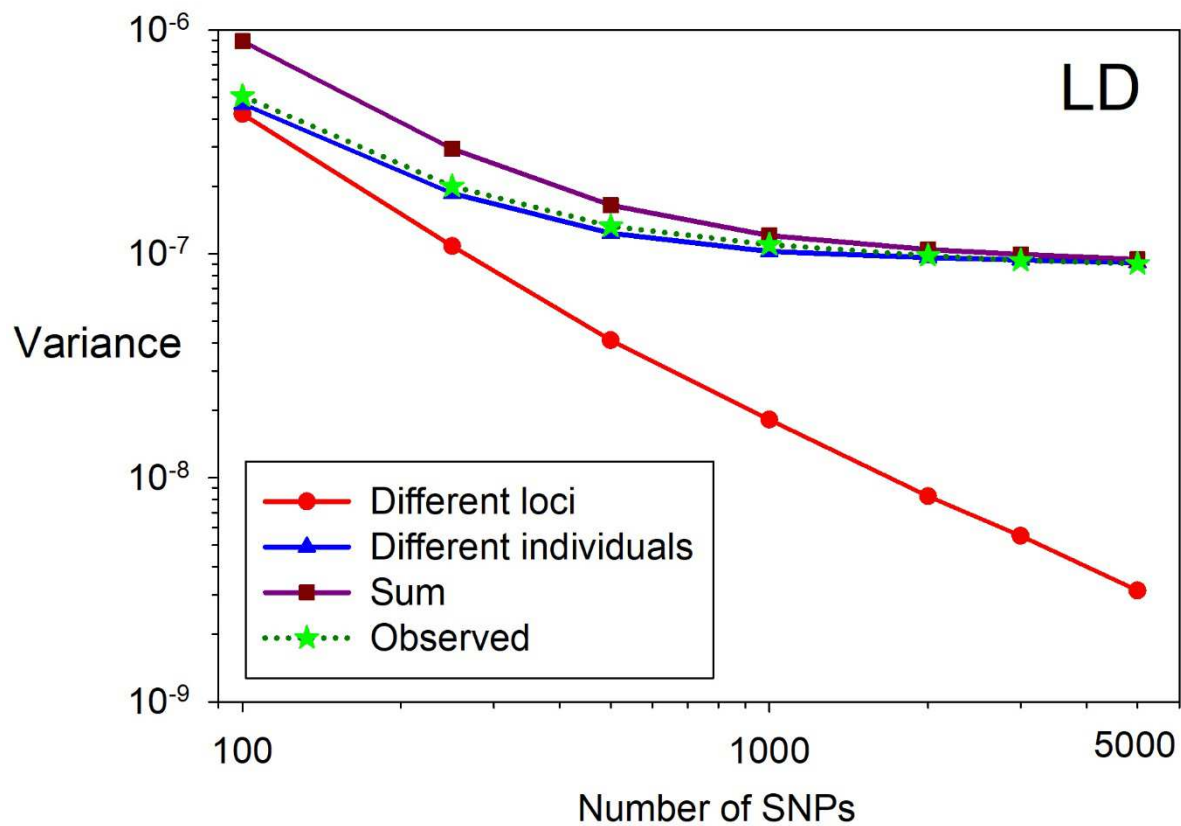


Figure 4. Variance components analysis for mean r^2 . As depicted in Table 3, V_1 is the variance of mean r^2 for the same individuals assayed for different, non-overlapping sets of loci, and V_2 is the variance of mean r^2 for different (potentially overlapping) sets of individuals assayed for the same loci. “Sum” = V_1+V_2 and “Observed” is the total observed variance of mean r^2 . Results are for $N_e = 200$, $C = 16$, $S = 50$, and using all pairwise comparisons of loci.

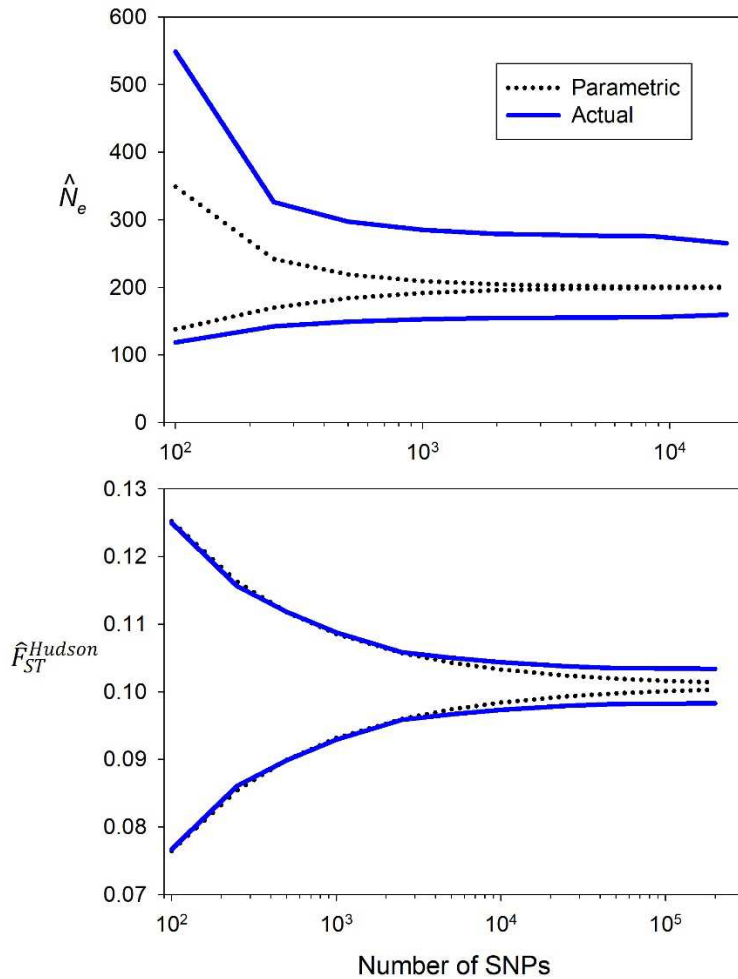


Figure 5. Comparison of parametric and actual 90% confidence intervals for \hat{N}_e based on LD (top) and $\hat{F}_{ST(L)}^{Hudson}$ (bottom). Parametric CIs use the nominal degrees of freedom ($L =$ the number of diallelic (SNP) loci for $\hat{F}_{ST(L)}^{Hudson}$; $n = L(L-1)/2$ for LD); actual CIs use the effective degrees of freedom calculated in this study (L' and n'). Results are for simulations with $N_e = 200$, $C = 16$, and $S = 50$. Note the different X-axis scales in the two panels.

