# Does $G_{ST}$ underestimate genetic differentiation from marker data?

J. Wang

*Institute of Zoology, Zoological Society of London, London NW1 4RY, United Kingdom*

*Corresponding author*

Jinliang Wang

Institute of Zoology

Regent's Park

London NW1 4RY

United Kingdom

Tel: 0044 20 74496620

Fax: 0044 20 75862870

Email: jinliang.wang@ioz.ac.uk

24 **Abstract**

25 The widely applied genetic differentiation statistics $F_{ST}$ and $G_{ST}$ have recently been criticised

26 for underestimating differentiation when applied to highly polymorphic markers such as

27 microsatellites. New statistics claimed to be unaffected by marker polymorphisms have been

28 proposed and advocated to replace the traditional $F_{ST}$ and $G_{ST}$. This study shows that $G_{ST}$

29 gives accurate estimates and underestimates of differentiation when demographic factors are

30 more and less important than mutations, respectively. In the former case, all markers,

31 regardless of diversity ($H_S$), have the same $G_{ST}$ value in expectation and thus give replicated

32 estimates of differentiation. In the latter case, markers of higher $H_S$ have lower $G_{ST}$ values,

33 resulting in a negative, roughly linear correlation between $G_{ST}$ and $H_S$ across loci. I propose

34 that the correlation coefficient between $G_{ST}$ and $H_S$ across loci, $r_{GH}$, can be used to distinguish

35 the two cases and to detect mutational effects on $G_{ST}$. A highly negative and significant $r_{GH}$,

36 when coupled with highly variable $G_{ST}$ values among loci, would reveal that marker $G_{ST}$

37 values are affected substantially by mutations and marker diversity, underestimate population

38 differentiation, and are not comparable among studies, species and markers. Simulated and

39 empirical datasets are used to check the power and statistical behaviour, and to demonstrate

40 the usefulness of the correlation analysis.

41

42 **Introduction**

43 A species rarely breeds at random throughout its whole range to form a homogenous unit.

44 Frequently a species is genetically structured in space, subdivided into subunits called demes,

45 races, subpopulations, … Delineating the spatial genetic structure by dividing a species into

46 subunits and quantifying the genetic differentiation among the subunits is important in many

47 biological fields such as evolution, conservation, human medicine and forensics. The

48 subdivision can be made based on natural (e.g. rivers) or artificial (e.g. dams or highways)

49 boundaries, on geographical locations, or on genetic data (e.g. Pritchard *et al.* 2000), and the

50 differentiation can be measured from marker data using Wright's (1943) $F_{ST}$, Nei's (1973)

51 $G_{ST}$ and related statistics such as Weir & Cockerham's (1984) $\theta$ and Slatkin's (1995) $R_{ST}$. The

52 development and wide application of highly polymorphic markers such as microsatellites

53 made these statistics ever more popular, but also caused some confusion and concern. The

54 most popular differentiation statistics, $F_{ST}$ and $G_{ST}$, are believed to underestimate population

55 differentiation when calculated from markers of high diversity (e.g. Nagylaki 1998; Hedrick

56  2005; Jost 2008), and for this reason alternative statistics were proposed and advocated to

57  replace them (Hedrick 2005; Jost 2008; Meirmans & Hedrick 2011). The new differentiation

58  statistics, however, are criticized for their lack of biological meaning and applications, their

59  marker dependency but drift independency, and so on (see Ryman & Leimar 2009; Whitlock

60  2011; Wang 2012).

61      The claim that $F_{ST}$ and $G_{ST}$ underestimate population differentiation is made from

62  both theoretical and empirical grounds. The mathematical definition of $G_{ST} = (H_T - H_S) / H_T$

63  suggests that it cannot take values larger than the average within subpopulation homozygosity,

64  $1 - H_S$ (Jin & Chakraborty 1995; Nagylaki 1998; Hedrick 1999, 2005). This constraint is true

65  both mathematically and biologically. Both $F_{ST}$ and $G_{ST}$ are inherently constrained by $H_S$, as

66  they signify the amount of genetic variation between populations ($V_B$) as a proportion of the

67  total variation $V_T$, which is composed of within ($V_W$) and between ($V_B$) population variation.

68  A high $H_S$ means a high $V_W$, and necessarily a low $V_B$ as a proportion of $V_T$ (i.e. low $F_{ST}$ and

69  $G_{ST}$). However, the constraint imposed on $F_{ST}$ and $G_{ST}$ by $H_S$ does not necessarily mean they

70  are always marker $H_S$ dependent and underestimate differentiation from markers of high $H_S$,

71  as claimed by some authors (e.g. Nagylaki 1998; Hedrick 1999, 2005; Jost 2008). On the

72  empirical grounds, some studies showed that $G_{ST}$ based on highly polymorphic

73  microsatellites is usually lower than $G_{ST}$ based on weakly polymorphic allozyme loci (e.g.

74  Sanetra & Crozier 2003), and is obviously too low for highly differentiated subspecies (e.g.

75  Balloux *et al*. 2000; Carreras-Carbonell *et al*. 2006). These empirical evidences are true for

76  these particular systems, but do not suggest that $F_{ST}$ and $G_{ST}$ calculated from highly

77  polymorphic markers must always underestimate population differentiation in all

78  circumstances.

79      Are $F_{ST}$ and $G_{ST}$ dependent on marker diversity? Do they always underestimate

80  population differentiation from markers of high diversity (e.g. microsatellites)? Under which

81  set of conditions do they provide marker dependent (and thus biased) and marker independent

82  (and thus accurate) estimates of population differentiation? Is it possible to detect whether

83  $F_{ST}$ and $G_{ST}$ values calculated from a set of markers underestimate differentiation or not? In

84  this paper, I will use a combination of analytical modelling, simulated data and empirical data

85  to answer these questions. I show $G_{ST}$ is independent of $H_S$ when mutation rate ($u$) is small

86  relative to migration rate ($m$) or drift ($1/2N$). Otherwise, $G_{ST}$ decreases nearly linearly with an

87  increase in $H_S$. The results suggest a test for the presence or absence of mutational effects on

88  $G_{ST}$. If single-locus $G_{ST}$ values are highly variable and the correlation between single-locus

89    $G_{ST}$ and $H_S$ values is significantly negative, then the observed $G_{ST}$ values are substantially

90    affected by mutations, are locus specific, and seriously underestimate the differentiation due

91    to population demography. If the correlation is insignificant, then the observed single-locus

92    $G_{ST}$ values are unaffected by mutations and are marker independent. They can then be

93    averaged to give an overall estimate of the genetic differentiation caused by demography only.

94    Simulations and empirical data are analysed to check the power and statistical properties of

95    the correlation and regression analyses.

96    **Method**

97    The relationship between $G_{ST}$ and $H_S$ is investigated by analyses of standard population

98    genetics models of migration, drift and mutation. The results are then verified by analyses of

99    simulated and empirical datasets.

100    *Theory*

101    Following most previous studies of $F_{ST}$, I assume a population under the finite island model

102    of migration (Wright 1931) and the infinite allele model of mutation (Kimura & Crow 1964)

103    for mathematical tractability. The results and conclusions are, however, applicable

104    qualitatively to populations under other migration models, such as Wright's (1943) isolation

105    by distance or neighbourhood model and Kimura & Weiss's (1964) stepping stone model,

106    and under other mutation models, such as stepwise mutation model for microsatellites or

107    allozymes (Ohta & Kimura 1973).

108    Under the finite island model with migration rate $m$ among $s$ subpopulations of

109    effective size $N$, and under the infinite allele model for a neutral locus with mutation rate $u$,

110    the recurrence equations for the expected homozygosity within a subpopulation, $J_0$, and

111    between two subpopulations, $J_1$, is (Nei 1975; Li 1976)

112    $$J_{0(t+1)} = d\big(a\big(c + (1-c)J_{0(t)}\big) + (1-a)J_{1(t)}\big), \tag{1}$$

113    $$J_{1(t+1)} = d\big(b\big(c + (1-c)J_{0(t)}\big) + (1-b)J_{1(t)}\big), \tag{2}$$

114    where $b = m(2-m)/s$, $a = (1-m)^2 + b$, $c = 1/(2N)$ and $d = (1-u)^2$. Equivalently,

115    $J_0$ and $J_1$ are the probabilities that two genes taken at random from within a subpopulation

116    and from different subpopulations, respectively, are identical in state. The complements, $H_S$

117    $=1- J_0$ and $H_1 =1- J_1$, give the expected (i.e. assuming random union of gametes)

118    heterozygosity or gene diversity (Nei 1973) within and between subpopulations. The total

119    expected heterozygosity or gene diversity in the entire population is $H_T = (H_S + (s -$

120    $1)H_1)/s = 1 - J_1 - (J_0 - J_1)/s$ (Nei 1975). Given $H_T$ and $H_S$, $G_{ST}$ is calculated by $G_{ST} =$

121    $1 - H_S/H_T$ (Nei 1973). Using (1) and (2), we can calculate recurrently the values of $H_S$, $H_T$

122    and $G_{ST}$ at each generation, given parameters $m, N, u, s$ and initial gene identities $J_{0(0)}$ and

123    $J_{1(0)}$.

