# Microhaplotypes provide increased power from short-read DNA sequences for relationship inference in fish and wildlife.

1    **Microhaplotypes provide increased power from short-read DNA sequences for**

2    **relationship inference in fish and wildlife**

3

4    Diana S. Baetscher[1,2], Anthony J. Clemento[2,3], Thomas C. Ng[2,4], Eric C. Anderson[2,3],

5    John Carlos Garza[1,2,3*],

6    [1] Department of Ocean Sciences, University of California, Santa Cruz, CA 95064, USA

7    [2] Southwest Fisheries Science Center, National Marine Fisheries Service, Santa Cruz, CA

8    95060, USA

9    [3] Institute of Marine Sciences, University of California, Santa Cruz, CA 95064, USA

10    [4] Department of Biomolecular Engineering, University of California, Santa Cruz, CA

11    95064, USA

12

13    Correspondence: 110 McAllister Way Santa Cruz CA 95060, USA; Tel.: 831-420-3903;

14    carlos.garza@noaa.gov

15

18

19    **Abstract**

20    The accelerating rate at which DNA sequence data is now generated by high-throughput

21    sequencing instruments provides both opportunities and challenges for population genetic

22    and ecological investigations of fish and wildlife. We show here how the common

23    practice of calling genotypes from a single SNP per sequenced region ignores substantial

24    additional information in the phased short-read sequences that are provided by high-

25    throughput sequencing instruments. We target sequenced regions with multiple SNPs in

26    kelp rockfish (*Sebastes atrovirens*) to determine "microhaplotypes" and then call these

27    microhaplotypes as alleles at each locus. We then demonstrate how these multi-allelic

28    marker data from 96 such loci dramatically increase power for relationship inference. The

29    microhaplotype approach decreases false positive rates by several orders of magnitude,

30    relative to calling bi-allelic SNPs, for two challenging analytical procedures, sibling and

31    single parent-offspring pair identification. The advent of phased short-read DNA

32    sequence data, in conjunction with emerging analytical tools for their analysis, promises

33    to improve efficiency by reducing the number of loci necessary for a particular level of

34    statistical confidence, thereby lowering the cost of data collection and reducing the

35    degree of physical linkage amongst markers used for relationship estimation. Such

36    advances will facilitate collaborative research and management for migratory and other

37    widespread species.

38

39   **Introduction**

40

41   The proliferation of individual-based population genetic methods in the study of ecology

42   and evolution has led to a commensurate demand for increasing analytical power. The

43   identification of first-order relatives, including parents and offspring, or full siblings, is

44   now commonplace in the study of fish and wildlife, with genotypes serving both to

45   identify relationships and as elements of larger data aggregations used in the estimation

46   of population-genetic parameter values. As the demands of such analyses grow, and

47   extend to more difficult problems of relationship estimation, making optimal use of the

48   data from high-throughput DNA sequencers is critical to achieving strong inference at

49   low cost and with wide availability.

50         High-throughput sequencing technologies have dramatically increased the rate of

51   data generation, making collection of data for genetic analysis cheaper and less time-

52   consuming. Methodological advances in both generating and analyzing these high-

53   throughput sequencing data have made it more feasible to address difficult biological

54   questions (McCormack et al. 2013, Kidd et al. 2014, Andrews et al. 2016, McKinney et

55   al. 2017). One such area of investigation that has benefited from these technological

56   developments is the identification of family relationships and pedigree reconstruction.

57   Since the inception of genetically informed relationship inference, a half-century ago,

58   researchers have applied a number of different molecular markers to the problem of

59   pedigree analysis, including allozymes, microsatellites, and most recently single-

60   nucleotide polymorphisms (SNPs). Key considerations for the utility of a molecular

61 marker include 1) variability, 2) ease of laboratory data generation, and 3) cost per

62 individual.

63     Initial studies using single-locus protein-based markers, such as allozymes, had

64 limited utility in species with low variability, with the added issue that data from these

65 markers may not be consistent with neutral expectations (Parker et al. 1998). Highly

66 polymorphic microsatellite loci quickly became the molecular marker of choice for

67 ecological studies with the widespread adoption of PCR in the early 1990s (Morin et al.

68 2004). These DNA-based markers can have large numbers of alleles and thus, high

69 information content, and became the dominant marker for exclusion-based pedigree

70 analysis (Parker et al. 1998). However, microsatellites also have many shortcomings,

71 including substantial homoplasy and high genotyping error rates (Garza and Freimer

72 1996; Morin et al. 2004, Hoffman and Amos 2005, Pemberton 2008). In addition,

73 measurement error between genotyping platforms and laboratories makes reproducibility

74 challenging (Seeb et al. 2007; Pemberton 2008) and identifying sufficiently variable

75 microsatellite loci, particularly in species with low diversity, has historically been

76 difficult (Parker et al. 1998, Pastor et al. 2008).

77     In contrast, single-nucleotide polymorphisms (SNPs) are the most abundant form of

78 variation in the genome of most species (Brumfield et al. 2003, Morin et al. 2004) and

79 their identification has become simple with the advent of high-throughput DNA

80 sequencing. In addition, SNP genotypes can be called with much less human interaction,

81 generally have low error rates, and facilitate data sharing and collaboration (Anderson

82 and Garza 2006, Seeb et al. 2009, Clemento et al. 2011). Despite the advantages of SNPs,

83 the vast majority are bi-allelic and do not provide the same per-locus power as

84      microsatellites. As such, many more SNPs than microsatellite loci are generally required

85      to provide similar power for population genetic and molecular ecological studies (e.g.,

86      Narum et al. 2008, Hauser et al. 2011, Weinman et al. 2015, Kaiser et al. 2016).

87          The huge amounts of data generated by high-throughput DNA sequencers are

88      transforming population biology, where they have helped to elucidate species

89      relationships, genetic connectivity, and ecological processes (Ekblom and Galindo 2011,

90      McCormack et al. 2013, Narum et al. 2013, Andrews et al. 2016). However, unlike

91      traditional Sanger sequencing, precise control over instrument output is challenging, so

92      most initial applications have involved collection of large amounts of data from one or a

93      small number of individuals, with sequencing reads either randomly sampling the

94      genome or a reduced fraction of it. However, many questions in population biology do

95      not require "whole genome" sequences or even the thousands of SNPs provided by most

96      reduced representation methods, such as RADseq. As such, much effort has been

97      expended to direct sequencing power to small numbers of genomic targets, allowing

98      more individuals to be studied in a single instrument run.

99          Here, we describe how data from multiple SNPs that occur within the same small

100     region, and which can be genotyped jointly from single reads from high-throughput DNA

101     sequencers, can be used to much more efficiently derive accurate relationship inference.

102     This method uses the phase information inherent in these short read DNA sequences to

103     derive multi-allelic microhaplotype markers from multiple, proximate SNPs (Kidd et al.

104     2013, 2014). We use data from a nearshore marine fish and simulation analysis to show

105     how utilizing the additional information that comes from considering all variation in

106  these short sequences provides large increases in inferential power for identifying kin

107  relationships from the same amount of DNA sequence data.

