1

2   DR. STUART C WILLIS (Orcid ID : 0000-0002-2274-1112)

3

4

9

10

# Haplotyping RAD loci: an efficient method to filter paralogs and account for physical linkage

13

14   Authors: Stuart C. Willis[1]*, Christopher M. Hollenbeck[1], Jonathan B. Puritz[1,2], John R.

15   Gold[1], & David S. Portnoy[1]

16

17   *To whom correspondence should be addressed: swillis4@gmail.com

18

19   Authors' Institutions:

20   [1]Marine Genomics Laboratory, Department of Life Sciences, Texas A&M University-

21   Corpus Christi, 6300 Ocean Drive, Corpus Christi, Texas 78412, USA.

22   [2]Marine Science Center, Northeastern University, 430 Nahant RD, Nahant, MA 01908,

23   USA

24   **Abstract:**

25   Next-generation sequencing of reduced-representation genomic libraries provides a

26   powerful methodology for genotyping thousands of single nucleotide polymorphisms

27   (SNPs) among individuals of non-model species. Utilizing genotype data in the absence

28   of a reference genome, however, presents a number of challenges. One major challenge is

29   the tradeoff between splitting alleles at a single locus into separate clusters (loci), creating

30   inflated homozygosity, and lumping multiple loci into a single contig (locus), creating

31   artifacts and inflated heterozygosity. This issue has been addressed primarily through the

32   use of similarity cutoffs in sequence clustering. Here, two commonly employed, post -

33   clustering, filtering methods (read depth and excess heterozygosity) used to identify

34   incorrectly assembled loci are compared with haplotyping, another post-filtering

35   clustering approach. Simulated and empirical data sets were used to demonstrate that

36   each of the three methods separately identified incorrectly assembled loci; more optimal

37   results were achieved when the three methods were applied in combination. The results

38   confirmed that including incorrectly assembled loci in population-genetic datasets

39   inflates estimates of heterozygosity and deflates estimates of population divergence.

40   Additionally, at low levels of population divergence, physical linkage between SNPs

41   within a locus created artificial clustering in analyses that assume markers are

42   independent. Haplotyping SNPs within a locus effectively neutralized the physical

43   linkage issue without having to thin data to a single SNP per locus. We introduce a Perl

44   script that haplotypes polymorphisms, using data from single or paired-end reads, and

45   identifies potentially problematic loci.

46   **Keywords:** population genomics, non-model species, single nucleotide polymorphisms

47   **Introduction**

48         The field of population genetics, empowered by high-throughput DNA

49   sequencing, is rapidly expanding the potential for high resolution demographic, genomic,

50   and evolutionary analyses of non-model organisms (Mardis 2008). The technology has

51   not yet reached the point where sequencing the full genome of many samples is cost or

52   labor-efficient, so most studies rely on reduced-representation libraries to provide a

53   manageable number of single-nucleotide polymorphisms (SNPs) to survey across

54   individuals (Altshuler *et al.* 2000). Currently, there are several library-preparation

55   approaches and bioinformatics procedures used to identify and genotype hundreds to

56   thousands of SNPs in a panel of individuals (e.g. Okou *et al.* 2007; Van Tassell *et al.*

57   2008). One form of library preparation (restriction-site associated DNA or RAD) takes

58   advantage of the relative frequency of restriction endonuclease sites to tailor the number

59   of fragments sequenced (Puritz *et al.* 2014b). The major challenge for most RAD

60   sequencing projects applied to non-model organisms is to assemble a high quality set of

61   homologous sequences with minimal missing data across the greatest number of

62   individuals, without use of a reference genome (Davey *et al.* 2011). This challenge has

63   been met with many solutions and mixed degrees of success (Puritz *et al.* 2014a; Puritz *et*

64   *al.* 2014b).

65         Assembling a RAD dataset requires separation of reads into clusters

66   corresponding to a single location on a haploid set of chromosomes (hereafter, single-

67   copy locus). The challenge, therefore, is to identify highly similar sequences that occupy

68   different chromosomal locations (hereafter, multi-copy loci). These multi-copy loci

69   include paralogs, transposons, and other, non-allelic similar sequences (Hohenlohe *et al.*

70   2011; Peterson *et al.* 2012) that may artificially cluster together during assembly. There

71   are several approaches to detect multi-copy loci such as quantitative PCR (e.g. D'haene

72   *et al.* 2010) or phylogenetic analysis of homologous sequences (e.g. Cannon & Young

73   2003), but none of these are cost-effective for the volume of data typical of a RAD

74   population genetics dataset. The problem is especially challenging for taxa with recent

75   whole genome duplications followed by partial "diploidization", such as salmonids

76   (Christensen et al. 2013).

77         Identification and elimination of multi-copy loci in SNP datasets begins during

78   bioinformatics assembly and filtering. An initial step in clustering reads is to select a

79   cutoff for the number of base differences allowed among reads that are assembled into a

80   contiguous sequence alignment (contig), what will thereafter be considered as

81   corresponding to a single-copy locus (e.g. Catchen *et al.* 2011). A stringent cutoff can be

82   applied at this step to restrict the number of multi-copy loci; however, divergent alleles

83   within a single locus may be split into different contigs (over-splitting) and this can

84   inflate observed homozygosity, compromising downstream analyses that depend on

85   unbiased estimates of heterozygosity (Catchen *et al.* 2011; Harvey *et al.* 2015; Ilut *et al.*

86   2014). Alternatively, a lower sequence similarity threshold can be used to avoid over-

87   splitting and post-assembly approaches can be employed to filter the dataset and identify

88   potential multi-copy loci (Ilut *et al*. 2014).

89         One post-assembly filtering approach is based on the observation that read depths

90   derived from single-locus clusters theoretically form a distribution around a mean read

91    depth (Emerson *et al.* 2010). Contigs with abnormally high read depth often signal the

92    presence of multi-copy loci (Emerson *et al.* 2010), meaning that secondary peaks or

93    outliers in a frequency distribution of read depth per contig may indicate suspect

94    alignments and can be used to choose thresholds for single- vs. multi-copy loci. A second

95    filtering approach (Hohenlohe *et al.* 2011) relies on the occurrence of fixed or near-fixed

96    differences between non-allelic loci which causes an excess of heterozygotes above the

97    expected 50% for bi-allelic SNPs. Filters that employ this approach tend to eliminate

98    SNP loci with proportions above this level or that deviate significantly from either

99    Hardy-Weinberg or binomial expectations (Hohenlohe *et al.* 2011; Parchman *et al.* 2012).

100   A third filtering approach, haplotyping, relies on the fact that closely linked SNPs can

101   constitute haplotypes of which a diploid individual can have no more than two (Ilut *et al.*

102   2014; Peterson *et al.* 2012). Consequently, contigs that contain reads with three or more

103   haplotypes within an individual can be flagged for inspection or removed (Parchman *et*

104   *al.* 2012; Peterson *et al.* 2012). Unlike a filter for excess heterozygosity, which relies on

105   significant divergence between alleles at multi-copy loci, identifying excess haplotypes

106   within individuals only requires that there are two or more variable SNP sites within a

107   contig. The number of individuals exhibiting reads with more than two haplotypes can

108   then be used as a cut off to eliminate possible multi-copy loci. These filters are designed

109   to eliminate multi-copy or artifactual contigs from population genomic datasets, and

110   though many researchers may wish to identify true paralogous loci (potential sites of

111   evolutionary innovation) in their data, the loci identified by these filters will often result

112   from a variety of assembly and scoring errors also.