124         Under the joint action of mutation, migration and drift at rates $u, m$, and $1/(2N)$

125    respectively, the gene diversity ($H_S$, $H_T$) and its distribution ($G_{ST}$) will reach equilibrium

126    values. $G_{ST}$ attains its equilibrium value much faster than $H_S$ and $H_T$, because it is determined

127    by the strongest (in terms of rate) while $H_S$ and $H_T$ are determined by the weakest among the

128    forces of mutation, migration and drift. The equilibrium gene identity values are (Nei 1975;

129    Li 1976)

130    $J_{0(\infty)} = cd(a - (a - b)d)/G,$ (3)

131    $J_{1(\infty)} = cdb/G,$ (4)

132    where $G = 1 - d(a(1 - c) + 1 - b) + d^2(a - b)(1 - c)$. The equilibrium gene diversity

133    and differentiation values, $H_{S(\infty)}$, $H_{T(\infty)}$ and $G_{ST(\infty)}$, can be calculated using (3) and (4). The

134    expression for $G_{ST(\infty)}$ is complicated, but can be simplified approximately to (Takahata & Nei

135    1984)

136    $G_{ST(\infty)} \approx 1/\left[1 + 2N\left(\frac{s}{s-1}\right)\left(\frac{1}{(1-m)^2(1-u)^2} - 1\right)\right].$ (5)

137    When $m, u \ll 1$, (5) is further simplified to (Takahata & Nei 1984)

138    $G_{ST(\infty)} \approx 1/\left[1 + 4N\left(\frac{s}{s-1}\right)(m + u)\right].$ (6)

139    When $s \rightarrow \infty$, (6) again reduces to the equilibrium $F_{ST}$ of the infinite island model of Wright

140    (1969, page 291), indicating that $F_{ST}$ and $G_{ST}$ are equivalent (Nei 1977; Takahata & Nei

141    1984).

142         Although several studies have used similar models to investigate the impact of

143    mutations on $F_{ST}$ and $G_{ST}$ (e.g. Ryman & Leimar 2008; Whitlock 2011), none has examined

144    the direct relationship between $G_{ST}$ and $H_S$. Herein I will use equations (1-6) to explore this

145    relationship in populations in both equilibrium and non-equilibrium conditions under

146    different parameter ($m$, $u$, $N$, $s$) combinations. This is important as both $G_{ST}$ and $H_S$ are

147    estimable from marker data, and examining the observed patterns of $G_{ST}$ and $H_S$ at a set of

148    marker loci sheds light on the possible impact of mutations on $G_{ST}$.

149    *Simulations*

150    Simulated data typical of those encountered in practice were generated to test whether the

151    correlation analysis of single locus estimates of $G_{ST}$ and $H_S$ could be used to detect the effect

152    of mutations on $G_{ST}$ when it is present, and whether the analysis does not falsely detect the

153    effect of mutations when it is absent. The behaviour and power of the correlation analysis

154    were investigated by analysing simulated data with varying sampling intensities (of

155    individuals from a subpopulation, of subpopulations, and of markers), different population

156    properties ($N$, $s$, $m$, $u$) and different mutation and migration models.

157        The simulations considered the finite island model as described above, and a one-

158    dimensional circular stepping stone model (Kimura & Weiss 1964). In the latter model, a

159    number of $s$ subpopulations are arranged in a circle and each subpopulation receives a

160    proportion $m/2$ of its individuals from each of its two neighbouring subpopulations. In both

161    models, each subpopulation is composed of $N$ diploid monoecious individuals. At each

162    discrete generation, the events are mutations, migrations and reproductions occurring in that

163    order. Mutations are assumed to follow either the infinite allele model or the stepwise

164    mutation model. For the former, a mutation always generates a novel allele the population has

165    never seen before. For the latter, the mutated allele increases or decreases in size by 1 repeat

166    with an equal probability of 0.5. For both models, the number of new mutations at a locus in

167    each subpopulation at each generation was sampled from a Poisson distribution with

168    parameter value $2Nu$. For each new mutation, a gene was drawn at random from the $2N$ genes

169    and was changed according to the mutation model. Reproduction is assumed to be random

170    union of gametes, such that selfing and outbreeding occur at rates $1/N$ and $1-1/N$ respectively,

171    and the effective size is equal to the census size for each subpopulation.

172        An ancestral population was assumed to be the same as the subdivided population

173    described above except for population size and structure. It was unsubdivided and had a size

174    $N_A = rsN$, where $r$=0.5, 1 and 2 such that it had equilibrium genetic diversity smaller than,

175    close to, and larger than the subdivided population respectively. The ancestral population was

176    maintained for a large number of generations for it to reach mutation-drift equilibrium at a

177    neutral locus with mutation rate $u$ (which was variable among a number of $L$ loci). It was then

178  subdivided into *s* subpopulations of size *N*, which were maintained as described above for *g*

179  (=100, 200, 400) generations or for a sufficiently large number of generations, in the order of

180  Max($1/u$, $1/m$, $2N$), to reach mutation-drift-migration equilibrium. A sample of *M* individuals

181  was then taken at random from each of *R* ($\leq s$) randomly selected subpopulations, and each

182  sampled individual was genotyped at a number of *L* loci.

183       The genotype data were then used to calculate Nei & Chesser's (1983) nearly

184  unbiased estimators of $H_S$, $H_T$, and thus $G_{ST}$,

185  $\widehat{H}_S = \frac{2\widetilde{M}}{(2\widetilde{M}-1)R} \sum_{j=1}^{R} \left(1 - \sum_{i=1}^{k} x_{ij}^2\right),$

186  $\widehat{H}_T = 1 - \sum_{i=1}^{k} \left(\frac{1}{R}\sum_{j=1}^{R} x_{ij}\right)^2 + \frac{\widehat{H}_S}{2\widetilde{M}R},$

187  $\widehat{G}_{ST} = 1 - \widehat{H}_S / \widehat{H}_T,$

188  where $x_{ij}$ is the frequency of allele *i* in the sample from subpopulation *j*, *k* is the number of

189  alleles observed in the set of samples from the *R* subpopulations, and $\widetilde{M}$ is the harmonic mean

190  sample sizes ($\equiv M$ in the simulations).

191       The estimates $\widehat{H}_S$ and $\widehat{G}_{ST}$ were then used to calculate their correlation coefficient $r_{GH}$

192  across loci. The significance of $r_{GH}$ was tested by a permutation analysis in which $\widehat{H}_S$ and $\widehat{G}_{ST}$

193  were both randomized across loci before calculating $r_{GH}$ in $10^6$ replicates. The proportion of

194  replicates in which $r_{GH}$ was smaller than the $r_{GH}$ value calculated from the original data was

195  taken as the *p* value. The correlation coefficient was taken as statistically significant when

196  $p<0.001$. A significant negative correlation $r_{GH}$ indicates that $\widehat{G}_{ST}$ has been affected by

197  mutations and thus underestimates the differentiation caused purely by demography (drift and

198  migration). Otherwise, markers with different levels of diversity $\widehat{H}_S$ are equally differentiated,

199  they all give the same $G_{ST}$ expected from the impact of drift and migration only, and the

200  single locus $G_{ST}$ estimates can be averaged to give a better (in precision) overall estimate of

201  differentiation.

202       Too many parameter combinations, due to the numerous parameters and the numerous

203  plausible values of each parameter, are involved in determining $\widehat{H}_S$ and $\widehat{G}_{ST}$ that a realistic

204  simulation study can only consider a small fraction of them. I studied the effect of each

205  parameter in isolation of others each time by varying the values of the focal parameter only

206  (see Table 1). For each parameter combination, a number of 100 replicate datasets were

207    generated and analysed. The analysis results were reported as the mean correlation coefficient

208    between $\hat{G}_{ST}$ and $\hat{H}_S$, $\bar{r}_{GH}$, and the proportion of replicates with a statistical significant (at

209    $p<0.001$) $r_{GH}$ among the 100 replicates.