108      As sequencing instruments are limited in the total number of sequencing reads

109  produced in a single run, finding the optimal trade-off between the number of samples

110  analyzed and the number of genomic targets sequenced becomes critically important for

111  population biological studies. For questions that are extremely data-intensive, or are

112  focused specifically on genomic questions, whole genome sequencing or reduced

113  representation methods may be necessary and appropriate, but they will be prohibitive

114  when it is also necessary to analyze a large number of individuals. For projects that

115  require analysis of thousands of samples, it is important to utilize data collection methods

116  that make the most efficient use of sequencing technology, so that a modest number of

117  loci, or genomic regions, are targeted, with these loci chosen to possess high information

118  content. Multi-allelic microhaplotype markers meet this criterion and allow genotyping of

119  many more individuals in a sequencing run, since many fewer such loci are necessary to

120  achieve the same power than when just calling SNPs from such DNA sequence data.

121      Kidd et al. (2013, 2014) provided a proof of concept that microhaplotype markers

122  exist in the human genome and are useful for forensic and pedigree-type questions.

123  Gattepaille and Jakobsson (2015) showed analytically and empirically that such

124  microhaplotypes increase the power for assignment of individuals to population of origin,

125  a result that was extended by McKinney et al. (2017) for natural populations of salmon.

126      We expand on this concept by describing a set of microhaplotype loci in an organism

127  without a reference genome, kelp rockfish (*Sebastes atrovirens*), a Pacific Ocean

128  nearshore species of ecological and cultural importance. We then show how targeting

129 gene regions with abundant natural variation allows development of a 96 locus

130 microhaplotype panel with sufficient power for difficult relationship inference problems,

131 including accurately matching single parents and offspring, and identifying full-sibling

132 pairs. We show how these microhaplotypes have significantly higher heterozygosity than

133 96 SNPs from the same data set and provide much more power for pedigree inference.

134 While hundreds of SNP loci would be necessary to achieve similar accuracy, the panel of

135 96 microhaplotypes provides sufficiently low error rates for even the largest studies. We

136 highlight how microhaplotypes will substantially increase the power for population

137 genetic and ecological applications, and will be particularly useful for studies that require

138 genetic markers that are easily genotyped and portable among laboratories that use

139 benchtop sequencers to generate data. Microhaplotypes will substantially increase the

140 efficiency of genotyping, provide greater analytical power, lowering costs and potentially

141 enhancing collaboration and coordination in the study, management and conservation of

142 fish and wildlife species.

143

144 **Methods**

145 *Samples*

146 Tissue samples were obtained from field collections of rockfishes sampled at sites

147 throughout Carmel and Monterey Bays, CA. Adult kelp rockfish were sampled by hook-

148 and-line capture and removal of small caudal fin clip samples or non-lethal, underwater

149 pole-spear biopsy (J. Smith, pers. comm.), and were subsequently dried on blotting paper.

150 Genomic DNA was extracted from the dried tissue using DNeasy 96 Blood and Tissue

151    kits on a BioRobot 3000 (Qiagen, Inc.) using an elution volume of 200 µl, with DNA

152    extracts stored at 4°C until analysis.

153

154    *SNP discovery and amplicon design*

155    To identify sufficient nucleotide variation in kelp rockfish for design of microhaplotype

156    markers, we used reduced-representation genome sequencing to generate data from

157    which we could design small amplicons (100-130 bp) containing multiple SNPs. We

158    performed double-digest restriction site-associated DNA sequencing (ddRADseq;

159    Peterson et al. 2012) on 20 adult kelp rockfish. DNA concentration was normalized

160    across individuals and samples were digested with two restriction enzymes, Sph1 and

161    EcoR1, with all other details of the library preparation as in Peterson et al. (2012). We

162    selected 350 bp genomic fragments using a Pippin Prep (Sage Science) and sequenced 12

163    samples in one run and eight samples in a second run on a MiSeq (Illumina, Inc.) using

164    600-cycle paired-end sequencing kits.

165        Additionally, several loci were identified from publicly available expressed sequence

166    tags (ESTs) in an approach analogous to that used to discover and validate SNPs in other

167    fish species (e.g., Clemento et al. 2011, Abadía-Cardoso et al. 2011). We selected 192

168    ESTs for screening by PCR to determine those that effectively amplified. We then

169    generated Sanger sequence data for each locus from two kelp rockfish individuals to

170    identify variants.

171        Initial analysis of the ddRAD Illumina sequencing data with Stacks v1.34 (Catchen et

172    al. 2013) identified 17,991 gene regions in the 20 kelp rockfish samples, where each

173    region should correspond to a unique DNA sequence. We then filtered the Stacks-

8

174    assembled gene regions according to two criteria: 1) the presence of at least one SNP, and

175    2) genotyping data present in at least eight samples. This filtering reduced the dataset to

176    3,517 gene regions. To ensure that amplicon design targeted unique gene regions (e.g., no

177    repetitive elements), we used BLAT — the BLAST-Like Alignment Tool (Kent 2002) —

178    to perform pairwise comparisons of each gene region with every other region and

179    removed likely duplicates (those with greater than 95% similarity).

180        We then filtered the remaining sequences for 1) multiple SNPs within 100-130 bases

181    and 2) presence of multiple haplotypes observed across the 20 kelp rockfish sequenced.

182    From the remaining 2,333 gene regions, we selected 192 small gene regions (< 200 bp)

183    for amplicon design. We targeted regions < 200 bp because such short regions appear to

184    amplify more uniformly in multiplex reactions than larger DNA fragments (unpublished

185    data). We then designed PCR primers for candidate microhaplotype markers using

186    Primer3 software in Geneious v7.1.7 (Kearse et al. 2012). Eight of these gene regions

187    came from ESTs and 184 from our genomic sequencing data.

188

189    *Amplicon sequencing*

190    We used Genotyping-in-Thousands by Sequencing (GT-seq; Campbell et al. 2015) to

191    generate sequence data for haplotype calling. Briefly, we used an initial multiplex PCR to

192    select amplicon sequences from genomic DNA in each sample. We performed multiplex

193    PCR with primers for 96 amplicons targeting DNA from 96 adult kelp rockfish in each

194    reaction. The locus-specific primers were designed to include priming sites for the

195    sequencing reactions, which allows the sequencing instruments to recognize start

196    locations for sequencing. A second PCR added individual-specific indexes (DNA

197    barcodes) that allow sequences to be identified to individual samples during

198    bioinformatic analysis. After both PCRs, DNA concentration was normalized across

199    samples to minimize variation in number of sequencing reads per individual. Post-

200    normalization, indexed samples were combined and the sequencing library was

201    quantified by Qubit Fluorometer (Thermo Fisher Scientific) and then by qPCR with the

202    Illumina Library Quantification Kit (Kapa Biosystems). Finally, we sequenced the library

203    on a MiSeq instrument using a paired-end approach and 150-cycle sequencing kit.

204    We tested 192 loci in two sets of 96 amplicons per sequencing run, with 96 DNA

205    samples each. We replicated the first sequencing run with 48 of the same samples to

206    evaluate consistency across sequencing runs and substituted half of the samples with 48

207    different individuals from the same collection to check for consistency of loci across

208    samples. For the second set of 96 amplicons, we dropped three of the loci in the replicate

209    run due to high read depth. These four sequencing runs provided variation information

210    for a total of 144 individuals and each run produced 23.8-27.6 million reads that passed

211    filter.