113        Closely linked SNPs can also pose complications in data analysis when

114   associations due to linkage are treated as a statistical association among loci resulting

115   from consanguinity, selection, or population structure (Kaeuffer *et al.* 2007). Over time

116   scales of population-level processes, SNPs within a fragment of a few hundred base pairs

117   in length are expected to exhibit background linkage disequilibrium (LD), and thus

118   should not be considered independent markers (Falush *et al.* 2003; Kaeuffer *et al.* 2007).

119   This presents a dilemma for researchers who wish to glean as much information as

120   possible from their data as the total observed SNPs will be greater than the number of

121   segregating loci. In addition, considering that SNPs contain less information per-locus

122 than multi-allelic markers such as microsatellite loci (Morin *et al.* 2009), thinning the

123 dataset to one SNP per locus reduces the total information content. Fortunately, the

124 information content of all SNPs in a dataset can be preserved and physical linkage

125 artifacts removed by haplotyping SNPs within segregating loci.

126 Here, we explore the efficacy of using read depth, excess heterozygosity, and

127 haplotyping, sequentially, separately, and in combination to identify multi-copy loci for

128 elimination from a SNP dataset. We evaluated filter performance by using four simulated

129 RAD datasets containing multi-copy loci, generated with a combination of either high or

130 low mutation rate and either simple or complex evolutionary history. We also evaluated

131 an empirical data set generated from a marine fish with low population structure and high

132 genetic diversity. Finally, we examined bias and precision in estimating population-

133 genetic parameters by retaining and considering all SNPs as independent loci, thinning to

134 a single SNP per contig, or haplotyping SNPs within contigs.

135

136 **Methods**

137 *Simulated RAD data*

138 Sequence reads from a double-digest RAD library (i.e., paired reads of fixed-

139 length, allelic sequences) were simulated using the *simrrls* Python script (D. Eaton,

140 Yale), creating reads of a user-specified library type. The *EggLib* library (De Mita & Siol

141 2012) was used to specify demographic parameters that affect allelic coalescence and

142 simulate sequences under those conditions. Two large, randomly mating populations that

143 diverged from a common, homogenous population $4N$ generations in the past, followed

144 by bi-directional gene exchange ($m = 0.01$) until $0.1N$ generations in the past (after this,

145 $m = 0$) were simulated, and 1,000 loci from 40 individuals (20 per population) were

146 sampled. To introduce multi-copy loci (in this case double-copy) another pair of

147 populations, with the same demographic history but which had diverged from the first

148 pair of populations $20N$ generations in the past, followed by zero gene exchange, were

149 simulated. From this second pair of populations, sequences from 50 of the 1000 loci (5%)

150 were sampled and combined with reads from the first pair of populations. The resulting

151 dataset contained 950 single-copy and 50 multi-copy loci. Simulated sequences consisted

152 of paired 100 base pair (bp) forward and reverse reads, with the number of reads per

153    locus per individual specified with a gamma distribution ($k = 1.6$, $\theta = 20$) with a mean of

154    $k*\theta = 32$, mode of $(k-1)*\theta = 12$, and a 95% probability interval of 2.6 - 97.2. These

155    simple, multi-copy datasets also included sequencing errors and insertion-deletion

156    mutations introduced at default rates ($P = 0.001$ per site). Data were simulated at lower

157    ($N = 35,000$) and higher ($N = 70,000$) population sizes, with a constant mutation rate ($\mu = $

158    $7x10^{-9}$), thus creating low and high genetic diversity, simple datasets. Simulations for the

159    larger population also included a low but positive rate of recombination within fragments

160    ($\rho = 4Nr = 10$, sites = 100). Complex multi-copy datasets also were generated to explore

161    the performance of filtering for older, more divergent multi-copy loci which may feature

162    fixed-site or nearly-fixed differences. Both sequence datasets (low/high diversity) were

163    duplicated, and for reads from the 50 multi-copy loci derived from the second pair of

164    populations the 5th G of every odd read was changed to an A and the fourth G of every

165    even read was changed to a T. While this procedure did not create fixed differences

166    between locus copies from each population pair, it increased the likelihood of divergent

167    haplotypes over *in situ* mutation alone.

168

169    *Empirical data*

170        Empirical data consisted of a reduced-representation, genomic library of red

171    drum, *Sciaenops ocellatus*, created using a modified version of the double-digest,

172    restriction-associated DNA sequencing (ddRAD) protocol of Peterson et al. (2012).  The

173    data set was composed of 100 bp paired-end reads for 40 individuals sampled from two

174    localities (Lower Laguna Madre and Sabine Lake, Texas). These localities, while

175    demographically independent over a single generation, are part of the same western

176    "regional population" of red drum (Hollenbeck 2016), and could thus be considered to

177    consist of one or two clusters of individuals. Details of library construction can be found

178    in Puritz et al. (2014a) and data be obtained from NCBI's Short Read Archive (SRA)

179    under Accession SRP041032.