210          The simulation program was checked by comparing the simulated against the

211    predicted values of several quantities to make sure it worked properly. First, the effective size

212    of the entire population in the finite island model is $N_e = sN/(1-F_{ST})$ (Wright 1943; Wang &

213    Caballero 1999), where $F_{ST}$ can be replaced by $G_{ST}$. This theoretical prediction was compared

214    with that estimated from the simulated pedigrees, using the formula $\frac{1}{2N_e} = \frac{\theta_{t+1}-\theta_t}{1-\theta_t}$ where $t$ the

215    generation is large and $\theta_t$ is the average coancestry at generation $t$ for all individuals in the

216    entire population. Second, the predicted values of $H_S$, $H_T$ and $G_{ST}$ by (3-4) were compared

217    with the corresponding observed values for an equilibrium population under infinite allele

218    and finite island models. In all situations investigated, the predicted and estimated (observed)

219    values fitted very well.

220    *Empirical data*

221    The simulation model may be too simple to reflect the reality. In a real population, both *m*

222    and *N* may vary over space and time, and migrations and mutations may not follow the ideal

223    models assumed in the simulations. Supplementing simulations, therefore, I also analysed

224    several recently published empirical datasets to demonstrate the use of the proposed

225    correlation analysis.

226    *Atlantic Salmon*: To investigate the genetic structure of Atlantic salmon populations in the

227    entire North American range of the species, Moore *et al*. (2014) sampled 9142 individuals

228    from 153 populations and genotyped each individual at 15 microsatellite loci. They also

229    sampled 1080 individuals from 50 populations and genotyped each individual at 3192 SNP

230    loci. The two datasets were analysed separately in the present study of the relationship

231    between $G_{ST}$ and $H_S$.

232    *Blacknose sharks*: Using 23 microsatellites and mtDNA sequences, Portnoy *et al*. (2014)

233    investigated the genetic structure and barriers to gene flow of 10 blacknose shark populations

234    sampled (651 individuals in total) from the western North Atlantic Ocean. It was found that

235    the $F_{ST}$ values at the 23 microsatellite loci between the Bahamas and any of the other

236    populations were more than an order of magnitude greater than the values between any two

237    of the other populations. Therefore, $G_{ST}$ and $H_S$ values were calculated for each locus in the 2

238    alternative population structures, the 10- and 2-population (Bahamas and the rest) models in

239    the present study.

240    *Mediterranean shore crab*: Schiavina *et al*. (2014) investigated the genetic structure of the

241    Mediterranean shore crab (*Carcinus aestuarii*) in the Adriatic Sea (central Mediterranean),

242    using 11 polymorphic microsatellites in 431 individuals collected from eight sites. One locus,

243    Cae30, has only 5 alleles and a gene diversity of $H_S = 0.1$, much lower than the locus with the

244    2nd lowest diversity, which has 13 alleles and a $H_S = 0.77$. So Cae30 was excluded as an

245    obvious outlier from the $G_{ST}$ and $H_S$ correlation analysis.

246    *Blacktip reef sharks*: To understand the genetic structure of blacktip reef sharks

247    (*Carcharhinus melanopterus*), Vignaud *et al*. (2014) sampled 758 individuals from 15 sites (4

248    widely separated locations in the Indo-Pacific and 11 islands in French Polynesia) widely

249    distributed in the Indian and Pacific Oceans. Each sampled individual was genotyped at 17

250    microsatellite loci. Three loci (cil169, cli107 and cli12) were found to deviate significantly

251    from Hardy-Weinberg equilibrium and were suspected to contain null alleles (Vignaud *et al*.

252    2014). The three loci were excluded from their original genetic analysis. Herein I investigated

253    the impact of mutations on the estimated differentiation among these shark populations by

254    analysing the relationship between $G_{ST}$ and $H_S$, using both the entire set of 17 loci and the

255    selected subset of 14 loci.

256    *Copper rockfish*: Using 17 microsatellite DNA loci, Dick *et al*. (2014) assessed the genetic

257    diversity of and the differentiation among ten populations of copper rockfish (*Sebastes*

258    *caurinus*) representing paired samples of outer coast and the heads of inlets in five replicate

259    sounds on the west coast of Vancouver Island, British Columbia. The sample size per

260    population varies between 30 and 105. I calculated the $G_{ST}$ and $H_S$ values at each of the 17

261    loci among the 10 populations, and tested whether the marker differentiation is affected by

262    mutations or not.

263    **Results**

264    *Analytical results*

265    Equation (6) suggests that $G_{ST}$ at neutral loci is determined by the joint action of migration,

266    mutation and drift occurring at rates *m, u*, and $1/(2N)$ respectively. The relative impact of

267    each evolutionary force on $G_{ST}$ is determined by its rate as a proportion of the total rate,

268    $m+u+1/(2N)$. When subpopulations are small such that drift is the dominating force (i.e.

269     $1/(2N) \gg u+m$), then $H_S \to 0$ (i.e. fixation) and $G_{ST(\infty)} \to 1$ in equilibrium conditions. When

270     mutation is weak relative to drift and migration (i.e. $u \ll 1/(2N) + m$), then $G_{ST(\infty)} \approx$

271     $1/\left[1 + 4N\left(\frac{s}{s-1}\right)m\right]$, which suggests that $G_{ST(\infty)}$ reflects demography only and all loci with

272     varying but small $u$ have the same expected $G_{ST}$. In contrast, for loci with a high $u$ in a

273     population with a large $N$ and a small $m$ (i.e. $u \gg 1/(2N) + m$), $G_{ST(\infty)}$ becomes locus (or

274     mutation) dependent and covaries with locus specific $H_S$ (below). In such a case, marker

275     based $G_{ST(\infty)}$ has little bearing on population demography, the $G_{ST(\infty)}$ value calculated from

276     one set of loci can hardly be congruent with that from another set of loci, and it is

277     incomparable among studies, species and loci.

278          Figure 1 plots the equilibrium $G_{ST}$ as a function of $H_S$, calculated by (5) and (3)

279     respectively, for different parameter combinations of $u$, $m$ and $N$, assuming $s=10$. When

280     differentiation is expected to be small due to either strong migration ($m \geq 0.01$) or weak drift

281     ($N \geq 2500$), $G_{ST}$ keeps constant and does not vary with $H_S$ in its entire range of [0, 1] caused by

282     widely varying $u$ values in range of $[10^{-6}, 10^{-2}]$. The observation disproves the belief that $G_{ST}$

283     underestimates differentiation and becomes $H_S$ dependent when $H_S$ is high (e.g. Nagylaki

284     1998; Hedrick 1999, 2005; Jost 2008). High $H_S$ values (say 0.95) do constrain $G_{ST}$ to small

285     values with a maximum of 1- $H_S$, but do not necessarily lead to underestimated and locus-

286     varying $G_{ST}$. What is relevant is the main mechanism (determined by the relative strengths of

287     mutation, drift and migration) leading to the observed high $H_S$, not the observed high $H_S$ *per*

288     *se*. A high $H_S$ is usually due to a high $u$ or/and a high $N$. However, as long as $m$ is much

289     higher than $u$, $G_{ST}$ is virtually independent of $H_S$.

290          When drift is strong (i.e. $N$ small) and migration is weak relative to mutations, $G_{ST}$

291     decreases almost linearly with an increasing $H_S$ due to an increasing $u$ (Figure 1). Only in this

292     situation is the belief that $G_{ST}$ covaries with $H_S$ (e.g. Nagylaki 1998; Hedrick 1999, 2005; Jost

293     2008) certified. For the parameter combination $N=250$, $m=0.001$, and $s=10$ in Figure 1, for

294     example, $G_{ST}$ keeps almost a constant value of 0.45 when $u$ varies between $10^{-6}$ and $3 \times 10^{-6}$

295     that leads to a $H_S$ varying between 0 and 0.5. With $u > 3 \times 10^{-6}$ and thus $H_S > 0.5$, $G_{ST}$ begins to

296     decrease linearly with an increasing $H_S$ (or $u$). Similar results are obtained with other values

297     of the number of subpopulations ($s$).

298          Many generations, in the order of $1/m$, $1/u$ or $2N$ whichever is the smallest, are

299     required for a subdivided population to reach the equilibrium differentiation. Natural

300    populations may never reach such equilibrium as $m$ and $N$ are constantly changing. It is thus

301    important to check whether the above observations (Figure 1) also apply to non-equilibrium

302    populations. Figure 2 plots $G_{ST}$ as a function of $H_S$ at generations 50, 200 and 1000 since the

303    subdivision. Mutation rate ($u$) is assumed to vary from $10^{-6}$ to $10^{-2}$, and the initial gene

304    diversity is assumed to be $J_{0(0)} = J_{1(0)}$ and to take values $rJ_{0(\infty)}$, where $r=1$, 0.5 and 0.25. The

305    relationship between $G_{ST}$ and $H_S$ in a non-equilibrium population is similar to that in an

306    equilibrium population (Figure 1). Whenever $u \ll 1/(2N) + m$, $G_{ST}$ does not vary with $H_S$ (or

307    $u$). Depending on $u$ as well as $N$ and $m$, $H_S$ can freely vary in almost the entire range of [0,1]

308    without affecting the value of $G_{ST}$. Otherwise, $G_{ST}$ decreases nearly linearly with an

309    increasing $H_S$ (or $u$). The further away a population departs from the equilibrium, the less

310    affected it is by mutations because the latter require time to accumulate. When $N=250$, for

311    example, mutations start to have a substantial impact on $G_{ST}$ at generations 50, 200 and 1000

312    when $H_S \geq 0.8$ ($u \geq 0.0015$), $H_S \geq 0.5$ ($u \geq 0.00001$) and $H_S \geq 0.3$ ($u \geq 0.00003$) respectively.