212

213    *Bioinformatic processing*

214    Sequencing reads for each sample were grouped by index with the MiSeq Analysis

215    Software (Illumina), paired-end reads were combined using the Fast Length Adjustment

216    of SHort reads (FLASH; Magoc and Salzberg 2011) and then mapped to a reference file

217    of consensus sequences using the Burrow-Wheeler Aligner (BWA-MEM; Li and Durbin

218    2009). Mapped reads were converted from Sequence Alignment/Map (SAM) files to

219    Binary Alignment/Map (BAM) files with SAMtools (Li et al. 2009) and then FreeBayes

220    (Garrison and Marth 2012) was used to call variants with settings that did not include an

221    input set of variants, multi-nucleotide polymorphisms, or complex variation (composites

222    of other types of variation). FreeBayes outputs a variant call format (vcf) file with

223    cumulative information about the position of each SNP in each locus from the 144

224    rockfish evaluated.

225        Existing software was unable to reliably assemble haplotypes from specified variants,

226    primarily due to the large number of reads per locus. Accordingly, we developed

227    MICROHAPLOT, a novel program that easily imports amplicon data containing

228    microhaplotypes and allows filtering based on read depth and allelic ratio (the ratio of the

229    most frequent haplotype to the second-most-frequent haplotype in an individual at a

230    locus) before outputting individual haplotypes (Ng et al., DOI: 10.5281/zenodo.820110).

231    MICROHAPLOT uses a reference vcf file containing target sites, a Sequence Alignment

232    Map (SAM) file for each sample, and "population" information for each sample, if

233    multiple populations or species are included. We filtered data to retain genotypes with a

234    minimum of 20 reads per individual per locus and an allelic ratio of 0.1. We then

235    excluded any locus that generated data for less than 75% of samples or produced more

236    than two haplotypes per individual at a read depth threshold of 50. Loci with obvious

237    deviations from Hardy Weinberg equilibrium, as determined by plots of observed and

238    expected haplotype frequencies, were noted and removed. Finally, we removed

239    monomorphic loci; those with only one haplotype present in the 144 test samples.

240    Individual haplotypes from the 165 remaining loci were then exported from

241    MICROHAPLOT for downstream analyses.

242        To determine the utility of microhaplotypes for pedigree analyses and compare their

243    performance with bi-allelic SNPs, we generated five datasets to assess power in both

244    marker types across all 165 gene regions and with sets of 96 gene regions. These datasets

245    are as follows: microhaplotypes in all 165 loci (m165); the single SNP with the highest

246    heterozygosity in each of the 165 gene regions (s165); the 96 microhaplotypes with the

247    highest heterozygosities (m96); 96 SNPs, one each from amongst the 165 gene regions,

248    having the highest heterozygosity (s96_top); and finally, the single SNP with the highest

249    heterozygosity within the gene regions containing the best 96 microhaplotypes (s96_m).

250    We then used Monte Carlo simulation to evaluate the power to accurately identify single

251    parent-offspring pairs, and full- and half-sibling pairs using these five datasets.

252        The Monte Carlo simulations were made using CKMRsim (Anderson, DOI:

253    10.5281/zenodo.820162), an R (R Core Development Team 2016) package that

254    implements a variant of the importance-sampling algorithm of Anderson and Garza

255    (2006) tailored to pairwise relationship inference and multi-allelic markers. Briefly, in

256    CKMRsim, the genotypes of related pairs of individuals are simulated from the estimated

257    allele frequencies and the probabilities of those genotype pairs are calculated to compute

258    a log-likelihood ratio of the true relationship versus the hypothesis of no relationship.

259    Similarly, genotypes of unrelated pairs are also simulated and their log-likelihood ratios

260    computed. The simulated distributions of these log-likelihoods are used to compute the

261    false negative rates (the per-pair rate at which truly related pairs are deemed unrelated)

262    and the false positive rates (the per-pair rate at which unrelated individuals are incorrectly

263    inferred to be related) to be expected when any particular log-likelihood ratio threshold is

264    used as a criterion for classifying a pair into a given relationship, versus unrelated. The

265    importance sampling algorithm permits accurate estimation of very small per-pair false

266    positive rates ($< 10^{-10}$) which cannot be accurately estimated using conventional Monte

267    Carlo.

268        Simulations and likelihood calculations in CKMRsim were made using a genotyping

269    error model that includes allelic dropout and sequencing errors. We set the rates of the

270    errors so that, with both microhaplotypes and SNPs, the per-locus rate of calling an

271    incorrect genotype was between 0.005 and 0.01. False positive rates for parent-offspring

272    and full- and half-sibling relationships were calculated for a range of false negative rates

273    from 0.01 to 0.3. In addition, to further evaluate the power to identify half-sibling pairs,

274    we replicated two of the three 96-locus datasets (m96, s96_top) providing data sets that

275    included 1, 2, 4, 8, and 16 times as many loci (i.e., providing allele frequencies for

276    between 96 and 1536 markers) and assessed power at a single FNR value of 0.01. Finally,

277    as physical linkage between markers (even if they are not in linkage disequilibrium)

278    results in a reduction in realized power for inference of siblings (relative to using entirely

279    unlinked markers), and because close physical linkage becomes more likely with a larger

280    number of markers, we evaluated the effects of physical linkage on the power of the

281    replicated data sets for half-sibling inference. This was done by assuming a "typical

282    vertebrate genome" (25 chromosomes of between 1 and 2 Morgans in recombinational

283    length) into which loci were randomly positioned. Simulations in CKMRsim were then

284    performed assuming physical linkage using the package's ability to call the software

285    MENDEL (Lange et al. 2013).

286

287    **Results**

288    Three of the 192 loci were removed because they collectively accounted for nearly

289    73% of reads in one of the sequencing runs. For the 96 loci sequenced in replicate runs,

290    the ordinal rank of loci by number of reads from the same 48 individuals was extremely

291    strongly correlated (Spearman's coefficient = 0.99), demonstrating consistency of results

292    for individual loci across runs. In addition, genotype call rates were consistently high,

293    ranging between 92% and 96% of all locus/individual combinations (at read depth of 20)

294    in the four runs. After manually curating the remaining 189 loci in MICROHAPLOT using

295    the criteria described above, 165 loci remained for analysis. These loci contained 825

296    unique haplotypes across 144 kelp rockfish, with between two and 13 haplotypes per

297    locus (Figure 1).

298    Observed heterozygosities of the 165 microhaplotypes were substantially higher than

299    those of the single most variable SNP in each of the 165 loci (Figure 2). Mean

300    heterozygosity of the microhaplotypes was 0.41, versus 0.22 when just the SNP with the

301    highest minor allele frequency (MAF) in each locus was called. The 96 most informative

302    microhaplotype loci had a mean of 5.64 alleles (haplotypes) per locus and mean

303    heterozygosity of 0.54 (range = 0.37-0.82), whereas for the 96 most variable SNPs it was

304    0.33 (range = 0.17-0.49). Mean heterozygosity of the most variable SNPs in each of the

305    96 best microhaplotype loci was very similar to that of the 96 best SNP loci and, as such,

306    that set of polymorphisms was not evaluated further.