180

181    *Reference construction, read mapping, variant calling, and preliminary filtering*

182        Both simulated and empirical data were processed using the *dDocent* pipeline *v*.2

183    (Puritz *et al.* 2014a) which facilitates efficient construction of a reference genome

184 (catalog of putatively-orthologous sequences), quality trimming of sequence reads,

185 alignment-based mapping of trimmed reads to the reference, and calling of polymorphic

186 positions by using a probabilistic model and considering *a priori* sampling units. For both

187 simulated and empirical data, the reference set was created from unique, untrimmed

188 sequences that were present at least twice within individuals (K1=2) and at least twice

189 among individuals (K2=2), and then clustered at no less than 80% sequence similarity

190 (c=0.8), from which a consensus sequence was derived. These parameters are expected to

191 bypass the majority of sequencing errors, which are expected to occur in only a single

192 sequence, and provide effective clustering of even divergent alleles within loci, with the

193 possibility of clustering reads from multi-copy loci with similar sequences (Ilut *et al.*

194 2014). Quality-trimmed reads were mapped by alignment to the reference consensus

195 sequences, using mapping parameter values of 1, 3, and 5 for match score, mismatch

196 cost, and gap-opening penalty, respectively. Variant calling was performed with

197 FREEBAYES (Garrison & Marth 2012) on BAM files of aligned reads. Polymorphisms

198 (which initially included complex, insertion-deletion, multi-allelic, and bi-allelic variants)

199 were filtered for quality and missing data with a combination of VCFTOOLS (Danacek *et*

200 *al.* 2011) and vcflib (E. Garrison Boston College) in addition to the filtering below (see

201 Supplemental Information).

202

203 *Multi-copy locus elimination by variant filtering and haplotyping*

204        Three approaches for post-clustering, filtering of multi-copy loci (read depth,

205 excess heterozygosity, and haplotyping) were investigated using both empirical and

206 simulated data. Full details of filtering routines are described in Supplemental

207 Information. The first (Scheme 1) was applied to individual SNPs and employed the three

208 filtering approaches sequentially in the order read depth (a), excess heterozygosity (b),

209 and haplotyping (c). In this scheme each filtering step received only data remaining after

210 a previous filtering step. Schemes 2 and 3 were applied jointly to all the SNPs in a contig

211 rather than to individual SNPs. Scheme 2 employed the three filtering approaches

212 separately (a-c); while Scheme 3 employed the three approaches separately but then

213 combined results from all three. For comparison, a fourth dataset (Scheme 4) was

214 generated with no filtering for multi-copy loci.

215    To filter multi-copy loci based on read depth (Schemes 1, 2a, and 3), SNPs were

216    filtered by mean read depth across individuals, with cutoffs determined empirically for

217    both simulated and empirical datasets (see Results and Discussion). In Scheme 1, only

218    high depth SNPs were removed; in Schemes 2a and 3, entire SNP-containing contigs

219    were removed if any of the constituent SNPs failed to pass the filter. To filter paralogs

220    based on excess heterozygosity (Schemes 1, 2b, and 3), the proportion of heterozygotes at

221    each SNP locus was estimated using VCFTOOLS. For SNPs with >50% heterozygotes, a

222    $\chi^2$ test was used to assess whether each conformed to expectations of Hardy-Weinberg

223    equilibrium (HWE) and a correction for multiple tests (Benjamini & Hochberg 1995) was

224    applied. In Scheme 1, SNPs significantly in excess of 50% heterozygotes were removed;

225    in Schemes 2b and 3, any contig with one or more SNPs in excess of 50% heterozygotes

226    and not in HWE was removed. To filter multi-copy loci based on haplotyping (Schemes

227    1, 2c, and 3), a custom Perl script was employed (Supplemental Information).  The script

228    identifies multi-SNP genotypes for each individual at each contig, compares this to a

229    catalog of haplotypes (spanning both read pairs) for each individual at each contig, and

230    flags homozygotes errantly called heterozygotes, based on genotyping error, and true

231    heterozygotes with more than two haplotypes. In addition, the script discards variants

232    observed in only one or two reads as sequencing errors. The user is able to set a cut-off

233    for the number of genotyping errors and for extra haplotype-containing individuals

234    allowed per contig, and for missing data. In this study cut-offs were set such that if one or

235    more individuals had >2 haplotypes at a contig, that contig was removed.

236    For simulated data, the number of multi-copy loci that were eliminated at each

237    step and in each filtering scheme were recorded (Table 1). For empirical data, where the

238    true number of multi-copy loci was unknown, the total number of contigs eliminated with

239    each filter was recorded (Table 2).

240

241    *Population statistics and effects of physical linkage*

242    To examine possible effects of filtering multi-copy loci and physical linkage on

243    estimates of population-genetic parameters, the empirical dataset was filtered using

244    Schemes 1, 3, and 4. For Schemes 1 and 3, the haplotyping filter was run on data with no

245    minor allele frequency (MAF) cut-off because rare alleles, while not necessarily desirable

246  for many population genetic analyses, are quite useful in identifying excess haplotypes at

247  a locus within individuals. After initial haplotype filtering, SNPs were filtered using a

248  MAF cut-off where the least common allele had to be observed at least twice in a given

249  dataset (MAF $\geq 2/2N$ alleles), and then the data were re-haplotyped (without further

250  filtering). For schemes 1 and 3, filtered datasets were thinned to a single SNP per contig

251  (the first SNP, by default) for comparison to data sets containing all filtered SNPs

252  (unthinned) and haplotypes (Table 3). For Scheme 4, only thinned and unthinned data

253  sets were compared.

254      Two simulated, simple datasets, one of low and one of high genetic diversity,

255  were generated for comparison with the empirical dataset. For both of these simulated

256  datasets, SNP loci were filtered for $\leq 95\%$ missing data for consistency with the

257  empirical dataset and then filtered using a MAF of $\geq 2/2N$. Analyses for each dataset were

258  run with and without simulated multi-copy loci (removed manually), and thinned datasets

259  were compared to unthinned datasets. After filtering with greater stringency for missing

260  data (50% vs. 95%; Supplemental Information), these datasets consisted of ~5-10% of the

261  original 1,000 contigs. Additionally, for datasets where multi-copy loci had been

262  removed, data were haplotyped for comparison to thinned and unthinned data sets as

263  above (Table 4).

264      GENODIVE (Meirmans & Van Tienderen 2004) was used to generate estimates of

265  the effective number of alleles ($A_E$) and the inbreeding coefficient ($G_{IS}$) for each of the

266  three datasets (one empirical, two simulated) and an estimate of unbiased population

267  divergence ($G''_{ST}$) between pairs of samples within datasets. $G''_{ST}$ is a measure of

268  divergence, calibrated to the maximum possible divergence given the number of alleles at

269  a locus, and consequently permits a direct comparison between bi-allelic loci (i.e., SNPs)

270  and multi-allelic loci (haplotyped contigs) (Hedrick 2005; Meirmans & Hedrick 2011).

271  Confidence intervals for $G_{IS}$ and $G''_{ST}$ were generated using 10,000 bootstrap replicates

272  across loci. Population assignment probability to two clusters (K=2) were calculated

273  using the program STRUCTURE, with the admixture model and correlated allele

274  frequencies (Pritchard *et al.* 2000). No *a priori* population membership information was

275  specified; runs consisted of 100,000 samples after 100,000 generations of burn-in.

276    Because there were two simulated populations, and two localities (from a single regional

277    population) from which empirical data were generated, assignment was estimated at K=2.