313    Initial gene identities (or diversities) seem to have little effect on the relationship between $G_{ST}$

314    and $H_S$ at any generation.

315    *Simulation results*

316    Confirming the analytical results presented above, simulations show that, when mutations are

317    strong relative to migrations ($m=0.001$), $G_{ST}$ estimates vary among loci that have different $u$

318    and thus different $H_S$, and are negatively correlated with $H_S$ (Figure 3). This is true for the

319    finite island and stepping stone migration models, and for the infinite allele, finite allele and

320    stepwise mutation models. This is also true no matter the population is at mutation-drift-

321    migration equilibrium (Figure 3) or not (data not shown). The negative correlation in stepping

322    stone migration model and infinite allele mutation model is stronger than that in other

323    migration and mutation models. In contrast, when mutations are weak relative to migrations

324    ($m=0.01$), $G_{ST}$ estimates are small and are almost constant among loci with different $u$ and

325    thus different $H_S$. This is shown for an equilibrium population under different migration and

326    mutation models (Figure 3), but is also true for non-equilibrium populations (data not shown).

327    When migrations are weak relative to mutations such that $G_{ST}$ is substantially affected

328    by $u$ and becomes negatively correlated with $H_S$, a modest sampling effort is needed to detect

329    the correlation for different migration and mutation models (Figure 4). This is also true for

330    populations that have not reached mutation-drift-migration equilibrium (data not shown).

331    Setting the statistical significance at a conservative level of $p<0.001$, the false detection rate

332    of mutational effects is low (generally below 7%), while the power is generally above 60%

333    except when less than 10 loci and less than 4 subpopulations are used in the analysis. In

334    agreement with the results in Figure 3, the correlation analysis is less powerful for the island

335    migration model coupled with the stepwise or 2-allele mutation model than other models.

336    While the power increases with the numbers of sampled loci and sampled subpopulations

337    (Figure 4), it is little affected by the number of sampled individuals per subpopulation, $M$, as

338    long as $M > 10$. This is not surprising because the population is highly differentiated for the

339    parameter combinations and just a few individuals per subpopulation would allow for a good

340    estimate of $G_{ST}$.

341    *Empirical analysis*

342    The Atlantic salmon data clearly show an extremely strong negative correlation ($r = -0.953$)

343    between $G_{ST}$ and $H_S$ estimates among the 15 microsatellites (Figure 5A), with a $p$ value of

344    $0.0 \times 10^{-6}$. These markers are highly polymorphic, with $H_S$ varying between 0.66 and 0.94 and

345    with the number of observed alleles varying between 15 and 91. Compatible with a

346    substantial impact of mutations, these markers have low but highly variable $G_{ST}$ values,

347    varying between 0.02 and 0.09 with a mean of 0.045 and a coefficient of variation of 0.629.

348    These single locus $G_{ST}$ values are all highly significant, as determined by permutation

349    (permuting individuals among subpopulations) tests.

350          In contrast, the correlation between $G_{ST}$ and $H_S$ estimates of the 3192 SNPs (Figure

351    5B) is positive and small ($r=0.044$), with a $p$ value of 0.993 which is insignificant. $H_S$ values

352    distribute nearly uniformly in the range [0, 0.5]. While most SNPs have $G_{ST}$ values of about

353    0.1, quite a few outliers show $G_{ST}$ values well above 0.4. The mean $G_{ST}$ is 0.099 for the 3192

354    SNPs and is 0.091 when the outlier SNPs with $G_{ST} > 0.3$ are removed. Both values are much

355    larger than the mean $G_{ST}$ across the 15 microsatellites which is 0.045. The comparison

356    between SNPs and microsatellites further verifies that the differentiation at microsatellites is

357    greatly impacted by mutations and thus underestimates the underlying population

358    differentiation due to demography.

359          The blacknose sharks have highly variable single-locus $G_{ST}$ values, with the highest

360    being 0.35 and 0.18 and the lowest being 0 and 0 for the 2- and 10-population models

361    respectively (Figure 5C). Among the 23 microsatellites, $G_{ST}$ and $H_S$ estimates are moderately

362    negatively correlated, with a correlation coefficient of -0.41 ($p=0.017$) and -0.43 ($p=0.007$)

363    for the 2- and 10-population models respectively. None of the correlations are significant at

364   $p$=0.001, but there is a clear trend of less differentiation at more polymorphic marker loci

365   which indicates that mutations might have reduced the $G_{ST}$ values at these loci.

366   The differentiation calculated from each of the 10 microsatellites is low ($G_{ST}$ <0.04)

367   among the 8 Mediterranean shore crab populations (Figure 5D). Nevertheless, $G_{ST}$ and $H_S$

368   estimates are highly negatively correlated, with a correlation coefficient of -0.80 and a small

369   $p$ value (0.010). It is likely that mutations have substantially impacted on the $G_{ST}$ estimates

370   from these microsatellites, and thus the underlying population differentiation due to

371   demography may well be underestimated by these microsatellites.

372   The 17 microsatellites in blacktip reef sharks are highly variable in diversity, with the

373   number of observed alleles varying from 4 to 48 and the $H_S$ varying from 0.15 to 0.89. The

374   $G_{ST}$ values among the 15 populations estimated from the 17 loci are also highly variable,

375   from 0.04 to 0.41 (Figure 5E). The 3 loci showing deviation from Hardy-Weinberg

376   equilibrium are apparently not outliers in terms of both diversity and differentiation. The

377   single locus $G_{ST}$ and $H_S$ estimates are highly negatively correlated, with a correlation

378   coefficient of -0.890 ($p$=0.000×10$^{-6}$) and -0.913 ($p$=0.000×10$^{-6}$) for the entire set of 17 loci

379   and the subset of 14 loci respectively. In this system, mutations are highly likely to have

380   reduced the differentiation of the microsatellites; the underlying population differentiation

381   due to drift and migration should be higher than the average $G_{ST}$ value calculated from these

382   microsatellites.

383   The differentiation measured by $G_{ST}$ at each of the 17 microsatellites is low among the

384   10 copper rockfish populations (Figure 5F). Except for locus Sra11-103 which has a $G_{ST}$ =

385   0.09, single locus $G_{ST}$ values are below 0.05. The overall mean $G_{ST}$ across loci is 0.027, very

386   close to the $F_{ST}$ value 0.031 obtained by Dick $et\ al.$ (2014). Single locus $G_{ST}$ and $H_S$ estimates

387   are not correlated, with a correlation coefficient of 0.011 and a $p$ value of 0.649. It can be

388   concluded confidently that mutations have no impact on these $G_{ST}$ estimates, and all markers,

389   regardless of polymorphisms, should have the same expected differentiation which is

390   equivalent to the population differentiation. The average $G_{ST}$ across loci, 0.027, should be an

391   unbiased estimate of the population differentiation due to demography.

392   **Discussion**

393   The claim that $F_{ST}$ and $G_{ST}$ are dependent on marker $H_S$ and underestimate population

394   differentiation when calculated from highly polymorphic (i.e. high $H_S$) markers (e.g.

395      Nagylaki 1998; Hedrick 2005; Jost 2008) can be misleading. It has led to the conclusion that

396      these traditional statistics should be either "corrected" for $H_S$ (e.g. Hedrick 2005) or replaced

397      by new statistics such as $D$ (Jost 2008). The claim creates lots of confusions, as if $F_{ST}$ and $G_{ST}$

398      should be independent of $H_S$ to be correct measures of differentiation. As Wright (1978, p.82)

399      explicitly stated, however, $F_{ST}$ (the same for $G_{ST}$) measures "the amount of differentiation

400      among subpopulations, relative to the limiting amount under complete fixation". Complete

401      fixation means each subpopulation is fixed with an allele (i.e. all individuals in a

402      subpopulation have the same homozygous genotype), which is not necessary to be unique

403      among subpopulations. Fixation results in $H_S$=0, and the maximal differentiation of $F_{ST}$ =1 is

404      achieved only at $H_S$ =0. For this reason, Wright (1951) also called his $F_{ST}$ a fixation index,

405      among other fixation indices of $F_{IS}$ and $F_{IT}$. The quantity $H_S$ measures the *absolute* distance

406      from complete fixation (i.e. $H_S$ =0), and naturally constrains $F_{ST}$, which measures the *relative*

407      (to total diversity $H_T$) or standardized distance from complete fixation. The definition of

408      $G_{ST} = 1 - H_s/H_T$ (Nei 1973) makes the functional relationship between absolute (i.e. $H_S$)

409      and relative (i.e. $G_{ST}$) differentiations explicit. Therefore, both $F_{ST}$ and $G_{ST}$ legitimately

410      depend on, or more precisely, are constrained by $H_S$. This relationship is true both

411      mathematically and biologically, and does not inherently cause $F_{ST}$ and $G_{ST}$ to underestimate

412      differentiation for markers with high $H_S$.