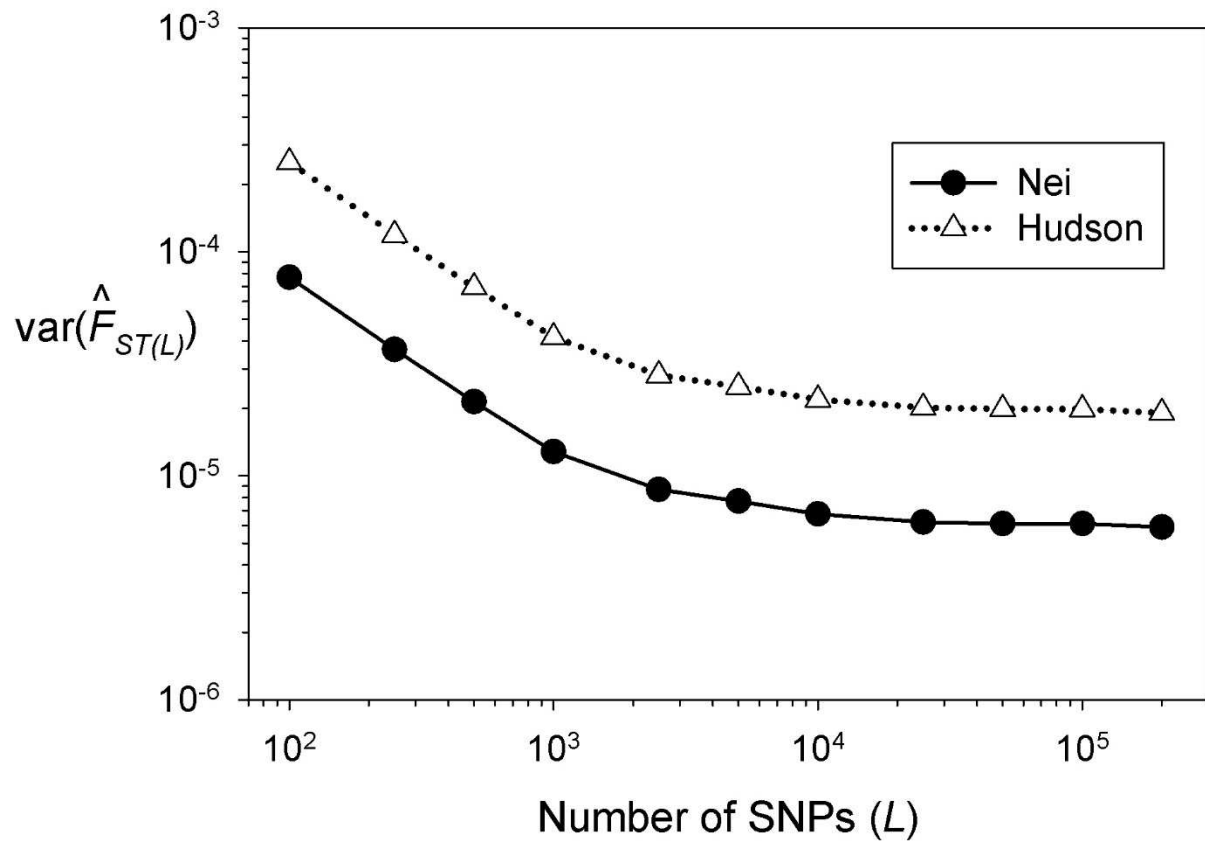


Figure 6. Rate of decline in the variance of multilocus $\hat{F}_{ST(L)}$ as more diallelic loci (SNPs, L) were used in the analysis. Results are for simulations with $N_e = 200$, $C = 4$, and $S = 25$ and are shown for the estimators of Nei (\hat{F}_{ST}^{Nei}) and Hudson (\hat{F}_{ST}^{Hudson}). Figure S15 shows comparable results for another scenario with different values of N_e , C , and S .

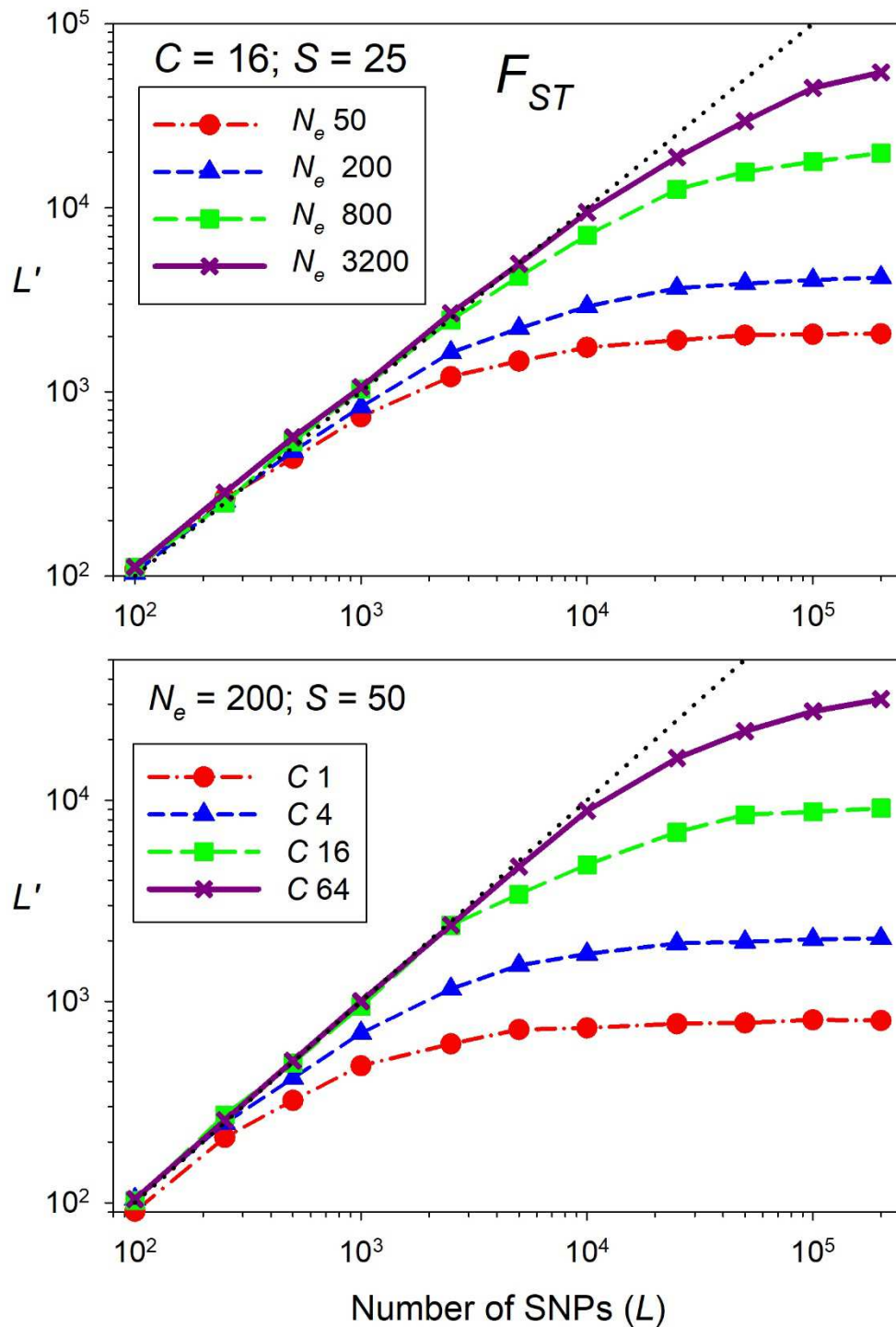


Figure 7. Influence of N_e (top panel, with number of chromosomes, C , fixed at 16) and C (bottom panel, with N_e fixed at 200) on the effective degrees of freedom (L') for $\hat{F}_{ST(L)}^{Nei}$ computed between pairs of populations. Black dotted line represents $L' = L =$ the number of SNPs.