307    False positive rates (FPR) for identifying parent-offspring pairs and full-sibling pairs,

308    estimated using simulations, were much smaller with microhaplotypes than with SNPs

309    (Figure 3). The false negative rate (FNR) is inversely correlated with FPR, so that

310    increasing FNR decreases FPR. At FNR = 0.01, matching single parents with offspring

311    using 96 microhaplotypes resulted in an FPR of 8.43 x $10^{-11}$, whereas with the top 96

312    SNPs it was 2 x $10^{-4}$ (Figure 3a). For identifying full-siblings, also at FNR = 0.01, the

313    FPR for 96 microhaplotypes was 9.62 x $10^{-8}$, and with the 96 SNPs was 2.54 x $10^{-3}$

314    (Figure 3b). In contrast, for identifying half-siblings, considerably more power than

315    provided by the set of either 96 microhaplotypes or 96 SNPs is needed to achieve

316    acceptable false positive rates (Figure 4). With 96 microhaplotypes, the FPR, again at

317    FNR = 0.01, is 0.065, which means that more than one out of twenty comparisons of non-

318    siblings would result in a false positive identification. For the 96 SNPs, and with the

319    same FNR, FPR = 0.44, indicating an almost complete lack of power to discriminate half

320    siblings from unrelated individuals.

321        Even when the SNP dataset is expanded by a factor of four (for a total of 384 loci),

322    the half sibling FPR for SNPs decreases to only 4.6 x $10^{-3}$ (Figure 4). In contrast, when

323    the microhaplotype dataset is expanded by a factor of four, the resulting FPR at a FNR of

324    0.01 is 6.8 x $10^{-9}$, which would be adequate for all but very large studies. Moreover,

325    when taking into account physical linkage, which is unavoidable when the number of

326    markers exceeds the number of chromosome arms and reduces the independence between

327    markers for sibling inference, the apparent increase in power when adding markers is

328    reduced, relative to unlinked markers, with the reduction increasing with the number of

329    markers (Figure 4). Although the reduction is not extreme, to achieve an FPR of 1 x $10^{-9}$

330    at FNR = 0.01 in half sibling analysis, about 50 more microhaplotypes are necessary than

331    would be predicted without taking into account a typical pattern of linkage. In contrast,

332    approximately 350 additional SNPs would be necessary to achieve such additional power

333    in the face of physical linkage.

334

335 **Discussion**

336 As population genetic and molecular ecology research transitions to use of data from

337 high-throughput DNA sequencers, it is critical to determine which data collection

338 methods provide the optimal balance between the necessary amount of data per

339 individual and the maximum number of individuals that can be accommodated in each

340 instrument run. Many population genetic questions, including elucidation of patterns of

341 population structure and most relationship inference analyses, require many fewer genetic

342 markers than provided by popular reduced-representation genome sequencing approaches

343 (e.g., RAD-seq, ddRADseq, etc.). In addition, such approaches usually yield different,

344 albeit overlapping, sets of SNP markers in different sequencing runs for different

345 investigators, complicating collaboration and replication. In contrast, directed methods,

346 such as amplicon-sequencing and capture array approaches, offer the ability to optimize

347 the trade-off between the number of samples and the amount of data acquired per sample,

348 yielding datasets that are predictable and easily replicated.

349    Here, we identified short gene regions containing multiple SNPs segregating as

350 haplotypes and designed amplicons that can be easily multiplexed and sequenced using

351 such targeted protocols. These microhaplotypes contain more information than single bi-

352 allelic SNPs and offer the benefit of providing much more inferential power per locus

353 than the SNP data typically derived from high-throughput DNA sequencers. The

354 microhaplotype information is provided directly in such data, without the need for

355 probabilistic phasing, because the sequences are replicated from single molecules and

356 therefore preserve phase information for variants located in the same sequencing read.

357    This phase information allows much higher inferential and statistical power for

358    examining population biology questions, including data-intensive inference of pedigree

359    relationships, by calling multi-allelic microhaplotypes from the same sequence data

360    typically used to call bi-allelic SNPs.

361        The value of utilizing incomplete linkage disequilibrium (LD) between proximate

362    SNPs (Pakstis et al. 2012) and leveraging the phase information that comes from high-

363    throughput DNA sequencing instruments (Kidd et al. 2013, 2014) has been previously

364    recognized. Kidd et al. (2013) demonstrated that areas of the human genome where two

365    or more SNPs occur within ~200 bp are common and that the SNPs were generally not in

366    complete LD, with recombination, genetic drift and/or selection creating population

367    ancestry-informative alleles (Kidd et al. 2014).

368        Here, we extend the documentation of microhaplotype utility, by showing how

369    selecting gene regions with multiple SNPs in close proximity for use with targeted

370    sequencing approaches, including the amplicon-sequencing approach we employ, allows

371    much more power for relationship inference to be derived from the same amount of high-

372    throughput sequencing data. In the example rockfish data set, simulations demonstrated

373    that 96 microhaplotype loci generate false positive rates for single parent-offspring

374    identification on the order of $10^{-11}$, at FNR = 0.01, whereas the most informative single

375    SNPs from each of 96 loci provided false positive rates of $10^{-4}$ (Figure 3a). Similarly,

376    power for the more challenging problem of full sibling identification was substantially

377    higher with the 96 microhaplotypes (FPR on the order of $10^{-8}$) than with the 96 best SNPs

378    (FPR on the order of $10^{-3}$).

379        The much higher mean heterozygosity of the microhaplotypes compared with the

380    SNPs is indicative of their greater information content for population genetic analyses,

381    particularly relationship inference. Indeed, the simulations demonstrated that the

382    microhaplotype markers substantially outperformed the corresponding SNP loci in all

383    cases (Figure 3). While 96 SNP loci with modest mean MAF have been shown to be

384    sufficient to identify parent pair/offspring trios (Anderson and Garza 2006; Abadía-

385    Cardoso et al. 2013), single parent/offspring pair identification is considerably more

386    challenging---there is greater separation between the likelihood ratio distributions of

387    parent pair-offspring trios and unrelated trios than there is between parent-offspring pairs

388    and unrelated pairs. While the false positive rates estimated for the 96 SNPs might seem

389    low, even exceedingly small rates can lead to a large number of false positive errors,

390    because these are *per-pair* rates. The expected number of false positive errors is found by

391    multiplying the FPR by the total number of pairwise comparisons necessary. Many

392    studies, particularly in natural populations, involve very large numbers of pairwise

393    comparisons. For example, with samples from 5,000 adults and 5,000 juveniles, a single-

394    parent-offspring matching analysis involves a total of $2.5 \times 10^7$ pairwise comparisons.

395    Thus, even false positive rates of $10^{-6}$ could result in dozens of incorrectly inferred

396    pedigree relationships. With the best 96 microhaplotype loci from kelp rockfish, the false

397    positive rate for single-parent/offspring pair analysis is $8.43 \times 10^{-11}$ at a false negative

398    rate of 1%. This means that even with 100,000 parents and 100,000 offspring genotyped

399    (for a total of $10^{10}$ pairwise comparisons), less than one falsely inferred parent-offspring

400    relationship between unrelated individuals would be expected. In contrast, with this same

401    sampling scheme, and a dataset with the 96 best SNPs (FPR = 2 x $10^{-4}$ at FNR of 1%), we

402    would expect thousands of false positives (Figure 3a).