413      More precisely, $F_{ST}$ and $G_{ST}$ become marker dependent and underestimate population

414      differentiation only when migration rate is lower than mutation rate. Otherwise, they provide

415      accurate estimates of population differentiation regardless of marker $H_S$. In a population with

416      low migration rates (i.e. $m < u$), a marker with a higher $u$ is expected to have a higher $H_S$ (or

417      absolute differentiation) and a correspondingly lower $G_{ST}$ (or relative differentiation) in both

418      equilibrium and many non-equilibrium conditions (Whitlock 2011; this study). It should be

419      emphasized that a high $u$ does not necessarily lead to a high $H_S$, and *vice versa*. This is

420      because it is the quantity $uN$ rather than $u$ that determines $H_S$. A marker with a small $u$ in a

421      population with a large $N$ can still harbour a high $H_S$, and a marker with a large $u$ in a

422      population with a small $N$ can still have a low $H_S$. The statement that microsatellites, because

423      of their high allelic polymorphisms and high $H_S$, must always underestimate differentiation is

424      imprecise. Such markers show less differentiation than less polymorphic markers (e.g. SNPs)

425      only when migration is weak ($m<u$), as illustrated by Figures 1 and 2.

426      This study reveals that whenever $m<u$ and thus mutations have a substantial impact,

427      single locus $G_{ST}$ values decrease almost linearly with single locus $H_S$. This is true in both

428   equilibrium (Figure 1) and non-equilibrium populations, as verified by simulations under

429   different migration and mutation models (Figure 3). It is not surprising that the pattern

430   observed under the ideal island migration model and the infinite allele mutation model

431   applies to other migration and mutation models, because $G_{ST}$ and $F_{ST}$ are defined as

432   descriptive statistics without any predefined demographic and mutation models. Mutations

433   act to increase genetic diversity ($H_S$ and $H_T$) and thus to decrease differentiation among

434   subpopulations, no matter they occur in the finite or infinite allele model or in the stepwise

435   mutation model (Wright 1943). Migrations, in contrast, tend to redistribute genetic diversity

436   evenly among subpopulations. Thereby they tend to reduce the difference between $H_S$ and $H_T$

437   and thus to reduce $G_{ST}$, no matter they occur in the island model, stepping stone model or the

438   isolation-by-distance model.

439         The simulations confirm that a correlation analysis between single locus $G_{ST}$ and $H_S$

440   estimates can be used to detect the mutational effects on differentiation. Under typical

441   sampling intensities, the analysis has sufficient power to identify the mutational effect when

442   it is present, and it does not falsely detect the mutational effect when it is absent (Figure 4),

443   when the significance level is chosen as $p=0.001$. A higher significance $p$ value (say, 0.05 or

444   0.01) leads to higher powers, but also higher false detect rates. Under the finite island and

445   infinite allele models (first row in Figure 4), for example, the power (when $m=0.001$) and

446   false detection rate (when $m=0.01$) increase to 86.7% and 11.8% respectively when $p=0.01$,

447   and to 90.7% and 30.0% respectively when $p=0.05$. A good balance between type I and II

448   errors is achieved at $p=0.001$, which leads to a false detection rate being always below 7%

449   irrespective of the widely varying sampling intensities of the number of subpopulations, the

450   number of individuals per subpopulation, and the number of loci and polymorphisms (Figure

451   4).

452         Two out of the five empirical microsatellite datasets (Figures 5A, 5E) show strong

453   evidence (a high negative $r_{GH}$ value and a small $p$ value) that mutations have reduced $G_{ST}$

454   estimated from microsatellites, two datasets (Figures 5C, 5D) indicate a similar trend with

455   higher uncertainties, and the remaining dataset (Figure 5F) shows no detectable effect of

456   mutations on $G_{ST}$. It is noticeable that the copper rockfish populations (Figure 5F) have high

457   and widely variable $H_S$ values across the 17 microsatellites, the highest $H_S$ being 0.936. These

458   $H_S$ values are similar to those of the microsatellites in Atlantic salmon populations (Figure 5A)

459   and the blacktip reef shark populations (Figure 5E). Yet, contrasting patterns of $G_{ST}$ and $H_S$

460   were observed among the three species. This again verifies the theory and simulation based

461 conclusion that a high $H_S$ does not necessarily lead to marker dependent $G_{ST}$, and does not

462 necessarily result in underestimation of population differentiation. In situations where the

463 correlation between $G_{ST}$ and $H_S$ has a high uncertainty (e.g. Figure 5C), collection of more

464 data (by sampling more subpopulations, loci, and individuals) may confirm or reject the

465 hypothesis that $G_{ST}$ in a study system is affected by $H_S$ or mutations. In contrast, the analysis

466 of a big SNP dataset (Figure 5B) does not detect any mutational effect. The correlation

467 between single locus $G_{ST}$ and $H_S$ values, 0.044, is small and positive, and clearly indicates no

468 mutational effects on $G_{ST}$. The results are understandable because the $u$ for SNPs can be

469 several orders smaller than that for microsatellites, and as a result is more likely to be smaller

470 than migration rate $m$.

471       The five empirical microsatellite datasets were taken from the most recent literature at

472 random with regard to the relationship between $G_{ST}$ and $H_S$, which was revealed only after the

473 correlation analyses. If this small sample of datasets represents the reality, then we may

474 conclude that underestimation of differentiation by microsatellites could be a common

475 problem (Hedrick 1999, 2005; Jost 2008). A meta-analysis of many more microsatellite

476 datasets as exemplified in this study is required for a solid conclusion. However, while

477 microsatellites do underestimate differentiation in some (or many) situations, they can also

478 yield unbiased estimates in situations where migration is high as shown for the copper

479 rockfish populations (Figure 5F). The assertion that all microsatellites of high $H_S$

480 underestimate differentiation and therefore all $G_{ST}$ estimates should be standardized (Hedrick

481 2005) or abandoned and replaced by new differentiation statistics (Jost 2008) is unjustified.

482 In addition to the problems shown before (Ryman & Leimar 2009; Whitlock 2011; Wang

483 2012), these new statistics are also marker diversity dependent as shown below.

484       It is notable that several authors have conducted a meta-analysis of the relationship

485 between $G_{ST}$ and $H_S$ across species/populations (Heller & Siegismund 2009; Meirmans &

486 Hedrick 2011). They found that the estimated $G_{ST}$ is always smaller than the maximum value

487 of 1- $H_S$, as expected, and shows a weak negative correlation with $H_S$. It should be pointed

488 out that the correlation analysis proposed in my study is fundamentally different from that in

489 these meta-analyses. In the latter, the correlation is at the species level, where $G_{ST}$ and $H_S$ are

490 average values across loci for each species. Because different species may have experienced

491 different evolutionary forces and demography such that their $G_{ST}$ values differ, it is unclear

492 what the hypothesis these meta-analyses are trying to prove or disapprove, except for the

493 functional relationship $G_{ST} < 1$- $H_S$ which should however always be true from the definition

494    of $G_{ST}$. The presence of a negative correlation between $G_{ST}$ and $H_S$ does not prove that $G_{ST}$ is

495    underestimated and is marker dependent because of mutational effects. The absence of the

496    correlation does not prove that mutations have negligible effects and $G_{ST}$ is unbiased and

497    marker independent. In my study, the correlation is between single locus values of $G_{ST}$ and

498    $H_S$ within a species (population). The hypothesis, clearly defined and supported by theory and

499    simulations, is that $G_{ST}$ values should be similar across markers of different $H_S$ if mutations

500    are unimportant (when $u<m$), resulting in an $r_{GH}$ not different from 0. Otherwise (i.e. $u>m$),

501    $G_{ST}$ values should decrease with markers showing an increasing $H_S$, resulting in a highly

502    negative correlation between $G_{ST}$ and $H_S$.

503        $G_{ST}$ calculated from a locus measures the genetic differentiation among

504    subpopulations at the locus due to the combined effect of all evolutionary forces (Nei 1973).