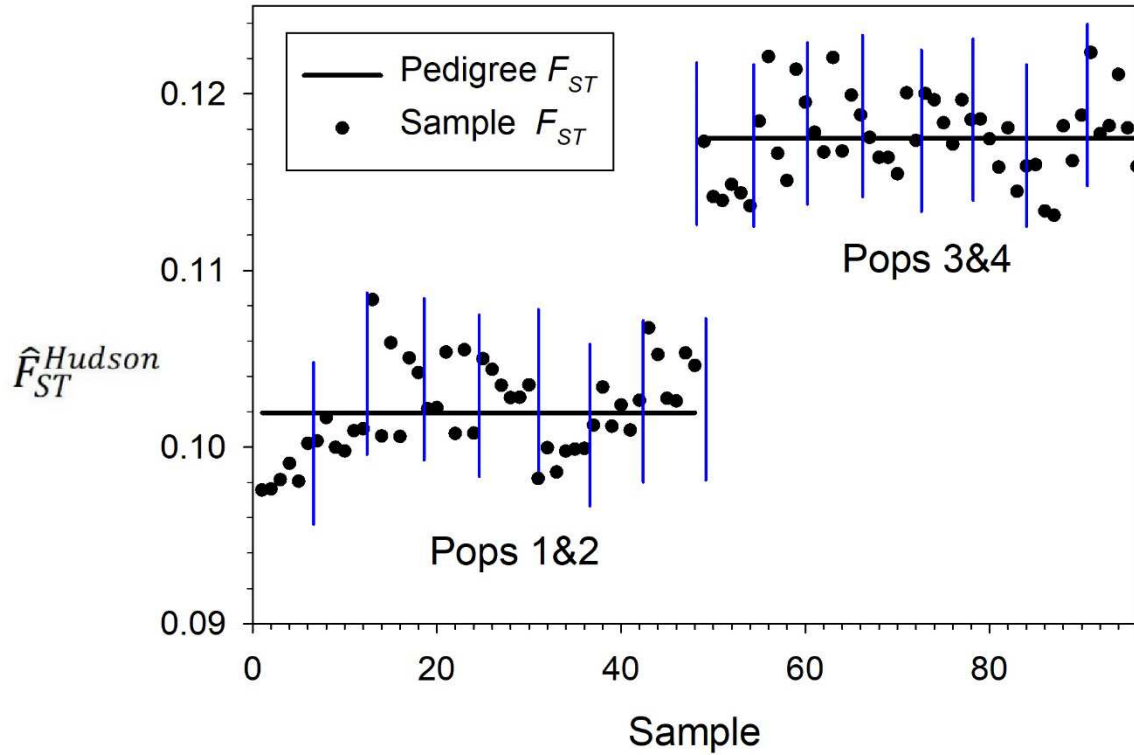


Figure 8. Effects of pedigree on variation in mean $\hat{F}_{ST(L)}^{Hudson}$. For each of two, 2-population pedigrees, 8 replicate samples (demarcated by vertical lines) were taken of $S = 100$ individuals. These results are for simulations with $N_e = 200$ and 4 chromosomes. Sampled individuals were drawn hypergeometrically from the N_e individuals in the final generation. For each sample, six mutational replicates generated non-overlapping sets of $L = 5000$ SNP loci that were used to compute mean $\hat{F}_{ST(L)}$. Solid horizontal lines (“Pedigree F_{ST} ”) represent mean \hat{F}_{ST}^{Hudson} across all $8 \times 6 = 48$ replicates within each pedigree. The first set of samples shows results for comparison of daughter populations 1 and 2 and the second set of samples shows results for comparison of daughter populations 3 and 4, all derived from the same ancestral population.