403         Similarly, Monte Carlo evaluation of the false positive rate for full-siblings

404    demonstrates the substantial increase in power obtained by using microhaplotypes rather

405    than SNPs. For example, searching for full siblings from amongst a dataset with 5,000

406    juveniles would involve nearly 1.25 x $10^{7}$ pairwise comparisons. False positive rates in a

407    study with this number of pairwise comparisons between potential full-sibling pairs and

408    using 96 microhaplotypes (FPR = 9.62 x $10^{-8}$ at FNR of 1%) is expected to produce less

409    than one false positive. In contrast, with the best 96 SNPs (FPR = 2.54 x $10^{-3}$ at FNR of

410    1%) such an analysis would potentially result in thousands of false positives, highlighting

411    how most published studies that attempt to identify pairs of full siblings have been

412    underpowered. While 96 SNPs can be sufficient for accurately identifying large full

413    sibling groups in a joint analysis (e.g., with COLONY [Wang 2004]), if only a small

414    number of sibling pairs are present in the sample, the joint analysis offers no increase of

415    power over a pairwise approach. This situation occurs frequently when sampling large

416    populations for the purpose of close-kin mark-recapture (Bravington et al. 2016).

417         Another analytical application that will benefit from increased power with

418    microhaplotypes is genetic stock identification (GSI) or individual assignment. The

419    microhaplotype panel described here is for a species without significant population

420    structure (Gilbert-Horvath et al. 2006), but utilizing the phased data from short-read

421    sequences for haplotype determination has recently been shown to increase power for

422    GSI as well, although by a much smaller margin than for relationship inference

423    applications (McKinney et al. 2017). However, with targeted ascertainment, it is feasible

424  to over-represent loci with haplotypes that have highly diagnostic frequencies across

425  different populations and with high power for identifying particular ancestry (Pakstis et

426  al. 2012; Kidd et al. 2014). Correspondingly, both initial ascertainment and

427  microhaplotype locus screening and validation would need to include appropriate

428  samples for this type of inference.

429  We used 96 loci to compare the power of microhaplotypes with SNPs, primarily

430  because of the standard 96 well microplate configuration, and the associated

431  standardization of laboratory equipment, including many traditional genotyping and

432  sequencing platforms, around this 96 well configuration. However, there is no inherent

433  constraint on the number of microhaplotype loci that can be included in a particular

434  study. One benefit of the targeted sequencing approaches, such as the GT-seq amplicon

435  protocol used here, is the ability to include any number of loci (Campbell et al. 2015), so

436  that panels of microhaplotype loci can be tailored to the study-specific requirements for

437  analytical power. In ecological and conservation studies that require genotyping a large

438  number of samples, extracting more information per sequence is easily achieved using

439  multi-allelic microhaplotypes rather than bi-allelic SNPs and will prove to be more

440  efficient and cost effective. In addition, we analyze 96 samples here for convenience, but

441  in subsequent work have shown that we can reliably generate genotypes for 384 fish at 96

442  of these microhaplotype loci in a single such sequencing run of an Illumina MiSeq—a

443  medium throughput benchtop sequencer—achieving call rates above 99%, at a minimum

444  read depth of 20, for all individual/locus combinations (unpublished data).

445  Markers with higher information content intuitively reduce the amount of genotyping

446  required for a set amount of inferential power (Rosenberg et al. 2003). Moreover, it is

447    particularly important to minimize the number of genetic markers used in the

448    identification of close kin, because of the challenge of physical linkage in some such

449    analyses. Because recombination does not occur between many loci that are on the same

450    chromosome during any single segregation event, such loci on the same chromosome do

451    not provide independent observations on relatedness for sibling relationship categories.

452    We show that the effect of physical linkage on relationship inference with a small number

453    of markers is minimal, and it is thus unlikely to greatly affect parent-offspring and full

454    sibling identification. However, with the much larger number of markers necessary for

455    half sibling analysis, linkage increases the false positive rate substantially and the

456    discrepancy becomes greater as the number of markers increases, so that hundreds of

457    additional SNP markers are necessary to account for this linkage and achieve FPR values

458    that might be necessary for studying natural populations. Furthermore, the reduced cost

459    per individual of genotyping a panel with a modest number (e.g., 96) of microhaplotype

460    loci compared to methods that target larger proportions of the genome will allow

461    enhanced monitoring and evaluation of lower-profile species, benefiting management and

462    conservation of fishes and wildlife.

463

464    **Acknowledgements**

470

## Data Accessibility

472    Consensus sequences and primer information for all 192 targeted amplicons, and

473    genotype files for all 144 kelp rockfish are deposited in Dryad, DOI:XXXXXX, and all

474    statistical analyses are documented in an R Notebook in the same repository.

475


## References

477    Abadía-Cardoso A, Anderson EC, Pearse DE, Garza JC (2013) Large-scale parentage
478         analysis reveals reproductive patterns and heritability of spawn timing in a hatchery
479         population of steelhead (*Oncorhynchus mykiss*). *Molecular Ecology*, **22**, 4733–4746.
480
481    Abadía-Cardoso A, Clemento AJ, Garza JC (2011) Discovery and characterization of
482         single-nucleotide polymorphisms in steelhead/rainbow trout, *Oncorhynchus mykiss*.
483         *Molecular Ecology Resources*, **11 Suppl 1**, 31–49.
484
485    Anderson EC, Garza JC (2006) The power of single-nucleotide polymorphisms for large-
486         scale parentage inference. *Genetics*, **172**, 2567–82.
487
488    Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA (2016) Harnessing the
489         power of RADseq for ecological and evolutionary genomics. *Nature Reviews*
490         *Genetics*, **17**.
491
492    Baird NA, Etter PD, Atwood TS *et al.* (2008) Rapid SNP discovery and genetic mapping
493         using sequenced RAD markers. *PLOS ONE*, **3**, e3376.
494
495    Bravington, MV, Skaug HJ, Anderson EC (2016) Close-kin mark-recapture. *Statistical*
496         *Science*, **31**, 259–274.
497
498    Brumfield RT, Beerli P, Nickerson DA, Edwards SV (2003) The utility of single
499         nucleotide polymorphisms in inferences of population history. *Trends in Ecology &*
500         *Evolution*, **18**, 249–256.
501
502    Campbell NR, Harmon SA, Narum SR (2015) Genotyping-in-Thousands by sequencing
503         (GT-seq): A cost effective SNP genotyping method based on custom amplicon
504         sequencing. *Molecular Ecology Resources*, **15**, 855–867.
505
506    Catchen JM (2013) Stacks: An analysis tool set for population genomics. *Molecular*
507         *Ecology*, **22**, 3124–3140.
508

509  Clemento AJ, Abadía-Cardoso A, Starks HA, Garza JC (2011) Discovery and
510      characterization of single nucleotide polymorphisms in Chinook salmon,
511      *Oncorhynchus tshawytscha*. *Molecular Ecology Resources*, **11 Suppl 1**, 50–66.
512