278

279    **Results**

280    *Multi-copy loci filtering of simulated data*

281         A total of 1,000 contigs from the low variability, simple and complex sequences

282    were reconstructed using *dDocent*, as were 1,000 contigs from the high variability,

283    simple sequences. In each case, the 50 multi-copy loci (contigs with reads from both

284    population pairs) were reconstructed into a single contig each, as expected (Table 1).

285    However, a total of 1002 contigs, including 950 single-copy loci and 47 of the multi-copy

286    loci, were reconstructed from the high variability, complex dataset. Of the three

287    remaining (expected) multi-copy loci, one contig contained only reads from the second

288    population pair (in effect becoming a single-copy locus). The other two expected, multi-

289    copy loci were divided into two contigs each (total of four). One was split into two

290    contigs but each contig contained reads from each population pair, while the other split

291    into two contigs where each contig contained only reads from the second population pair.

292    Hereafter these five are referred to as anomalous, multi-copy loci.

293         Results of filtering by Schemes 1-3 are shown in Table 1. Overall, filtering by

294    Scheme 3 (combined) was more effective than Scheme 1 (sequentially) and, in most

295    cases, than Scheme 2 (separately). When applied sequentially to individual SNPs

296    (Scheme 1), each filter removed data needed by the subsequent filter to identify multi-

297    copy loci, making overall filtering less effective. The three filters applied separately

298    (Scheme 2) were variously effective at eliminating multi-copy loci. The most effective

299    filter alone, excess heterozygosity, did achieve 100% success eliminating multi-copy loci

300    in the simulation involving the low-diversity, high-complexity dataset. When run

301    separately, haplotyping was the least effective filter in terms of removal of multi-copy

302    loci. However, haplotyping performed well in both high-complexity datasets and was

303    more effective than depth filtering in the high-diversity, high-complexity dataset. This is

304    due to the fact that multi-copy loci in high complexity datasets exhibited more divergent

305    haplotypes, increasing the chance of recognizing extra haplotypes within individuals.

306    Haplotyping also identified multi-copy loci not identified by the other two filters applied

307    under Scheme 3, including all five of the anomalous, multi-copy loci from the high

308    diversity, high complexity dataset.

309            Filtering in general was least effective in high-diversity datasets. This resulted

310    from less effective mapping of higher variability reads onto contigs, thus reducing clarity

311    of patterns needed to identify multi-copy loci. For example, mean depth for SNPs from

312    multi-copy loci was 48.2 (range 13.5-69.0) and 47.3 (10.1-69.0) for the simple and

313    complex, high-diversity datasets, respectively, versus 53.8 (18.3-72.6) and 53.2 (10.9-

314    72.6) for the simple and complex, low-diversity datasets, respectively. No substantial

315    difference was observed in depth for SNPs from single-copy loci (means 28.4, 28.5, 28.3,

316    28.4). This pattern can be better understood by inspecting frequency distributions of

317    mean depth across loci (Figure 1a). SNPs from multi-copy loci are shifted to the left in

318    high-diversity datasets relative to low-diversity datasets and into depth bins constituting

319    the first mode of the bi-modal distribution. Because of this shift more SNPs from multi-

320    copy loci fell below the selected depth cutoff (maximum mean of 45 reads/individual).

321    Similarly, values of and deviations between observed and expected heterozygosity were

322    smaller in high-diversity datasets (0.237/0.241 and 0.244/0.244 mean observed/expected

323    heterozygosity in simple and complex datasets, respectively) than low-diversity datasets

324    (0.272/0.255 and 0.287/0.262 mean observed/expected heterozygosity in simple and

325    complex datasets, respectively). Consequently, fewer loci exhibited excess

326    heterozygosity when tested for deviations from HWE. Finally, a higher proportion of

327    multi-copy loci with >2 haplotypes failed to be mapped within a single individual in

328    high-diversity datasets, resulting in decreased efficiency of the haplotyping filter (Table

329    1). More permissive mapping parameters were not explored here, but it is possible that

330    for datasets from populations with high genetic diversity (i.e., with a wide and

331    overlapping range of sequence divergence between and within multi-copy and single-

332    copy loci, respectively), less stringent initial mapping values would render these filters

333    more effective.

334

335    *Multi-copy loci filtering of empirical data*

336            Reference construction for the 40 red drum individuals resulted in 40,329 contigs

337    (Table 2). A total of 124,500 variants were scored from reads mapped to these reference

338 sequences, but only 79% of contigs contained variants. The average number of variants

339 per variable contig was 3.7, which made these data similar to the simulated, low-diversity

340 datasets (4.1 variants/contig) rather than simulated, high-diversity datasets (7.2

341 variants/contig). While the actual number of multi-copy loci in the empirical dataset was

342 unknown, it likely is comparable to other non-polyploid, bony fishes (e.g <5% in

343 stickleback, Ilut *et al.* 2014), and some results are still salient without this context. For

344 example, the distribution of read depth was unimodal and highly skewed (Figure 1b),

345 with some contigs exhibiting obvious depth excesses (e.g., mean 4,918 reads/individual,

346 versus an overall mode of 20). These contigs BLAST to known multi-copy loci such as

347 ribosomal RNA genes. However, the observation of a single mode made it difficult to

348 choose an effective read-depth threshold for discriminating multi-copy loci. Working

349 from the assumption that the majority of loci were single-copy, and that the observed

350 peak corresponds to the mean depth for these loci, several cutoffs meant to approximate

351 an upper confidence limit associated with the mode were examined: 2X the mode, the

352 mode plus the difference between the mode and the minimum mean depth (mode+mode-

353 min), and the $3^{rd}$ quartile. The first (2X the mode) proved to be the least stringent for this

354 dataset (read depth 40, approximately the $80^{th}$ percentile) and was chosen as the

355 experimental cutoff to potentially allow more multi-copy loci to remain in the data prior

356 to excess heterozygosity and haplotype-based filtering. As with the simulated data, these

357 filters removed fewer contigs than the depth filter, especially when applied sequentially

358 and not strictly across entire contigs (Table 2); when applied in a combined manner, the

359 heterozygosity and haplotype filters removed an additional 1,555 (of 5,912 total) contigs

360 not flagged by the depth filter. Subsequently, the frequency distribution of depth for

361 SNPs flagged by either excess heterozygosity or haplotyping was compared to the

362 unfiltered distribution in an attempt to estimate an effective cutoff for read depth. While

363 the depth distribution of flagged loci is shifted to the right as compared to the distribution

364 of all loci, and most loci with high depth are flagged by excess heterozygosity and

365 haplotyping filters (Figure 1b), 58.3% of SNPs the were below the selected experimental

366 cutoff (40). One strategy would be to remove only contigs flagged by multiple filters,

367 with the caveat that some multi-copy loci will remain (Table 1). The advantage of this

368    strategy, however, depends on the effect of retaining multi-copy loci on downstream

369    analyses.