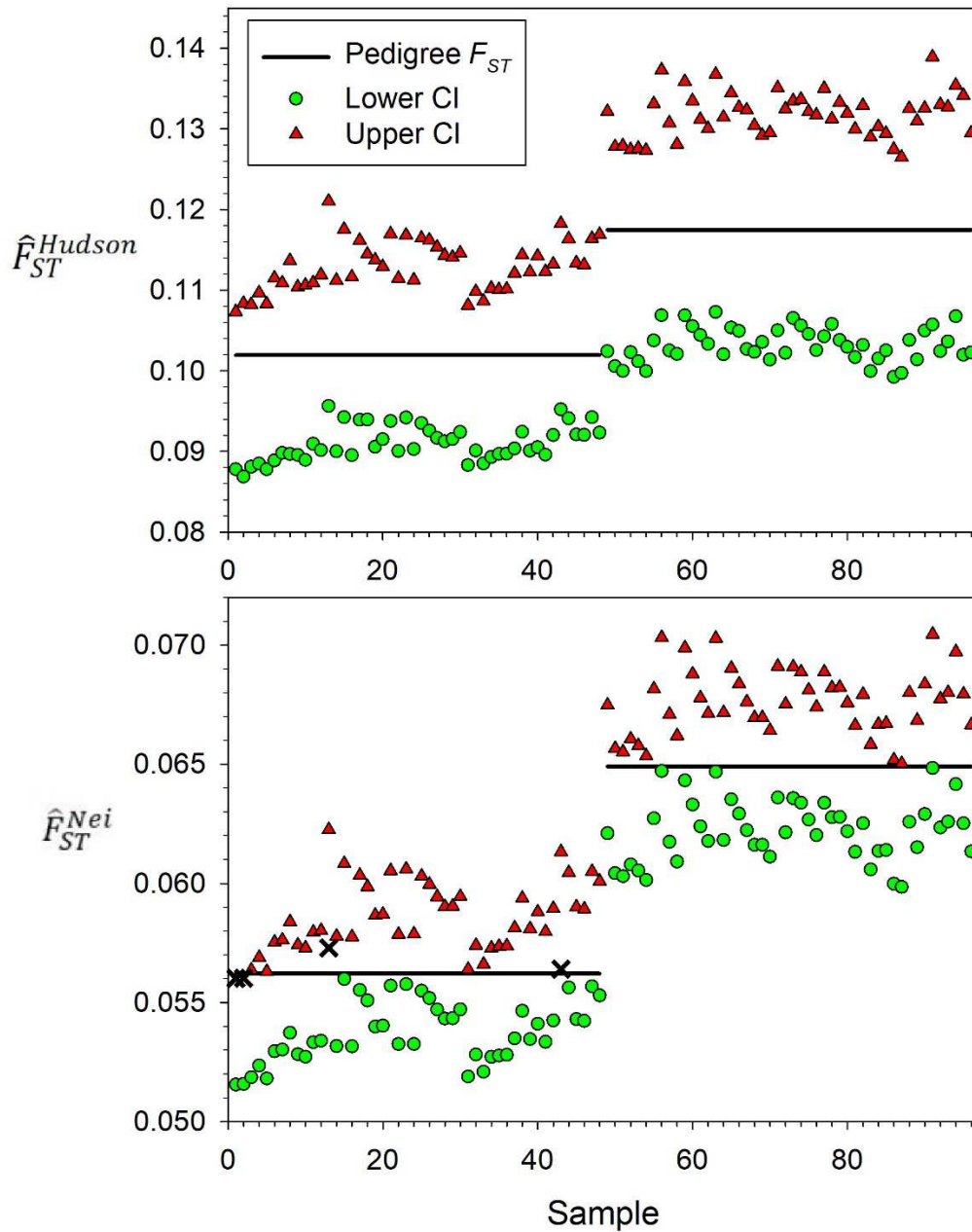
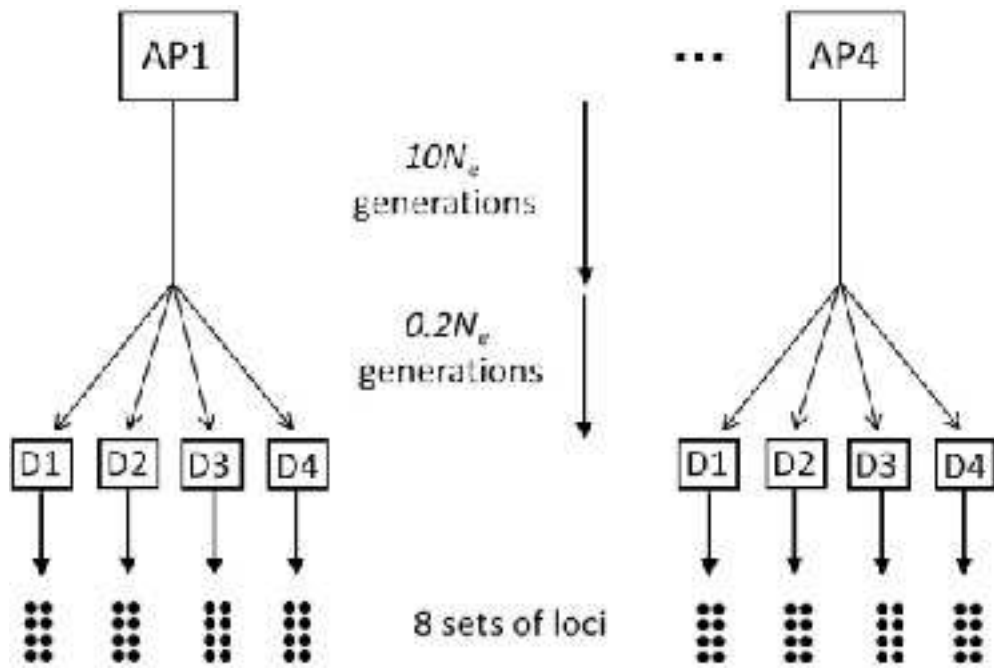
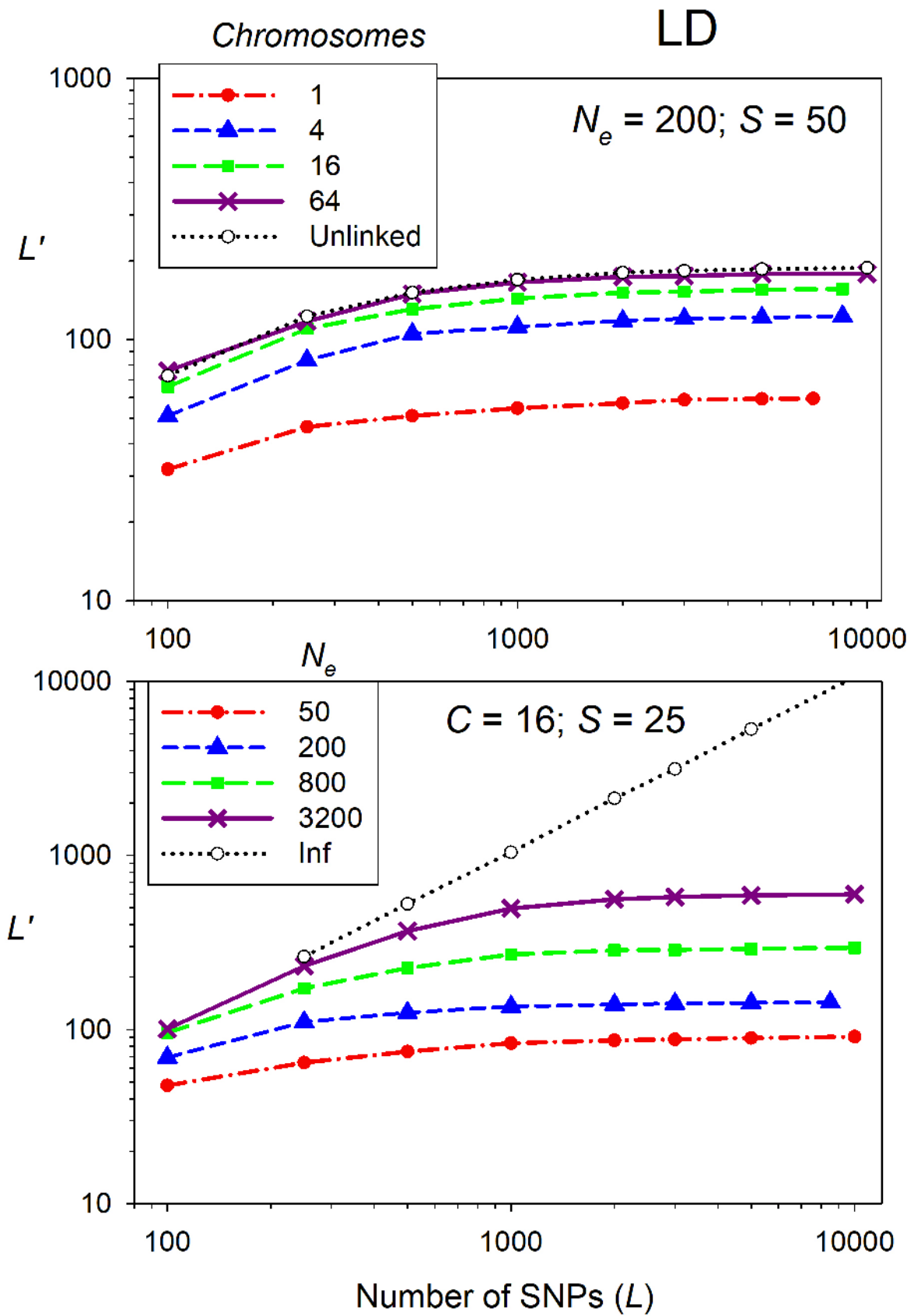