513  Ekblom R, Galindo J (2010) Applications of next generation sequencing in molecular
514      ecology of non-model organisms. *Heredity*, **107**, 1–15.
515

516  Garrison E, Marth G (2012) Haplotype-based variant detection from short-read
517      sequencing. *arXiv:1207.3907v2*, 9.
518

519  Garza JC, Freimer NB (1996) Homoplasy for size at microsatellite loci in humans and
520      chimpanzees. *Genome Research*, **6**, 211–217.
521

522  Gattepaille LM, Jakobsson M (2012) Combining markers into haplotypes can improve
523      population structure inference. *Genetics*, **190**, 159–174.
524

525  Gerber S, Mariette S, Streiff R, Bodénès C, Kremer A (2000) Comparison of
526      microsatellites and amplified fragment markers for parentage analysis. *Molecular
527      Ecology*, **9**, 1037–1048.
528

529  Gilbert-Horvath EA, Larson RJ, Garza JC (2006) Temporal recruitment patterns and gene
530      flow in kelp rockfish (*Sebastes atrovirens*). *Molecular Ecology*, **15**, 3801–3815.
531

532  Hauser L, Baird M, Hilborn R, Seeb LW, Seeb JE (2011) An empirical comparison of
533      SNPs and microsatellites for parentage and kinship assignment in a wild sockeye
534      salmon (*Oncorhynchus nerka*) population. *Molecular Ecology Resources*, **11**, 150–
535      161.
536

537  Hoffman JI, Amos W (2005) Microsatellite genotyping errors: detection approaches,
538      common sources and consequences for paternal exclusion. *Molecular Ecology*, **14**,
539      599–612.
540

541  Jones AG, Ardren WR (2003) Methods of parentage analysis in natural populations.
542      *Molecular Ecology*, **12**, 2511–2523.
543

544  Jones AG, Small CM, Paczolt KA, Ratterman NL (2010) A practical guide to methods of
545      parentage analysis. *Molecular Ecology Resources*, **10**, 6–30.
546

547  Kaiser SA, Taylor SA, Chen N *et al.* (2017) A comparative assessment of SNP and
548      microsatellite markers for assigning parentage in a socially monogamous bird.
549      *Molecular Ecology Resources*, **17**, 183–193.
550

551  Kearse M, Moir R, Wilson A *et al.* (2012) Geneious Basic: An integrated and extendable
552      desktop software platform for the organization and analysis of sequence data.
553      *Bioinformatics*, **28**, 1647–1649.
554

555    Kent WJ (2002) BLAT — The BLAST -Like Alignment Tool. *Genome Research*, **12**,
556        656–664.
557
558    Kidd KK, Pakstis AJ, Speed WC *et al.* (2013) Microhaplotype loci are a powerful new
559        type of forensic marker. *Forensic Science International: Genetics Supplement*
560        *Series*, **4**, e123–e124.
561
562    Kidd KK, Pakstis AJ, Speed WC *et al.* (2014) Current sequencing technology makes
563        microhaplotypes a powerful new type of genetic marker for forensics. *Forensic*
564        *Science International: Genetics*, **12**, 215–224.
565
566    Lange K, Papp JC, Sinsheimer JS, Sripracha R, Zhou H, Sobel EM (2013) Mendel: The
567        Swiss army knife of genetic analysis programs. *Bioinformatics*, **29**, 1568–1570
568
569    Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler
570        transform. *Bioinformatics*, **25**, 1754–1760.
571
572    Li H, Handsaker B, Wysoker A *et al.* (2009) The Sequence Alignment/Map format and
573        SAMtools. *Bioinformatics*, **25**, 2078–2079.
574
575    Magoč T, Salzberg SL (2011) FLASH: Fast length adjustment of short reads to improve
576        genome assemblies. *Bioinformatics*, **27**, 2957–2963.
577
578    McCormack JE, Hird SM, Zellmer AJ, Carstens BC, Brumfield RT (2013) Applications
579        of next-generation sequencing to phylogeography and phylogenetics. *Molecular*
580        *Phylogenetics and Evolution*, **66**, 526–538.
581
582    McKinney GJ, Seeb JE, Seeb LW (2017) Managing mixed-stock fisheries: genotyping
583        multi-SNP haplotypes increases power for genetic stock identification. *Canadian*
584        *Journal of Fisheries and Aquatic Sciences*, doi:cjfas–2016–0443.
585
586    Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA (2007) Rapid and cost-
587        effective polymorphism identification and genotyping using restriction site
588        associated DNA (RAD) markers. *Genome Research*, **17**, 240–248.
589
590    Morin PA, Luikart G, Wayne RK (2004) SNPs in ecology, evolution and conservation.
591        *Trends in Ecology & Evolution*, **19**, 208–216.
592
593    Narum SR, Banks M, Beacham TD *et al.* (2008) Differentiating salmon populations at
594        broad and fine geographical scales with microsatellites and single nucleotide
595        polymorphisms. *Molecular Ecology*, **17**, 3464–3477.
596
597    Narum SR, Buerkle CA, Davey JW, Miller MR, Hohenlohe PA (2013) Genotyping-by-
598        sequencing in ecological and conservation genomics. *Molecular Ecology*, **22**, 2841–
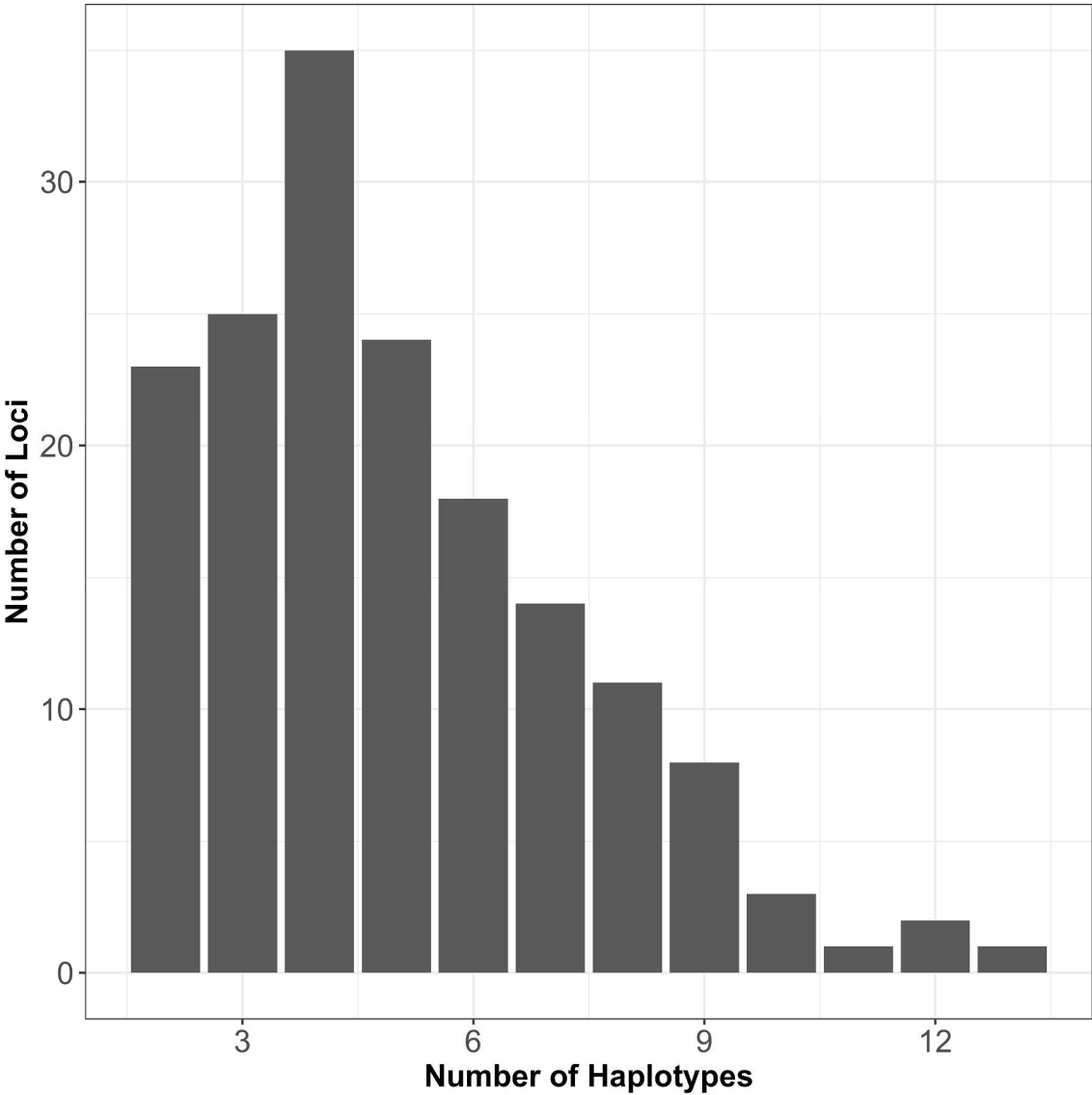599        7.
600