505    Selection directly influences $F_{ST}$ and $G_{ST}$, as Wright (1943) illustrated with several different

506    types of selection. In principle, a negative correlation between $H_S$ and $G_{ST}$ can also be

507    generated for markers closely linked with a locus under strong selection for spatially different

508    alleles (which causes a decrease in $H_S$ and an increase in $G_{ST}$) or/and for spatially different

509    allele combinations (which causes an increase in $H_S$ and a decrease in $G_{ST}$). Although my

510    correlation analysis assumes the absence of selection, it should be robust in most applications.

511    First, frequently only a few microsatellites (<30) are used in calculating $F_{ST}$ or $G_{ST}$, and the

512    chance of any of them being under selection or being linked to loci under selection strong

513    enough (compared with other evolutionary forces) for detection is slim. Second, with

514    genomic dense markers such as SNPs, it is highly likely that a small fraction of the loci are

515    under strong selection. The correlation analysis should however still be robust because the

516    vast majority of loci are neutral and a few selected loci should not affect the overall

517    relationship between $H_S$ and $G_{ST}$.

518        This study focusses on the widely applied differentiation statistic $G_{ST}$ (Nei 1973).

519    Other differentiation statistics or estimators such as $\theta$ (Weir & Cockerham 1984), $D$ (Jost

520    2008) and $G'_{ST}$ (Hedrick 2005) could also be affected by mutations and yield marker ($H_S$)

521    dependent estimates. All these statistics measure differentiation at marker loci due to the

522    collective actions of all evolutionary forces, including mutations. When mutations are

523    important (i.e. $u>m$), therefore, differentiation estimates are expected to be different among

524    loci. Some statistics, like $D$ which is claimed to outperform $G_{ST}$ for highly polymorphic

525    markers (Jost 2008), are even more problematic and produce marker dependent

526    differentiation estimates even when mutation rate is small relative to migration rate. For the

527  data simulated in finite island and infinite allele models, finite island and stepwise mutation

528  models, and stepping stone and stepwise mutation models shown in Figure 3, for example,

529  the correlation coefficient between $D$ and $H_S$, $r_{DH}$, is 0.43, 0.30, and 0.22 respectively when

530  $m=0.001$, and is 0.71, 0.24 and 0.26 respectively when $m=0.01$. The correlation is always

531  positive and substantially high, even in the situation where mutation is very weak relative to

532  migration and $G_{ST}$ is uncorrelated with $H_S$. Similarly highly positive $r_{DH}$ values are obtained

533  for all of the empirical datasets. For the Atlantic salmon SNP dataset, $r_{DH}$ is 0.73 while $r_{GH}$ is

534  only 0.04. This means $D$ always increases with $H_S$, even for markers with low mutation rate

535  (e.g. SNPs) and low diversity, and for a population with a high migration rate.

536      Slatkin's (1995) $R_{ST}$ provides unbiased estimates of population differentiation

537  regardless of the mutation rates or diversity of markers. A mutation does not erase the

538  evolutionary history of a gene when it occurs in some models such as the stepwise model.

539  Mutations occurring in these models are accommodated by $R_{ST}$, which therefore measures

540  differentiation purely due to population demography ($m$ and $N$). Unfortunately, however, $R_{ST}$

541  is sensitive to violations of the assumed mutation models and have a higher sampling

542  variance than $G_{ST}$ (Balloux & Lugon-Moulin 2002). Unless many (say in the hundreds)

543  markers are used, $R_{ST}$ may have a lower accuracy than $G_{ST}$.

544      What are the uses of a correlation analysis on $G_{ST}$ and $H_S$? What we are usually

545  interested are population level forces such as migration (or isolation) and drift, which have

546  roughly the same effect on all loci in the genome, and population differentiation, which

547  depends on population level forces and is estimated by all loci mainly controlled by

548  population level forces. $G_{ST}$ always faithfully reflects the differentiation at the marker loci, no

549  matter the loci are governed primarily by population demography ($m$ and $N$) or locus specific

550  forces such as selection and mutation. Marker $G_{ST}$ provides an unbiased and good estimate of

551  population differentiation only when these markers are not significantly affected by locus

552  specific forces. The correlation analysis essentially tests whether different markers give

553  replicated or different estimates of $G_{ST}$, or whether or not population level forces are much

554  more important than locus specific forces in shaping the marker diversity and distribution. A

555  highly negative correlation between $G_{ST}$ and $H_S$ values indicates that 1) the migration rate

556  must be low, lower than the mutation rate; 2) the marker $G_{ST}$ may well underestimate

557  population differentiation; 3) another set of markers with lower (higher) polymorphisms may

558  well yield a higher (lower) estimate of $G_{ST}$; 4) the marker $G_{ST}$ should be used cautiously in

559  comparisons across species, studies and sets of loci. If the correlation between $G_{ST}$ and $H_S$

560    values among loci is small and non-significant, then these single locus $G_{ST}$ estimates should

561    be marker (diversity) independent and can be averaged to provide a good estimate of

562    population differentiation.

563          A computer program, **CoDiDi** (**Co**rrelation between **Di**versity and **Di**ferentiation), is

564    written to calculate single locus $G_{ST}$ and $H_S$ values, to test whether a single locus $G_{ST}$ value is

565    significantly different from 0 or not by permutations, and to calculate and test the

566    significance of the correlation between $G_{ST}$ and $H_S$. The correlation analyses of all of the

567    simulated and empirical data presented in this study were conducted by this program, freely

568    available from the website: http://www.zsl.org/science/software/CoDiDi.

569

573

574    **References**

575    Balloux F, Lugon-Moulin N (2002) The estimation of population differentiation with

576        microsatellite markers. *Molecular Ecology*, **11**, 155-165.

577    Balloux F, Brunner H, Lugon-Moulin N, Hausser J, Goudet J (2000) Microsatellites can be

578        misleading: an empirical and simulation study. *Evolution*, **54**, 1414–1422.

579    Carreras-Carbonell J, Macpherson E, Pascual M (2006) Population structure within and

580        between subspecies of the Mediterranean triplefin fish *Tripterygion delaisi* revealed by

581        highly polymorphic microsatellite loci. *Molecular Ecology*, **15**, 3527–3539.

582    Dick S, Shurin JB, Taylor EB (2014) Replicate divergence between and within sounds in a

583        marine fish: the copper rockfish (*Sebastes caurinus*). *Molecular Ecology,* **23**, 575-590.

584    Hedrick PW (1999) Perspective: highly variable loci and their interpretation in evolution and

585        conservation. *Evolution*, **53**, 313–318.

586    Hedrick PW (2005). A standardized genetic differentiation measure. *Evolution*, **59**, 1633–

587        1638.

588    Heller R, Siegismund H (2009) Relationship between three measures of genetic

589        differentiation $G_{ST}$, $D_{EST}$ and $G'_{ST}$ : how wrong have we been? *Molecular Ecology*, **18**,

590        2080–2083.

591  Jin L, Chakraborty R (1995) Population structure, stepwise mutation, heterozygote deficiency
592      and their implications in DNA forensics. *Heredity*, **74**, 274-285.

593  Jost L (2008) $G_{ST}$ and its relatives do not measure differentiation. *Molecular Ecology,* **17**,
594      4015–4026.

595  Kimura M, Crow JF (1964) The number of alleles that can be maintained in a finite
596      population. *Genetics*, **49**, 725-738.

597  Kimura M, Weiss G (1964) The stepping-stone model of population structure and the
598      decrease of genetic correlation with distance. *Genetics,* **49**, 561-576.

599  Li W-H (1976) Effect of migration on genetic distance. *American Naturalist*, **110**, 841–847.

600  Meirmans PG, Hedrick PW (2011) Assessing population structure: $F_{ST}$ and related measures.
601      *Molecular Ecology Resources,* **11**, 5–18.

602  Moore JS, Bourret V, Dionne M *et al*. (2014) Conservation genomics of anadromous Atlantic
603      salmon across its North American range: outlier loci identify the same patterns of
604      population structure as neutral loci. *Molecular Ecology*, **23**, 5680-5697.

605  Nagylaki T (1998) Fixation indices in subdivided populations. *Genetics*, **148**, 1325–1332.

606  Nei M (1973) Analysis of gene diversity in subdivided populations. *Proceedings of the*
607      *National Academy of Sciences of the United States of America,* **70**, 3321–3323.

608  Nei M (1975) *Molecular Population Genetics and Evolution*. North-Holland, Amsterdam,
609      Netherlands.

610  Nei M (1977) *F*-statistics and analysis of gene diversity in subdivided populations. *Annals of*
611      *Human Genetics*, **41**, 225-233.

612  Nei M, Chesser R (1983) Estimation of fixation indices and gene diversities. *Annals of*
613      *Human Genetics*, **47**, 253–259.

614  Ohta T, Kimura M (1973) A model of mutation appropriate to estimate the number of
615      electrophoretically detectable alleles in a finite population. *Genetical Research*, **22**, 201-
616      204.