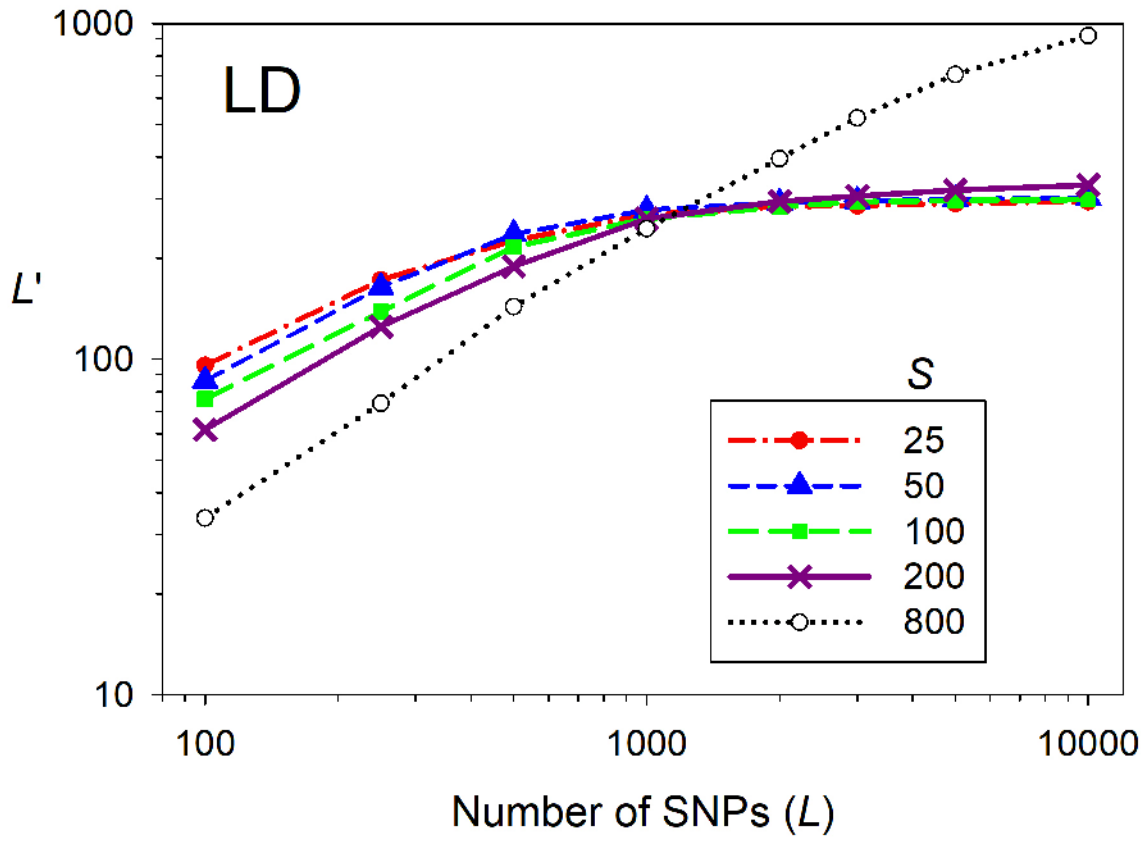
Figure 9. Coverage of 90% confidence intervals (CIs) around F_{ST} estimators for the population pedigrees and samples shown in Figure 8 ($N_e = 200$; $C = 4$; $L = 5000$; $S = 100$). Top: CIs generated from block-jackknife estimates of $\text{var}(\hat{F}_{ST}^{Nei})$. Bottom: CIs generated based on L' for \hat{F}_{ST}^{Nei} estimated from this study. CI coverage is evaluated with respect to mean \hat{F}_{ST}^{Nei} across all replicates within each pedigree (“Pedigree F_{ST} ”, horizontal lines). The black X symbols indicate an upper (or lower) bound that was below (or above) the mean pedigree F_{ST} .