370

371    *Linkage, haplotypes, and population parameters*

372           For the empirical dataset there was no clear difference among estimated

373    population-genetic parameters based on all SNPs, haplotypes, or thinned SNPs, despite

374    haplotypes having a higher effective number of alleles (greater heterozygosity) per locus

375    than SNPs (Table 3). Sequential versus combined filtering schemes also had little effect

376    on estimated values. Estimates of inbreeding ($G_{IS}$) were negative and of similar

377    magnitude with overlapping confidence intervals, reflecting high genetic diversity and

378    effective population size in red drum (Gold *et al.* 2001; Turner *et al.* 2002). Estimates of

379    population divergence ($G''_{ST}$) were similarly small, but confidence intervals did not

380    include zero.

381           There were larger differences among population statistics estimated from all

382    SNPs, haplotypes, and thinned SNPs for simulated datasets which had multi-copy loci

383    removed (Table 4). Population divergence estimated from haplotypes was larger than that

384    from all or thinned SNPs. This may reflect increased power to resolve divergence with

385    haplotypes or a sensitivity of $G''_{ST}$ to the number of alleles or heterozygosity (Kalinowski

386    2002; Meirmans & Hedrick 2011). $G_{IS}$ values, alternatively, while different, had wide

387    and overlapping confidence intervals, suggesting difficulty in accurately calculating a

388    precise genome-wide estimate for this parameter based on so few loci.

389           Another pattern appeared when assignment probabilities from STRUCTURE using

390    all SNPS, haplotypes, and thinned SNPs in the empirical dataset were compared. While

391    the mean level of assignment of samples into one of two clusters was small, reflecting

392    low levels of population divergence, the variance in probability of individual assignment

393    was much greater for the dataset of all SNPs than for haplotyped or thinned SNPs (Figure

394    2). This does not appear to result from the dataset of all SNPs being more informative, as

395    the thinned and all-SNPs datasets had similar $G''_{ST}$ values (0.0014 ±0.0499 vs. 0.0012

396    ±0.0484, mean ± standard deviation of thinned vs. all SNPs, respectively). Rather, when

397    the analysis was run with SNPs in tight physical linkage, artificial clusters were formed

398    on the mistaken interpretation that LD was the result of population structure. In contrast,

399  the simulated, low-diversity datasets did not show this pattern. Instead individuals were

400  assigned back to their correct group with considerably higher posterior probability (mean

401  >0.97). This reflects the higher degree of population divergence in simulated datasets

402  than in the empirical dataset, and suggests a greater opportunity for artifacts when the

403  level of population divergence is small.

404

405  **Discussion**

406  Haplotyping SNPs within a contig provides a method to remove additional multi-

407  copy loci or otherwise artifact-prone contigs from RAD datasets when used in

408  combination with depth and excess heterozygosity filters. Both simulated and empirical

409  datasets filtered with all three methods exhibited less heterozygosity than unfiltered

410  datasets, and without the added burden of splitting single-copy loci resulting from using

411  high similarity cutoffs for clustering sequences into contigs. When robust filtering, like

412  that demonstrated here, is not applied to RAD datasets without a full reference genome,

413  multi-copy loci (i.e. paralogs, transposons, and other, non-allelic similar sequences) will

414  often be retained in the final dataset and this can lead to biased results in population

415  genetic analyses. For example, there was higher heterozygosity (lower $G_{IS}$ values) in

416  datasets with no filtering of multi-copy loci as compared to those where multi-copy loci

417  had been filtered (Table 3) or manually removed (Table 4); this is likely due to SNPs

418  segregating independently in separate copies of multi-copy loci but being clustered into a

419  single contig. This artifactual heterozygosity deflated measures of overall population

420  divergence ($G''_{ST}$), although not substantially in the empirical datasets. This finding may

421  reflect a higher proportion of multi-copy loci in simulated data relative to the empirical

422  data, suggesting that artificially reduced heterozygosity is less of a problem for data

423  derived from genomes with fewer multi-copy loci. However, the percentage of multi-

424  copy loci falling below a given similarity cutoff, and therefore likely to be assembled

425  incorrectly, will generally be difficult to predict *a priori* for non-model species.

426  Nevertheless, the consequences of downward biases in estimates of inbreeding

427  and population divergence caused by retaining multi-copy loci are not easy to predict,

428  and depend on the intended purpose of the data. In situations of very low but non-zero

429  population divergence, an increase in total heterozygosity could conceivably mask

430　divergence, and would provide biased estimates of gene flow and dispersal. For analyses

431　that depend on unbiased and accurate estimates of heterozygosity or allele frequency

432　spectra, the retention of paralogous loci may be more serious. For example, analyses such

433　as genome scans depend on accurate estimates of neutral population divergence to

434　identify outliers. Artificial downward bias in estimates of global levels of divergence

435　might lead to more false positives for loci under directional selection, while multi-copy

436　loci might be identified as being under balancing selection (Foll & Gaggiotti 2008). This

437　prediction should be true regardless of the bioinformatic pipeline used to produce the

438　final marker dataset, although pipelines that reconstruct fewer multi-copy loci and less

439　often over-split alleles would naturally produce superior results in downstream analyses.

440　　　　The results indicated that haplotyping is also a straightforward way to manage

441　closely linked SNPs within a contig without loss of information content caused by

442　thinning. Ignoring linkage can produce misleading results in analyses that assume

443　observed LD is a result of demographic or evolutionary processes. This issue is

444　potentially problematic for datasets that feature high diversity within and among

445　populations and low divergence between populations, as was manifest in the clustering

446　results from STRUCTURE. These results suggest that caution is warranted when using

447　linked SNPs from populations with low expected genomic divergence to estimate

448　assignment probabilities.

449　　　　Finally, while it seems intuitive that haplotyped datasets retain more information

450　than thinned SNP datasets, population statistics in this study from filtered datasets were

451　quite similar between thinned SNP and haplotype datasets. In this case this may reflect

452　that the sheer number of SNPs recovered overcame any loss of signal associated with

453　thinning (Kalinowski 2002; Willing *et al.* 2012). However, analyses that rely on locus-

454　by-locus measures of divergence or linkage disequilibrium such as genetic mapping (e.g.

455　Ball *et al.* 2010), estimates of identity, parentage, or kinship (e.g. Lopéz Herráez *et al.*

456　2005), and LD based estimates of effective population size (e.g. Waples & Do 2010), will

457　find added benefit to haplotyping SNPs rather than thinning to a single SNP per contig

458　because of the increased discriminatory power of additional alleles per locus.