617  Portnoy DS, Hollenbeck CM, Belcher CN *et al*. (2014) Contemporary population structure
618      and post-glacial genetic demography in a migratory marine species, the blacknose shark,
619      *Carcharhinus acronotus*. *Molecular Ecology*, **23**, 5480-5495.

620  Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using
621      multilocus genotype data. *Genetics*, **155**, 945-959.

622  Ryman N, Leimar O (2008) Effect of mutation in genetic differentiation among
623      nonequilibrium populations. *Evolution*, **62**, 2250–2259.

624 Ryman N, Leimar O (2009) $G_{ST}$ is still a useful measure of genetic differentiation—a
625      comment on Jost's *D. Molecular Ecology,* **18**, 2084–2087.

626 Sanetra M, Crozier R (2003) Patterns of population subdivision and gene flow in the ant
627      *Nothomyrmecia macrops* reflected in microsatellite and mitochondrial DNA markers.
628      *Molecular Ecology*, **12**, 2281–2295.

629 Schiavina M, Marino IAM, Zane L, Melià P (2014) Matching oceanography and genetics at
630      the basin scale. Seascape connectivity of the Mediterranean shore crab in the Adriatic
631      Sea. *Molecular Ecology*, **23**, 5496-5507.

632 Slatkin M (1995) A measure of population subdivision based on microsatellite allele
633      frequencies. *Genetics*, **139**, 457–462.

634 Takahata N, Nei M (1984) $F_{ST}$ and $G_{ST}$ statistics in the finite island model. *Genetics*, **107**,
635      501-504.

636 Vignaud TM, Mourier J, Maynard JA *et al*. (2014). Blacktip reef sharks, *Carcharhinus*
637      *melanopterus*, have high genetic structure and varying demographic histories in their
638      Indo-Pacific range. *Molecular Ecology*, **23**, 5193-5207.

639 Wang J (2012) On the measurements of genetic differentiation among populations. *Genetics*
640      *Research*, **94**, 275-289.

641 Wang J, Caballero A (1999) Developments in predicting the effective size of subdivided
642      populations. *Heredity*, **82**, 212-226.

643 Weir BS, Cockerham CC (1984) Estimating *F*-statistics for the analysis of population
644      structure. *Evolution,* **38**, 1358–1370.

645 Whitlock MC (2011) $G_{ST}$ and *D* do not replace $F_{ST}$. *Molecular Ecology,* **20**, 1083–1091.

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

647 Wright S (1943) Isolation by distance. *Genetics,* **28**, 114-138.

648 Wright S (1951) The genetical structure of populations. *Annals of Eugenics*, **15**, 323–354.

649 Wright S (1969) *Evolution and the Genetics of Populations, Vol. II. The Theory of Gene*
650      *Frequencies.* University of Chicago Press, Chicago.

651 Wright S (1978) *Evolution and the Genetics of Populations, Vol. IV. Variability Within and*
652      *Among Natural Populations*. University of Chicago Press, Chicago.

653

654

655  J. Wang is interested in developing population genetics models and methods of analysis of

656  empirical data to address issues in evolutionary and conservation biology.

657  _____

658

659  **Data accessibility**

660  The computer program for simulating genotype data under different migration and mutation

661  models, and for estimating $H_S$, $G_{ST}$ and their correlation: Dryad DOI: 10.5061/dryad.733s9.

662  The 6 empirical datasets were retrieved from Dryad with DOIs:
663  http://dx.doi.org/10.5061/dryad.sb601; http://dx.doi.org/10.5061/dryad.vv277;
664  http://dx.doi.org/10.5061/dryad.r0d1q; http://dx.doi.org/10.5061/dryad.th4h5;
665  http://dx.doi.org/10.5061/dryad.s489b

666  The input files of the 6 empirical datasets for **CoDiDi** analysis: Dryad DOI:

667  10.5061/dryad.733s9.

668

669 **Table 1** Parameter ranges in simulations

| Migration model | Mutation Model | $t$ | $m$ | $u$ | $M$ | $R$ | $L$ |
|---|---|---|---|---|---|---|---|
| FIM, SSM | IAM, SWM, FAM | 200, $\propto$ | 0.01, 0.001 | $10^{-5}$~$10^{-3}$ | 10, 20, 40, 80, 160 | 2, 3, 4, 6, 8, 10, 12 | 5, 10, 15, 20, 30 |

670 The size ($N$) and number ($s$) of subpopulations are fixed at 250 (or 1000) and 20, respectively.

671 The finite island model (FIM) and circular stepping stone model (SSM) for migrations are

672 considered for neutral loci under infinite allele model (IAM), stepwise model (SWM) or

673 finite allele model (FAM) for mutations. For FAM, 2 alleles are considered to mimic SNPs.

674 Symbols $t, m, u, M, R, L$ represent number of generations when sampling occurs, migration

675 rate, mutation rate, number of individuals sampled from a subpopulation, number of sampled

676 subpopulations, and number of sampled loci, where $t=\propto$ indicates a population at mutation-

677 drift-migration equilibrium.
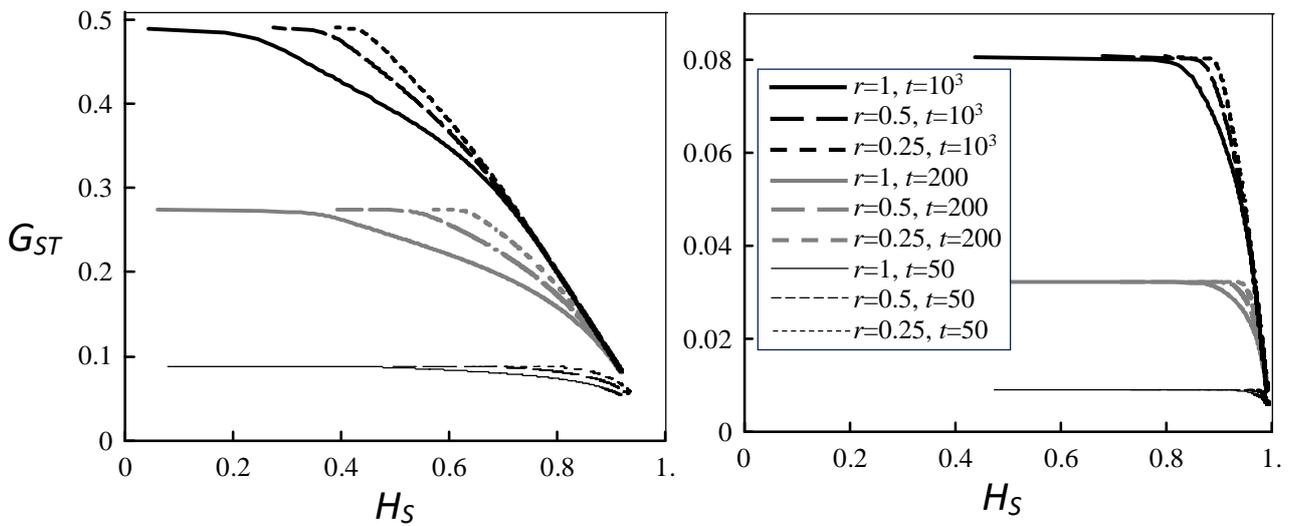
678

679



680

681

682

683

684

685

686

687

688     **Fig. 1** $G_{ST}$ as a function of $H_S$ in equilibrium populations. The $G_{ST}$ ($y$ axis) and $H_S$ ($x$ axis)

689     values for a population in a finite island model with $s=10$ subpopulations at mutation-drift-

690     migration equilibrium were calculated for various parameter values of subpopulation size ($N$),

691     migration rate ($m$), and mutation rate ($u$), where $u$ ranges from $10^{-6}$ (left side of $x$ axis) to $10^{-2}$

692     (right side of $x$ axis).