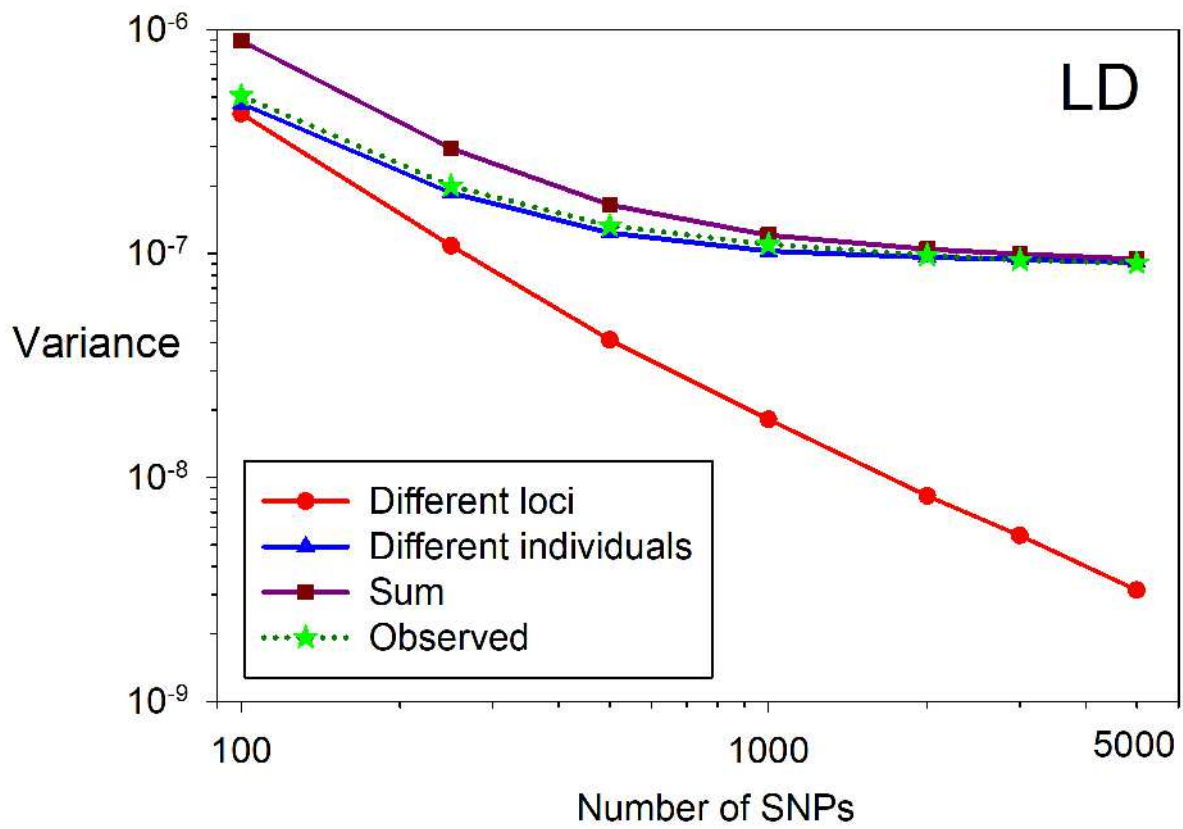
men_13482_f1.jpg



men_13482_f2.jpg

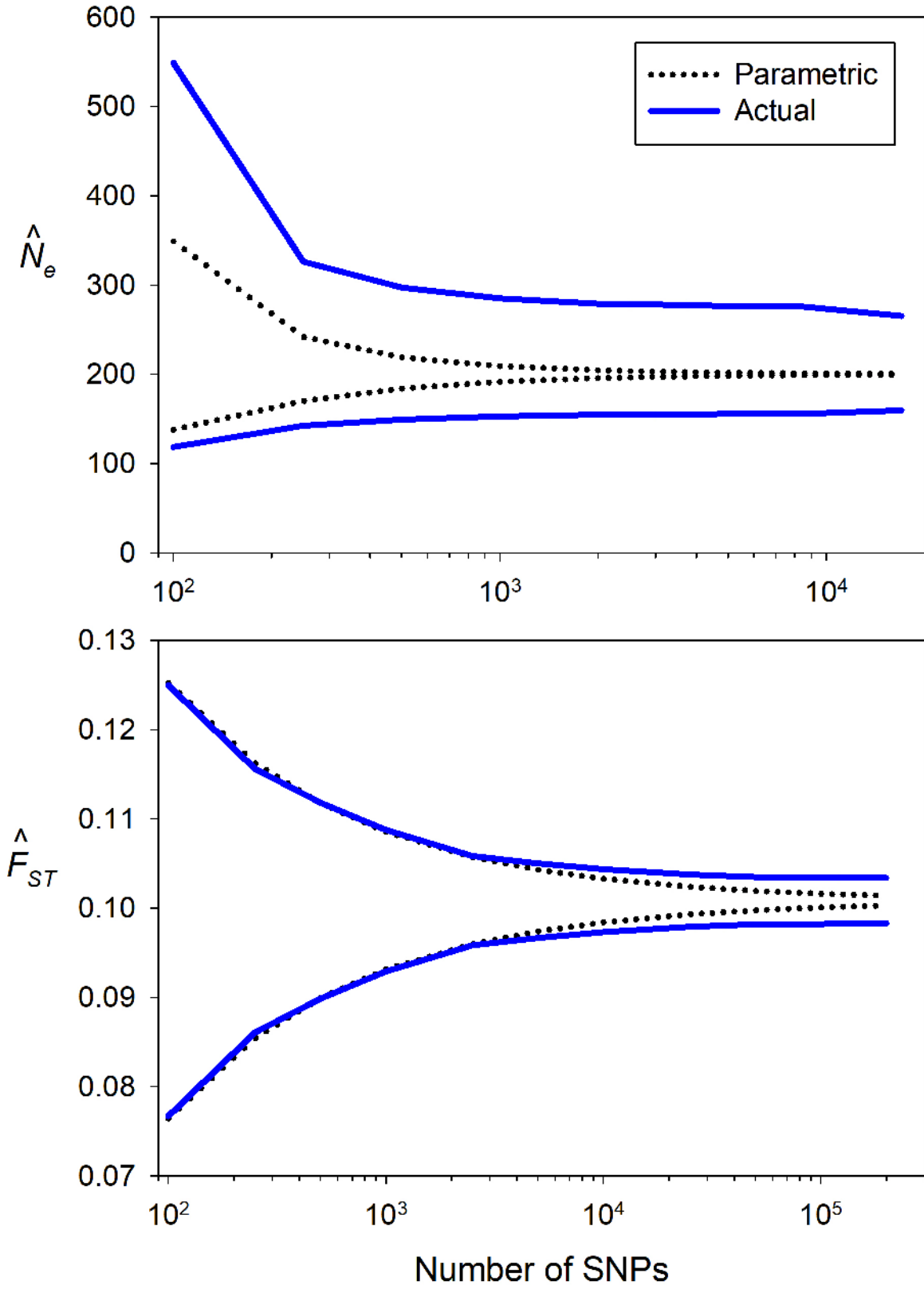


men_13482_f3.jpg