459

460　**Acknowledgements**

**References**

Altshuler DL, Pollara VJ, Cowles CR, *et al.* (2000) A SNP map of the human genome generated by reduced representation shotgun sequencing. *Nature* **407**, 513-516.

Ball AD, Stapley J, Dawson DA, *et al.* (2010) A comparison of SNPs and microsatellites as linkage mapping markers: lessons from the zebra finch (Taeniopygia guttata). *BMC Genomics* **11**.

Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B* **57**, 289-300.

Cannon SB, Young ND (2003) OrthoParaMap: Distinguishing orthologs from paralogs by integrating comparative genome data and gene phylogenies. *BMC Bioinformatics* **4**, 35.

Catchen JM, Amores A, Hohenlohe PA, Cresko WA, Postlethwait JH (2011) Stacks: Building and Genotyping Loci De Novo From Short-Read Sequences. *G3: Genes, Genomes, Genetics`* **1**, 3171-3182.

489    Christensen KA, Brunelli JP, Lamberg MJ, *et al.* (2013) Identification of single

490         nucleotide polymorphisms from the transcriptome of an organism with a whole

491         genome duplication. *BMC Bioinformatics* **14**.

492    D'haene B, Vandesompele J, Hellemans J (2010) Accurate and objective copy number

493         profiling using real-time quantitative PCR. *Methods* **50**, 262-270.

494    Danacek P, Auton A, Abecasis G, *et al.* (2011) The Variant Call Format and VCFtools.

495         *Bioinformatics* **27**, 2156-2158.

496    Davey JW, Hohenlohe PA, Etter PD, *et al.* (2011) Genome-wide genetic marker

497         discovery and genotyping using next-generation sequencing. *Nature Genetics* **12**,

498         499-510.

499    De Mita S, Siol M (2012) EggLib: processing, analysis and simulation tools for

500         population genetics and genomics. *BMC Genetics* **13**.

501    Emerson KJ, Merz CR, Catchen JM, *et al.* (2010) Resolving postglacial phylogeography

502         using high-throughput sequencing. *Proceedings of the Natinal Academy of*

503         *Sciences USA* **107**, 16196-16200.

504    Falush D, Stephens M, Pritchard JK (2003) Inference of Population Structure Using

505         Multilocus Genotype Data: Linked Loci and Correlated Allele Frequencies.

506         *Genetics* **164**, 1567-1587.

507    Foll M, Gaggiotti O (2008) A Genome-Scan Method to Identify Selected Loci

508         Appropriate for Both Dominant and Codominant Markers: A Bayesian

509         Perspective. *Genetics* **180**, 977-993.

510    Garrison E (2014) a simple C++ library for parsing and manipulating VCF files, + many

511         command-line utilities, Boston College.

512    Garrison E, Marth G (2012) Haplotype-based variant detection from short-read

513         sequencing. arXiv.

514     Gold JR, Burridge CP, Turner TF (2001) A modified stepping-stone model of population
515              structure in red drum, Sciaenops ocellatus (Sciaenidae), from the northern Gulf of
516              Mexico. *Genetica* **111**, 305-317.

517     Harvey MG, Duffie Judy C, Seeholzer GF*, et al.* (2015) Similarity thresholds used in
518              DNA sequence assembly from short reads can reduce the comparability of
519              population histories across species. *PeerJ* **895**.

520     Hedrick PW (2005) A standardized genetic differerntiation measure. *Evolution* **59**, 1633-
521              1638.

522     Hohenlohe PA, Amish SJ, Catchen JM, Allendorf FW, Luikart G (2011) Next-generation
523              RAD sequencing identifies thousands of SNPs for assessing hybridization
524              between rainbow and westslope cutthroat trout. *Molecular Ecology Resources* **11**,
525              117-122.

526     Hollenbeck CM (2016) *Genomic studies of Red Drum (Sciaenops ocellatus) in US waters*
527              Dissertation, Texas A&M University.

528     Ilut DC, Nydam ML, Hare MP (2014) Defining Loci in Restriction-Based Reduced
529              Representation Genomic Data from Nonmodel Species: Sources of Bias and
530              Diagnostics for Optimal Clustering. *BioMed Research International* **2014**, 1-9.

531     Kaeuffer R, Réale1 D, Coltman DW, Pontier D (2007) Detecting population structure
532              using STRUCTURE software: effect of background linkage disequilibrium.
533              *Heredity* **99**, 374-380.

534     Kalinowski ST (2002) How many alleles per locus should be used to estimate genetic
535              distances? *Heredity* **88**, 62-65.

536     Lopéz Herráez D, Schäfer H, Mosner J, Fries HR, Wink M (2005) Comparison of
537              Microsatellite and Single Nucleotide Polymorphism Markers for the Genetic
538              Analysis of a Galloway Cattle Population. *Verlag der Zeitschrift für*
539              *Naturforschung* **60**, 637-643.

540    Mardis ER (2008) Next-Generation DNA Sequencing Methods. *Annual Review of*
541         *Genomics and Human Genetics* **9**, 387-402.

542    Meirmans PG, Hedrick PW (2011) Measuring differentiation: Gst and related statistics.
543         *Molecular Ecology Resources* **11**, 5-18.

544    Meirmans PG, Van Tienderen PH (2004) GENOTYPE and GENODIVE: two programs
545         for the analysis of genetic diversity of asexual organisms. *Molecular Ecology*
546         *Notes* **4**, 792-794.

547    Morin PA, Martien KK, Taylor BL (2009) Assessing statistical power of SNPs for
548         population structure and conservation studies. *Molecular Ecology Resources* **9**,
549         66-73.

550    Okou DT, Steinberg KM, Middle C, *et al.* (2007) Microarray-based genomic selection for
551         high-throughput resequencing. *Nature Methods* **4**, 907-909.

552    Parchman TL, Gompert Z, Mudge J, *et al.* (2012) Genome-wide association genetics of
553         an adaptive trait inlodgepole pine. *Molecular Ecology* **21**, 2991-3005.

554    Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HH (2012) Double Digest
555         RADseq: An Inexpensive Method for De Novo SNP Discovery and Genotyping
556         in Model and Non-Model Species. *PLoS ONE* **7**, e37135.

557    Pritchard JK, Stephens M, Donnelly PJ (2000) Inference of populations structure using
558         multilocus genotype data. *Genetics* **155**, 945-959.

559    Puritz JB, Hollenbeck CM, Gold JR (2014a) dDocent: a RADseq, variant-calling pipeline
560         designed for population genomics of non-model organism. *PeerJ* **2**.

561    Puritz JB, Matz MV, Toonen RJ, *et al.* (2014b) Demystifying the RAD fad. *Molecular*
562         *Ecology* **23**, 5937-5942.