693



694

695

696

697

698

699

700

701

**Fig. 2** $G_{ST}$ as a function of $H_S$ in non-equilibrium populations. The $G_{ST}$ ($y$ axis) values are plotted against $H_S$ ($x$ axis) values at different generations ($t$=50, 200, 1000) for a population in a finite island model with $s$=10 subpopulations, assuming parameter values of $N$=250 (left panel) or 1000 (right panel), $m$=0.001, and a variable $u$ ranging from $10^{-6}$ (left side of $x$ axis) to $10^{-2}$ (right side of $x$ axis). The initial probability of gene identity is assumed to be $rJ_{0(\infty)}$, where $r$=1, 0.5 and 0.25 and $J_{0(\infty)}$ is the equilibrium value of $J_0$ given parameters $N, m, u, s$.
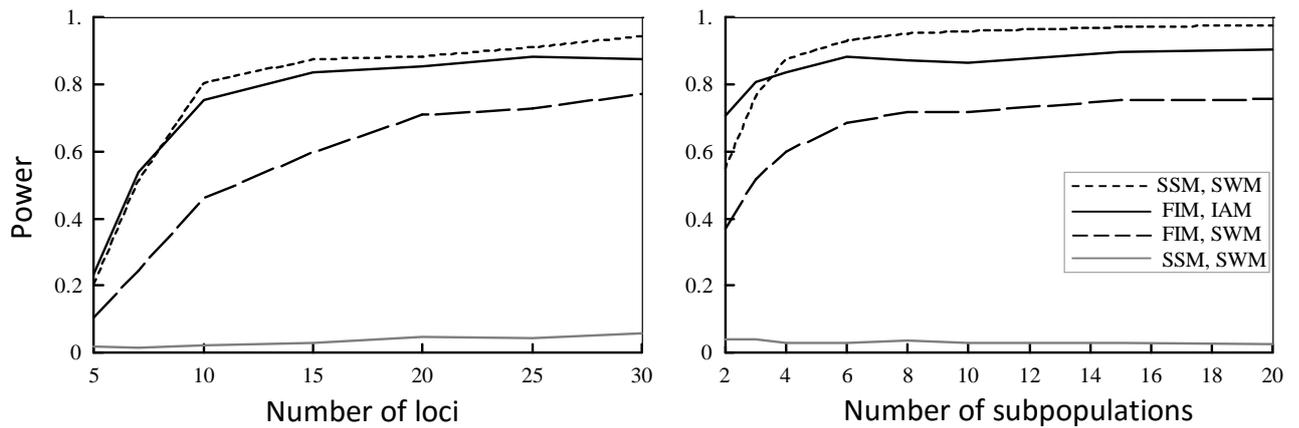
708

**Fig. 3** Scatter graphs of $G_{ST}$ ($y$ axis) and $H_S$ ($x$ axis) estimates in mutation-drift-migration equilibrium populations. The population parameters are $N$=250, $s$=20, $u$ is taken at random from a uniform distribution in the range [$10^{-5}$, $10^{-3}$], and migration rate is either $m$=0.001 (left column) or $m$=0.01 (right column). The population is assumed to follow the finite island and infinite allele models (first row), finite island and stepwise mutation models (second row), or stepping stone and stepwise mutation models (third row). For each graph, 5000 replicate simulated datasets (loci) were generated to estimate $G_{ST}$ and $H_S$, using $R$=4 (out of $s$=20) randomly sampled subpopulations and $M$=50 (out of $N$=250 or 1000) randomly sampled individuals per subpopulation. The correlation between the $G_{ST}$ and $H_S$ estimates for each graph is shown at the right corner of the graph.
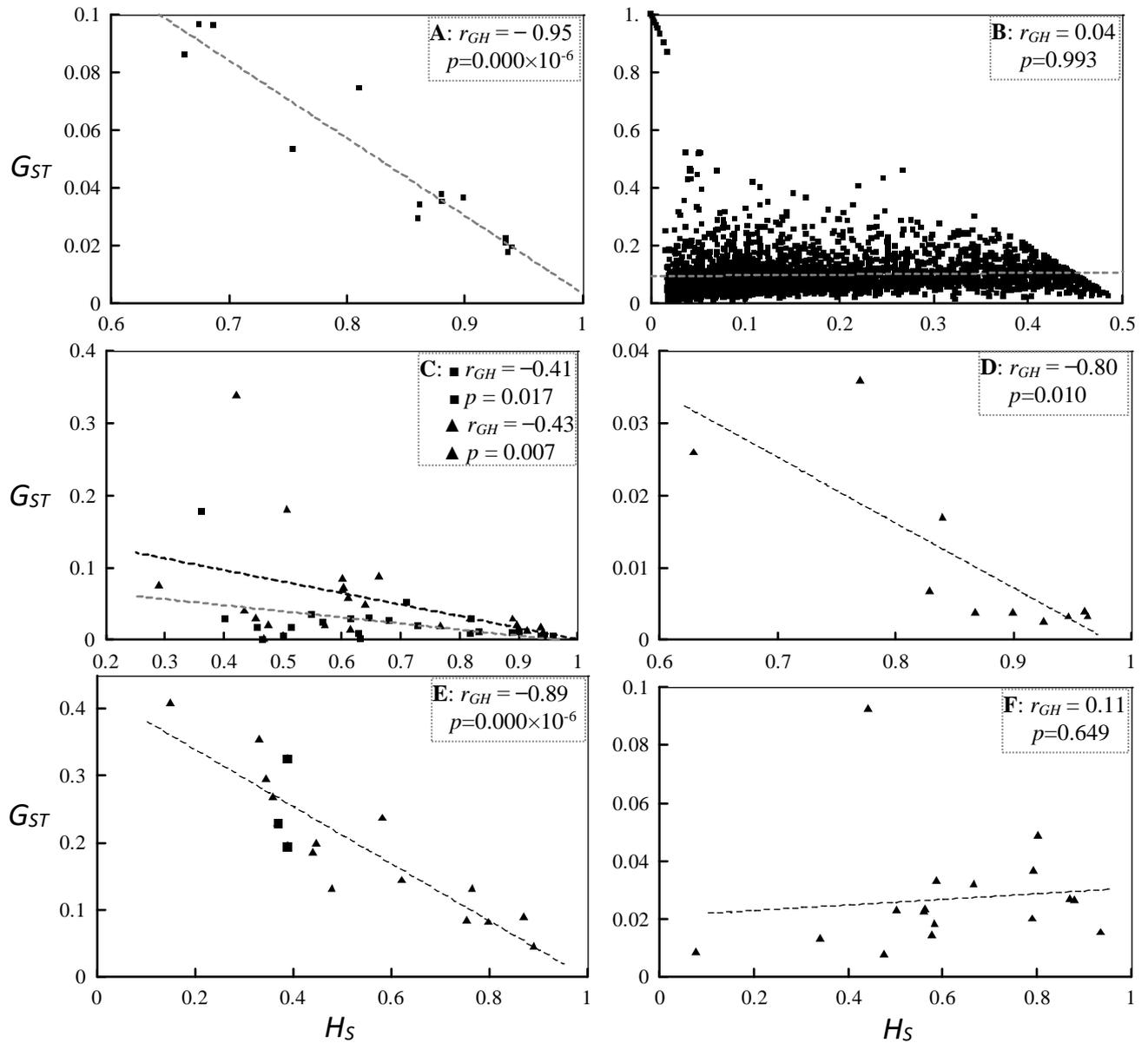
722



723
724

**Fig. 4** Power of correlation analysis of $G_{ST}$ and $H_S$ estimates in mutation-drift-migration

equilibrium populations. The population parameters are $N$=250, $s$=20, and $u$ is taken at

random from a uniform distribution in the range [$10^{-5}$, $10^{-3}$]. The numbers of sampled

subpopulations and loci are 4 and variable for the left panel, or variable and 15 for the right

panel. Migration rate is either $m$=0.001 (black continuous, black broken and black dotted

lines) or $m$=0.01 (grey continuous lines). The population is assumed to follow the finite

island model (FIM) and infinite allele model (IAM), finite island model and stepwise

mutation model (SWM), or stepping stone model (SSM) and stepwise mutation model. For

each parameter combination, the proportion of 1000 replicate datasets in which the

correlation coefficient between $G_{ST}$ and $H_S$, estimated using 40 individuals per sampled

subpopulation, is statistically significant at $p$<0.001 is plotted (on $y$ axis) as a function of the

number of sampled loci (left panel) or the number of sampled subpopulations (right panel)

(on $x$ axis). The black lines show the power in detecting mutational effects on $G_{ST}$ when such

effects exist (i.e. when migrations are weak relative to mutations, $m$=0.001), and the grey

lines show the false detection rates when mutational effects are absent (i.e. when migrations

are strong relative to mutations, $m$=0.01).

745

**Fig. 5** The relationship between single locus $G_{ST}$ and $H_S$ estimates in empirical datasets. The correlation coefficient between $G_{ST}$ and $H_S$ and the $p$ value for each dataset are shown at the top right corner of each graph, and the grey dotted lines show the fitted regression of $G_{ST}$ on

756  $H_S$. Graphs A and B show the results for the 15 microsatellites and 3129 SNPs respectively in

757  North American Atlantic salmon populations. Graph C shows the results for the 23

758  microsatellites in the blacknose shark populations, where each triangle and each square

759  shows the pair of $G_{ST}$ and $H_S$ values estimated from a single marker in the 2- and 10-

760  population models, respectively. Graph D shows the results for the 10 microsatellites in eight

761  Mediterranean shore crab populations. Graph E shows the results for the 17 microsatellites in

762  15 blacktip reef shark populations, where each triangle and each square represents a single

763  marker without and with deviation from Hardy-Weinberg equilibrium. Graph F shows the

764  results for the 17 microsatellites in 10 copper rockfish populations.

765