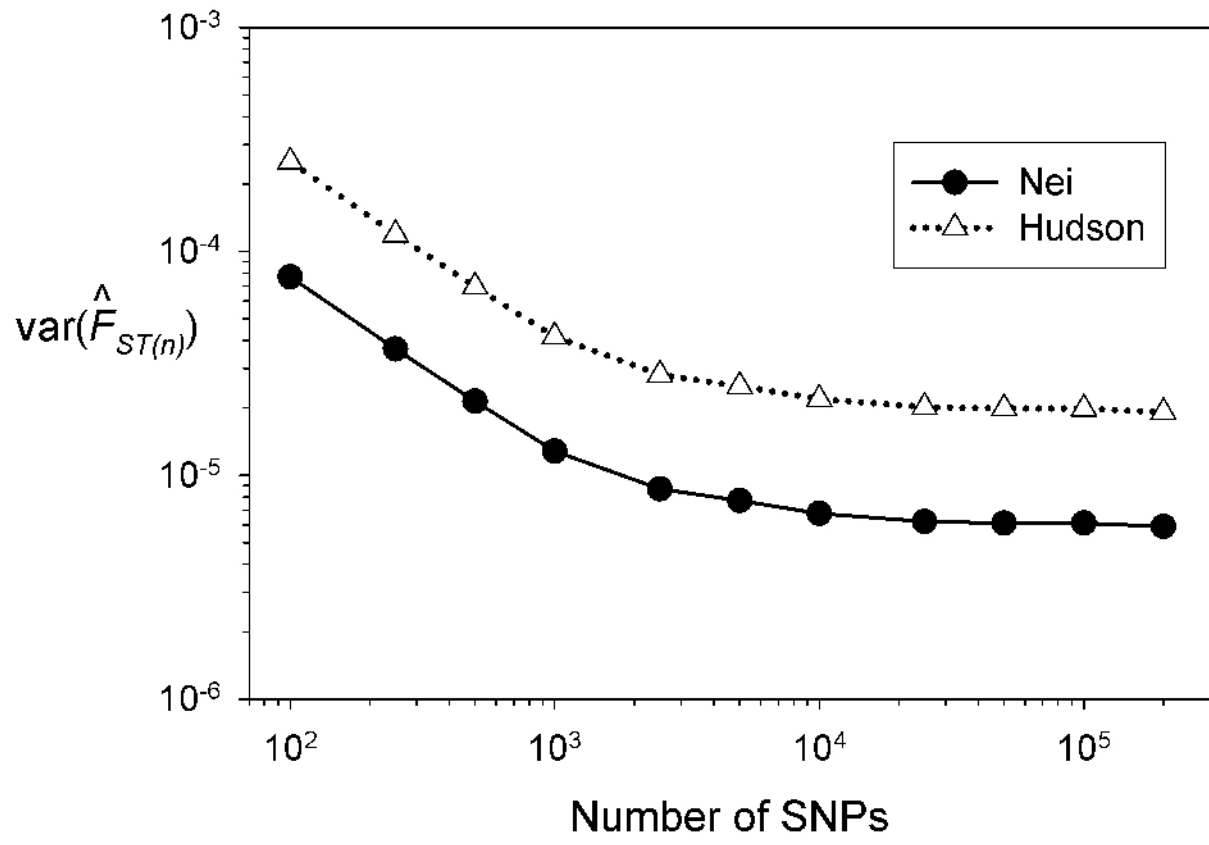


men_13482_f4.jpg

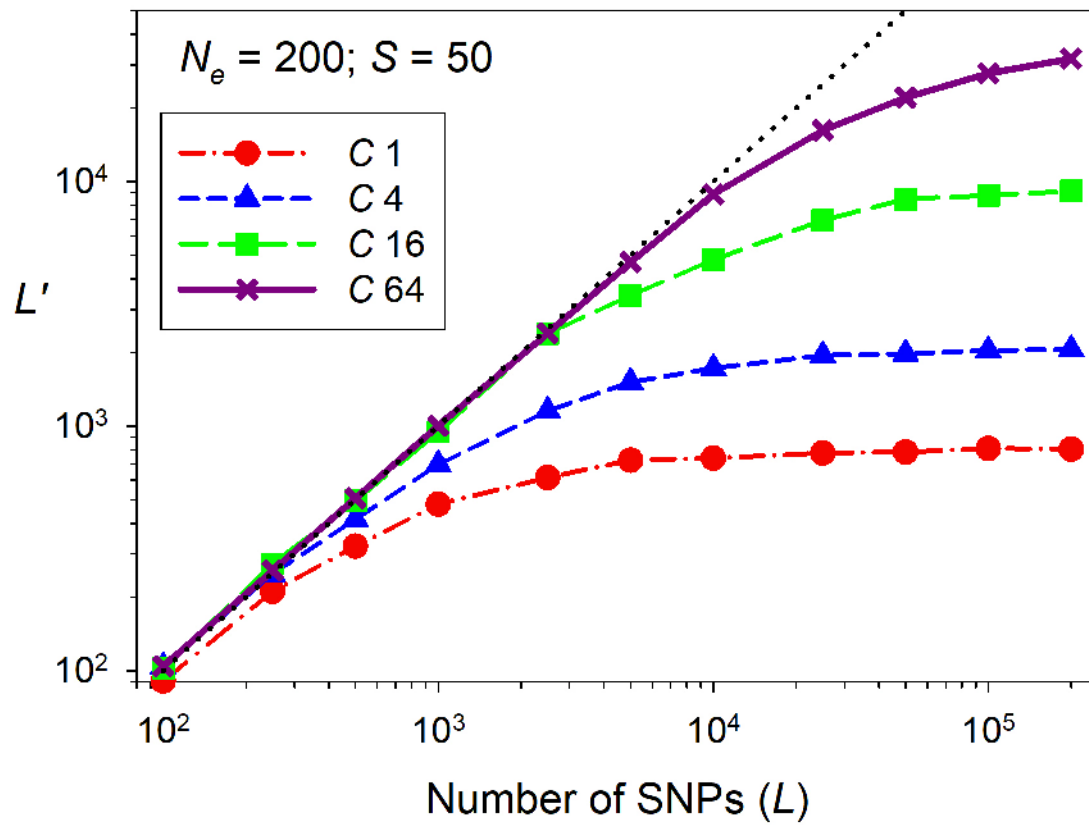
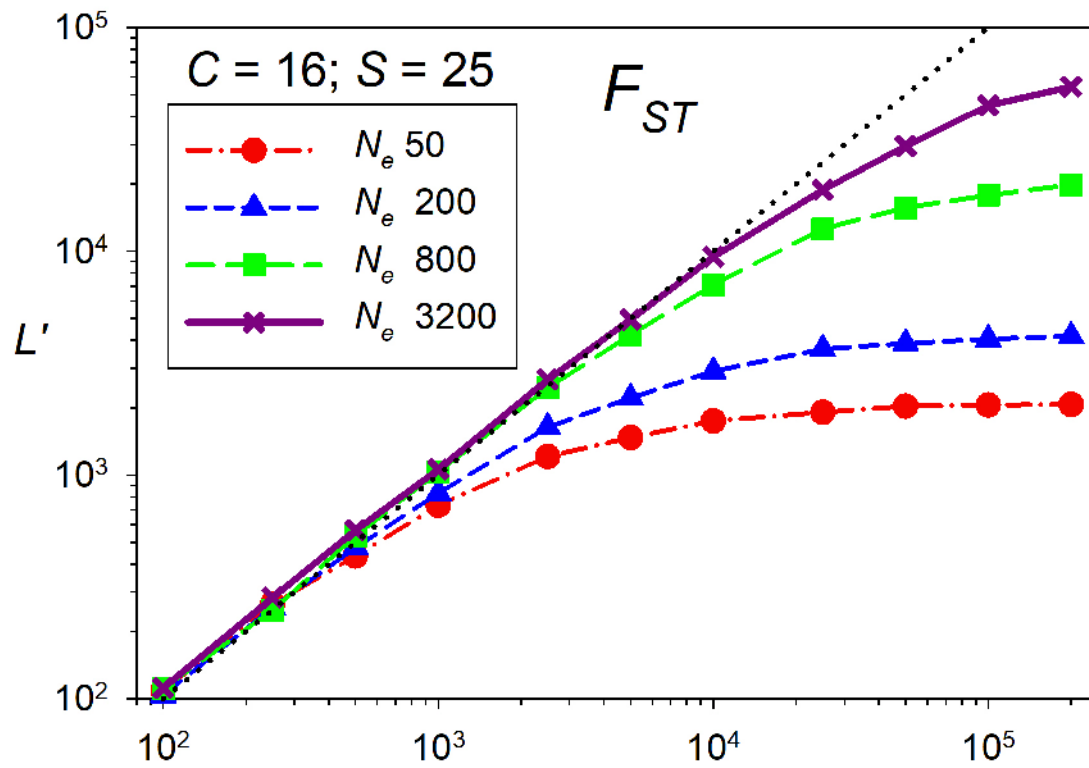
Confidence Intervals