601   Pakstis AJ, Fang R, Furtado MR, Kidd JR, Kidd KK (2012) Mini-haplotypes as lineage
602         informative SNPs and ancestry inference SNPs. *European Journal of Human*
603         *Genetics*, **20**, 1148–1154.
604
605   Parker PG, Snow AA, Schug MD, Booton GC, Fuerst PA (1998) What molecules can tell
606         us about populations: choosing and using a molecular marker. *Ecology*, **79**, 361–
607         382.
608
609   Pastor T, Garza JC, Allen P, Amos W, Aguilar A (2004) Low genetic variability in the
610         highly endangered Mediterranean monk seal. *Journal of Heredity*, **95**, 291-300.
611
612   Pemberton JM (2008) Wild pedigrees: the way forward. *Proceedings of the Royal Society*
613         *B: Biological Sciences*, **275**, 613–621.
614
615   Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE (2012) Double digest
616         RADseq: An inexpensive method for de novo SNP discovery and genotyping in
617         model and non-model species. *PLOS ONE*, **7**, e37135.
618
619   R Core Team (2016). R: A language and environment for statistical computing. R
620         Foundation for Statistical Computing, Vienna, Austria.
621
622   Rosenberg NA, Li LM, Ward R, Pritchard JK (2003) Informativeness of genetic markers
623         for inference of ancestry. *American Journal of Human Genetics*, **73**, 1402–1422.
624
625   Seeb JE, Pascal CE, Ramakrishnan R, Seeb LW (2009) SNP genotyping by the 5'-
626         nuclease reaction: advances in high-throughput genotyping with nonmodel
627         organisms. In: *Single Nucleotide Polymorphisms: Methods and Protocols* (ed Komar
628         AA), pp. 277–292. Humana Press, Totowa, NJ.
629
630   Seeb LW, Antonovich A, Banks MA *et al.* (2007) Development of a standardized DNA
631         database for Chinook salmon. *Fisheries*, **32**, 540–552.
632
633   Wang J (2004) Sibship reconstruction from genetic data with typing errors. *Genetics*,
634         **166**, 1963–1979.
635
636   Weinman LR, Solomon JW, Rubenstein DR (2015) A comparison of single nucleotide
637         polymorphism and microsatellite markers for analysis of parentage and kinship in a
638         cooperatively breeding bird. *Molecular Ecology Resources*, **15**, 502–511.
639

640    Figure 1. Distribution of microhaplotypes across 165 loci in 144 kelp rockfish samples.
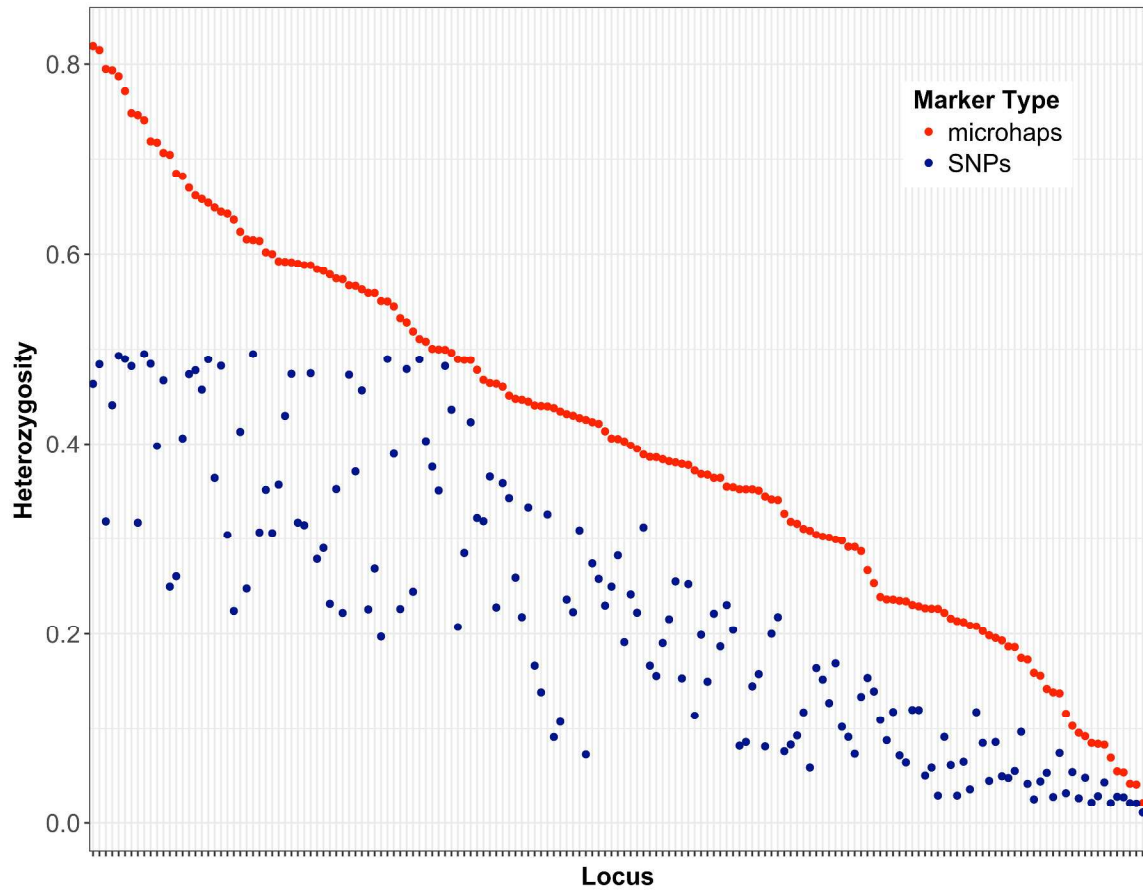641    The number of haplotypes per locus ranged from two to 13.