563     Turner TF, Wares JP, Gold JR (2002) Genetic Effective Size Is Three Orders of

564          Magnitude Smaller Than Adult Census Size in an Abundant, Estuarine-Dependent

565          Marine Fish (Sciaenops ocellatus). *Genetics* **162**, 1329-1339.

566     Van Tassell CP, Smith TP, Matukumalli LK*, et al.* (2008) SNP discovery and allele

567          frequency estimation by deep sequencing of reduced representation libraries.

568          *Nature Methods* **5**, 247-252.

569     Waples RS, Do C (2010) Linkage disequilibrium estimates of contemporary Ne using

570          highly variable genetic markers: a largely untapped resource for applied

571          conservation and evolution. *Evolutionary Applications* **3**, 244-262.

572     Willing E-M, Dreyer C, van Oosterhout C (2012) Estimates of Genetic Differentiation

573          Measured by FST Do Not Necessarily Require Large Sample Sizes When Using

574          Many SNP Markers. *PLoS ONE* **7**, e42649.

575

576 **Data Accessibility**

577

578 Empirical Illumina sequences data for red drum be obtained from NCBI's Short Read

579 Archive (SRA) under Accession SRP041032. Scripts for generating the simulated

580 sequence data as well as some automated filtering have been posted to github

581 (https://github.com/jpuritz/).

582 **Tables**

583

584 **Table 1. Results of filtering of simulated ddRAD datasets.** For each simulated

585 condition (low/high diversity, simple/complex), contigs were filtered sequentially by

586 depth, observed heterozygosity ($H_O$), and haplotyping (Scheme 1), filtered separately by

587 depth, heterozygosity, or haplotyping (Schemes 2a-c), or filtered in combination (Scheme

588 3). Values recorded in each filtering step are number of simulated, multi-copy loci

589 filtered divided by the total simulated, multi-copy loci available. The number of multi-

590 copy loci available to filter at each step may not necessarily match the number remaining

591 in a previous step because some number of multi-copy loci were eliminated in

592 intermediate filtering steps not directed towards multi-copy loci. The third through fifth

593 columns list the total number of contigs reconstructed by the *dDocent* pipeline, the

594 number of multi-copy loci clusters recovered, and the number of SNPs scored across all

595 clusters. The last columns are the number of simulated multi-copy loci remaining after

596 filtering and the number of those multi-copy loci observed to possess more than two

597 haplotypes.

598

599 **Table 2. Results of filtering of the empirical ddRAD dataset.** The number of reference

600 contigs and contigs containing variants ($\geq 1$ SNP) from the *dDocent* pipeline, as well as

601 the total SNPs before filtering, are shown. Rows list the number of contigs that were

602 filtered sequentially by depth, observed heterozygosity ($H_O$), and haplotyping (Scheme

603 1), filtered separately by depth, heterozygosity, or haplotyping (Scheme 2a-c), or filtered

604 in combination (Scheme 3). The number of contigs and SNPs retained with basic but no

605 multi-copy loci specific filtering also are shown (Scheme 4). For each scheme, the final

606 remaining number of contigs and SNPs with $\leq 5\%$ missing data are listed.

607

608 **Table 3. Dataset characteristics and population statistics for red drum from Lower**

609 **Laguna Madre and Sabine Lake, TX, USA.** Data were filtered for minor allele

610 frequency (MAF $>1/2N$ alleles). Results are shown from three multi-copy loci filtering

611 schemes: SNPs filtered by each method sequentially (Scheme 1), all SNPs from contigs

612 identified in combination (Scheme 3), or no multi-copy loci filtering (Scheme 4). Number

613 of remaining contigs (#contigs) and SNPs (#SNPS) for each filtering scheme are shown

614 for datasets of all SNPs, haplotypes, or thinned SNPs. Listed for each are number of

615 alleles recovered, effective number of alleles ($A_E$), and estimates and 95% confidence

616 intervals for the inbreeding coefficient ($G_{IS}$) and for population divergence ($G_{ST}$").

617

618 **Table 4. Dataset characteristics and population statistics for simulated data with**

619 **simple haplotypes.** Data from two simulations (low and high variability) are shown with

620 and without multi-copy loci removed from final datasets. Data were filtered for minor

621 allele frequency (MAF $> 1/2N$ alleles). The number of remaining contigs (#contigs) and

622 SNPs (#SNPs) are shown for datasets of all SNPs, haplotypes, or thinned SNPs. Listed

623 for each are number of alleles recovered, effective number of alleles ($A_E$), and estimates

624 and 95% confidence intervals for the inbreeding coefficient ($G_{IS}$) and for population

625 divergence ($G_{ST}''$).

626 **Figure Legends**

627

628 **Figure 1. Frequency distribution of mean number of reads per locus**

629 **(depth/coverage):** a) simulated ddRAD data with 'simple' haplotypes; and b) empirical

630 ddRAD data from red drum. Arrows in each figure indicate the chosen read-depth cutoff

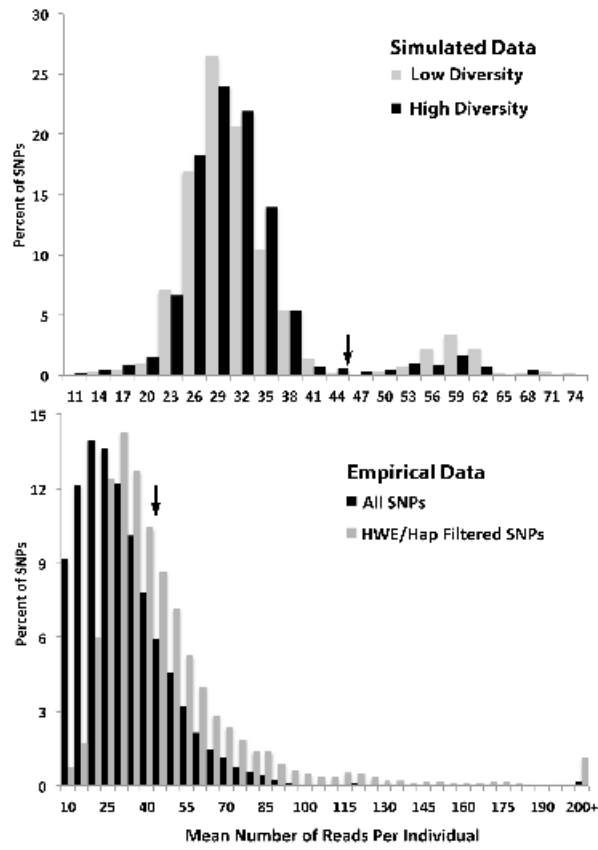631 above which contigs are flagged as multi-copy loci.