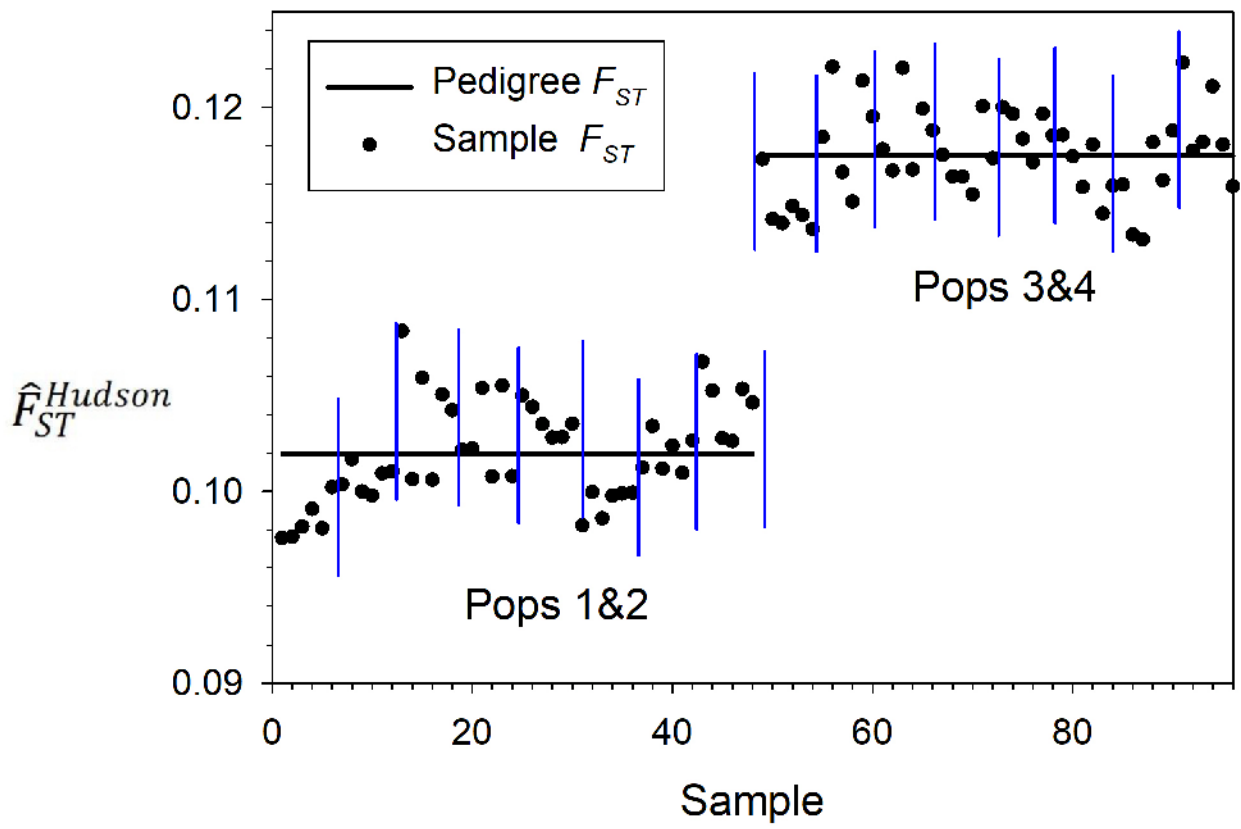
men_13482_f5.jpg



men_13482_f6.jpg

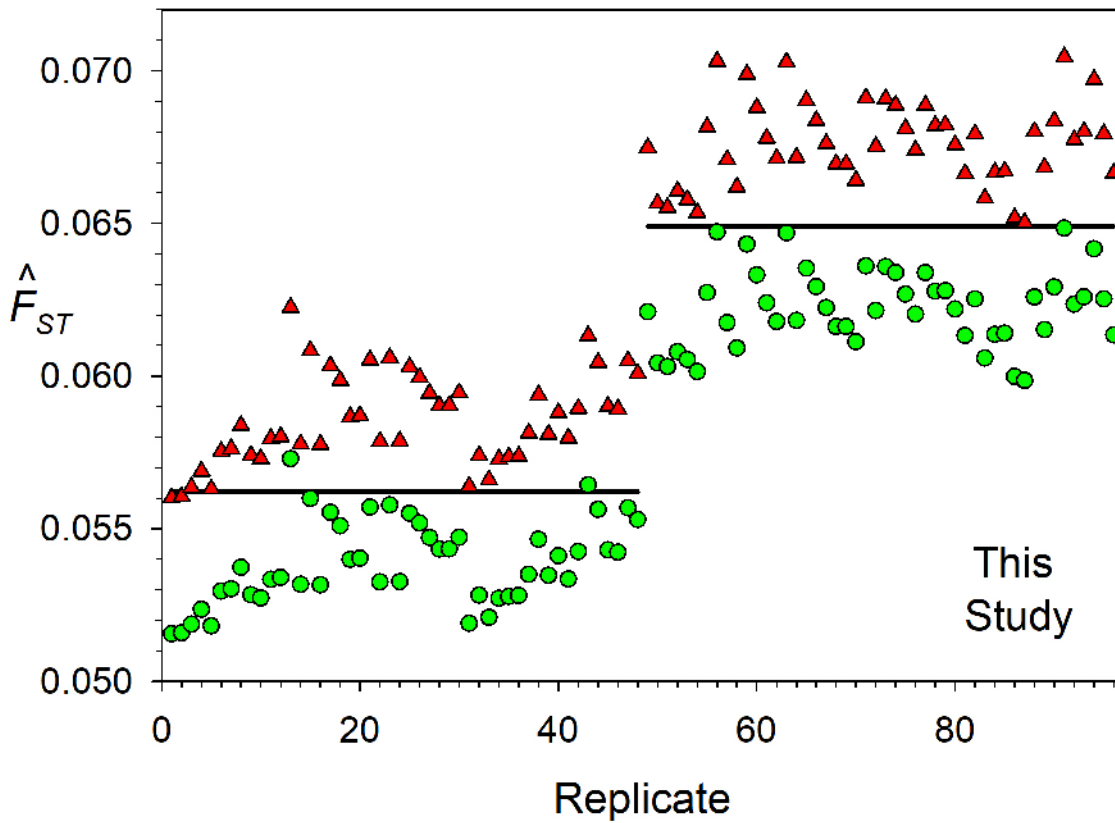
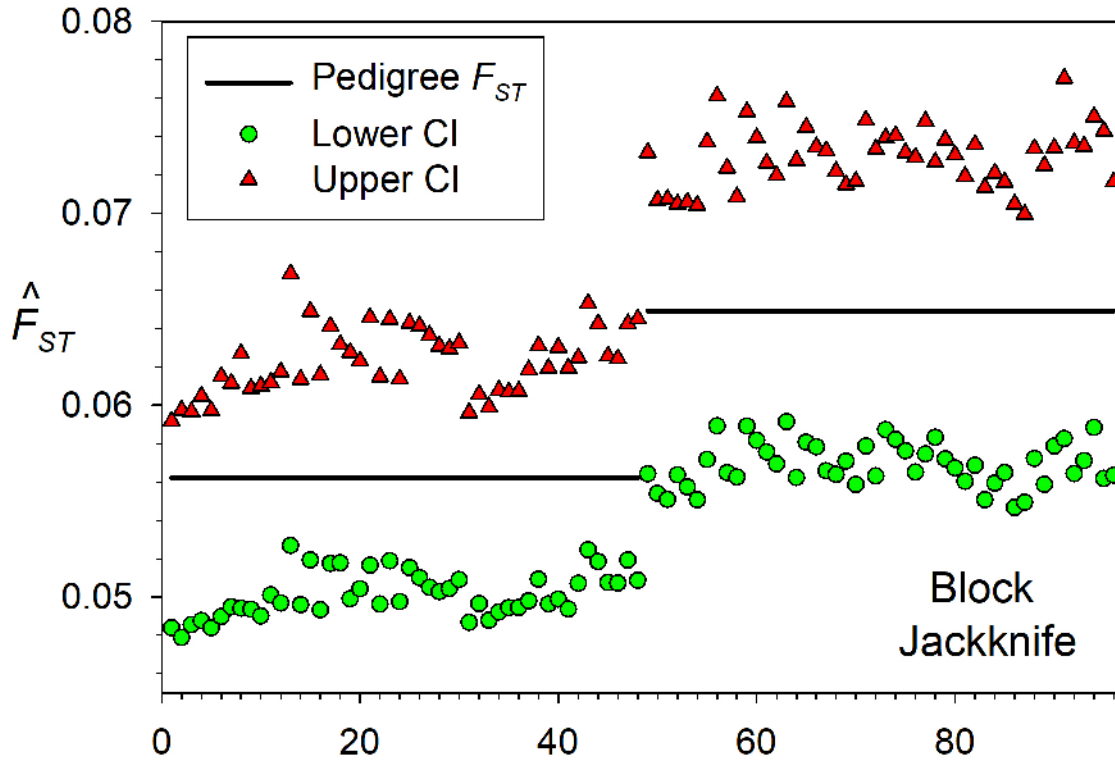


men_13482_f7.jpg



men_13482_f8.jpg

Confidence Intervals



men_13482_f9.jpg