642

643    Figure 2. Heterozygosity of 165 microhaplotypes comprised of all SNPs in a locus
644    compared to the single SNP with the highest minor allele frequency in that same locus.
645    Bi-allelic SNPs have a maximum heterozygosity of 0.5.

646

647    Figure 3. Simulated false positive rates for matching (a) single parents with offspring and
648    for (b) full-siblings at a given false negative rate using the four sets of markers: 165
649    microhaplotypes (m165), 165 SNPs (s165), 96 microhaps (m96) and 96 SNPs with the
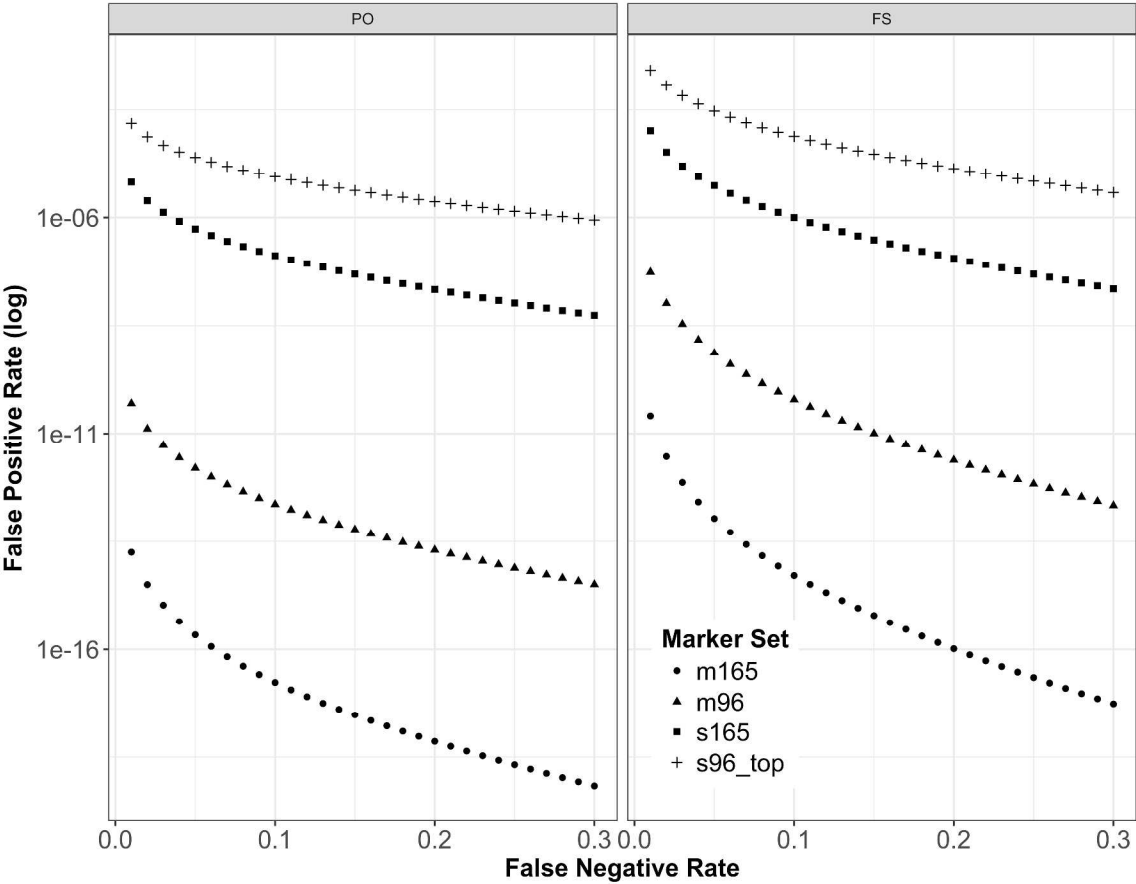650    highest heterozygosity (s96_top).

651

652    Figure 4. Simulated false positive rates for identifying half-siblings with a given number
653    of microhaplotype and SNP markers at a false negative rate of 0.1. Data are simulated
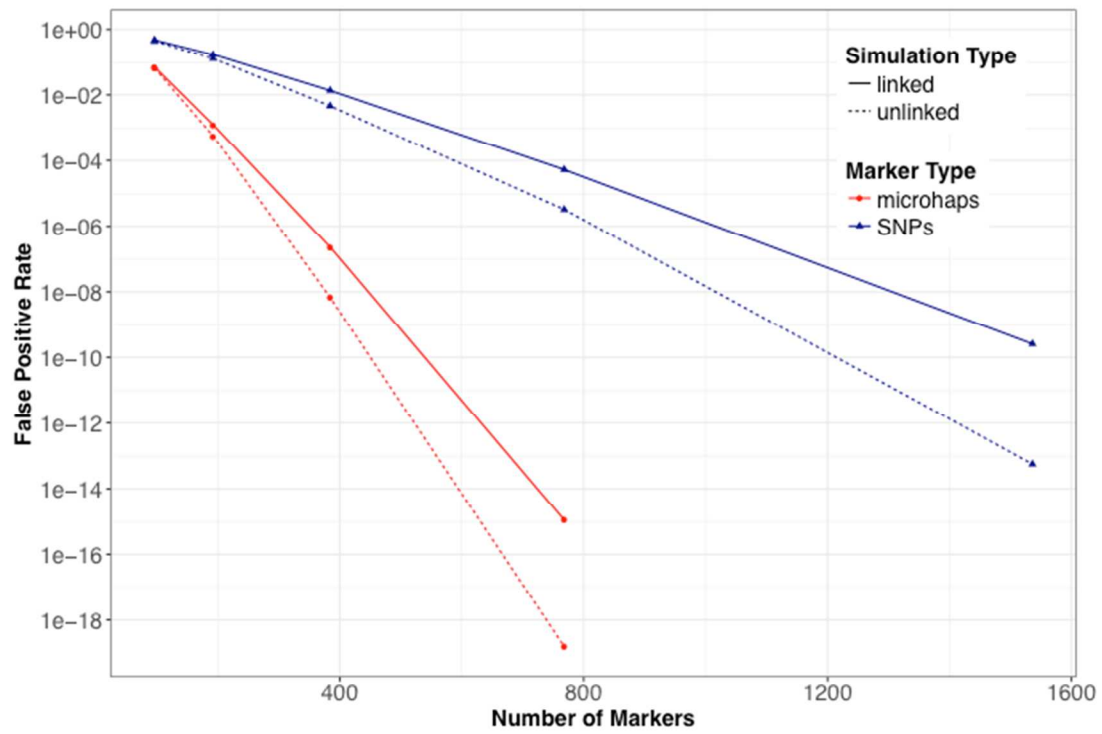654    both including and excluding physical linkage.
655

656
657

658

659
660

661