632

633 **Figure 2. Bar plots of posterior probability of individual assignment for 39 red**

634 **drum to K=2 clusters, using the program STRUCTURE for three versions of the**

635 **ddRAD dataset.**

| data diversity | multi-copy haplotypes | total contigs reconstructed | # multi-copy contigs | # SNPs | filtering scheme | filter by depth | filter by $H_o$ | filter by # haplotypes | multi-copy loci left | multi-copy loci >2 haps |
|---|---|---|---|---|---|---|---|---|---|---|
| low | simple | 1,000 | 50 | 3,641 | 1 | 30/50 (60%) | 0/15 (0%) | 2/15 (13%) | 13 | 5 |
| | | | | | 2a | 47/50 (94%) | -- | -- | 3 | -- |
| | | | | | 2b | -- | 49/50 (98%) | -- | 1 | -- |
| | | | | | 2c | -- | -- | 37/49 (76%) | 12 | 2 |
| | | | | | 3 | | | combined filters: | 0 | |
| low | complex | 1,000 | 50 | 3,714 | 1 | 28/50 (56%) | 3/18 (17%) | 1/15 (7%) | 14 | 5 |
| | | | | | 2a | 49/50 (98%) | -- | -- | 1 | -- |
| | | | | | 2b | -- | 50/50 (100%) | -- | 0 | -- |
| | | | | | 2c | -- | -- | 46/50 (92%) | 4 | 1 |
| | | | | | 3 | | | combined filters: | 0 | |
| high | simple | 1,000 | 50 | 7,097 | 1 | 17/50 (34%) | 0/32 (0%) | 4/32 (13%) | 28 | 16 |
| | | | | | 2a | 42/50 (84%) | -- | -- | 8 | -- |
| | | | | | 2b | -- | 40/50 (80%) | -- | 10 | -- |
| | | | | | 2c | -- | -- | 35/50 (70%) | 15 | 14 |
| | | | | | 3 | | | combined filters: | 7 | |
| high | complex | 1,002* | 52 (47)* | 7,187 | 1 | 16/52 (31%) | 5/36 (14%) | 7/31 (23%) | 24 | 17 |
| | | | | | 2a | 42/52 (81%) | -- | -- | 10 | -- |
| | | | | | 2b | -- | 47/52 (90%) | -- | 5 | -- |
| | | | | | 2c | -- | -- | 44/52 (85%) | 8 | 6 |
| | | | | | 3 | | | combined filters: | 2 | |

| reference contigs | # contigs ≥1 SNP | total SNPs before filtering | filtering scheme | filter by depth (2X mode) | filter by $H_o$ | filter by # haplotypes | remaining contigs (≤5%) | remaining SNPs (≤5%) |
|---|---|---|---|---|---|---|---|---|
| | | | **1** | 3,727 | 30 | 1,553 | 5,677 | 13,280 |
| | | | **2a** | 4,274 | -- | -- | 6,826 | 20,182 |
| | | | **2b** | -- | 353 | -- | 10,621 | 32,160 |
| 40,329 | 31,758 | 124,500 | **2c** | -- | -- | 2,554 | 8,332 | 20,647 |
| | | | **3** | combined filters: 5,912 | | | 5,271 | 12,664 |
| | | | **4** | no paralog filtering | | | 10,886 | 33,679 |

| multi-copy filtering | # contigs | markers | # SNPs | # alleles ($A_E$) | $G_{IS}$ (95%CI) | $G_{ST}$'' (95%CI) |
|---|---|---|---|---|---|---|
| 1. sequential | 4,932 | all SNPs | 9,964 | 19,928 (1.31) | -0.0103 (-0.0145:-0.0062) | 0.0032 (0.0019:0.0045) |
| | | haplotypes | 9,964 | 14,691 (1.61) | -0.0108 (-0.0155:-0.0060) | 0.0032 (0.0015:0.0049) |
| | | thin SNPs | 4,932 | 9,864 (1.31) | -0.0102 (-0.0162:-0.0043) | 0.0037 (0.0018:0.0057) |
| 3. combined | 4,590 | all SNPs | 9,476 | 18,952 (1.30) | -0.0094 (-0.0136:-0.0052) | 0.0030 (0.0017:0.0044) |
| | | haplotypes | 9,476 | 13,868 (1.62) | -0.0096 (-0.0142:-0.0049) | 0.0029 (0.0011:0.0047) |
| | | thin SNPs | 4,590 | 9,180 (1.31) | -0.0085 (-0.0145:-0.0025) | 0.0034 (0.0014:0.0055) |
| 4. none | 9,870 | all SNPs | 26,787 | 53,574 (1.36) | -0.0719 (-0.0764:-0.0675) | 0.0027 (0.0020:0.0035) |
| | | thin SNPs | 9,870 | 19,740 (1.34) | -0.0441 (-0.0505:-0.0377) | 0.0027 (0.0014:0.0039) |

| data diversity | multi-copy loci | # contigs | markers | # SNPs | # alleles ($A_E$) | $G_{IS}$ (95%CI) | $G_{ST}''$ (95%CI) |
|---|---|---|---|---|---|---|---|
| low | no | 55 | all SNPs | 151 | 302 (1.38) | -0.0142 (-0.0390:0.0107) | 0.2107 (0.1626:0.2591) |
| | | | haplotypes | 151 | 167 (1.68) | -0.0067 (-0.0422:0.0271) | 0.2677 (0.1972:0.3405) |
| | | | thin SNPs | 55 | 110 (1.35) | 0.0039 (-0.0428:0.0515) | 0.2656 (0.1729:0.3559) |
| low | yes | 99 | all SNPs | 474 | 948 (1.69) | -0.4592 (-0.4874:-0.4300) | 0.0782 (0.0595:0.0979) |
| | | | thin SNPs | 99 | 181 (1.58) | -0.3641 (-0.4361:-0.2891) | 0.1426 (0.0871:0.2037) |
| high | no | 80 | all SNPs | 359 | 718 (1.32) | 0.0205 (-0.0014:0.0421) | 0.2328 (0.2034:0.2617) |
| | | | haplotypes | 359 | 378 (2.16) | 0.0099 (-0.0170:0.0383) | 0.3272 (0.2736:0.3804) |
| | | | thin SNPs | 80 | 160 (1.34) | 0.0105 (-0.0303:0.0509) | 0.2148 (0.1652:0.2653) |
| high | yes | 123 | all SNPs | 753 | 1506 (1.59) | -0.3520 (-0.3755:-0.3275) | 0.1089 (0.0926:0.1256) |
| | | | thin SNPs | 123 | 246 (1.51) | -0.2582 (-0.3237:-0.1908) | 0.1360 (0.0978:0.1758) |

men_12647_f1.tif

men_12647_f2.tif