# Managing mixed-stock fisheries: genotyping multi-SNP haplotypes increases power for genetic stock identification

| | |
|---|---|
| Journal: | *Canadian Journal of Fisheries and Aquatic Sciences* |
| Manuscript ID | cjfas-2016-0443.R1 |
| Manuscript Type: | Rapid Communication |
| Date Submitted by the Author: | 07-Jan-2017 |
| Complete List of Authors: | McKinney, Garrett; University of Washington, School of Aquatic and Fishery Sciences<br>Seeb, James; University of Washington, School of Aquatic and Fishery Sciences<br>Seeb, Lisa; University of Washington, School of Aquatic and Fishery Sciences |
| Keyword: | Genetic stock identification, mixture analysis, haplotype, RADseq, Chinook salmon |

SCHOLARONE™
Manuscripts

1     **Managing mixed-stock fisheries: genotyping multi-SNP haplotypes increases**

2     **power for genetic stock identification**

3     Garrett J. McKinney[1*], James E. Seeb[1], Lisa W. Seeb[1]

4     [1] School of Aquatic and Fishery Sciences, University of Washington, 1122 NE Boat Street, Box

5     355020, Seattle WA 98195-5020, USA. (email: gjmckinn@uw.edu, jseeb@uw.edu,

6     lseeb@uw.edu)

7

8     [*]**Corresponding author:** Garrett J. McKinney (email: gjmckinn@uw.edu, phone: 1-765-430-

9     3272)

10   *Abstract*

11   A common challenge for fishery applications is resolving the relative contribution of closely

12   related populations where accuracy of genetic assignment may be limited. An overlooked

13   method for increasing assignment accuracy is the use of multi-SNP haplotypes rather than

14   single-SNP genotypes.  Haplotypes increase power for detecting population structure, and loci

15   derived from next-generation sequencing methods often contain multiple SNPs.  We evaluated

16   the utility of multi-SNP haplotyping for mixture analysis in Western Alaska Chinook salmon.

17   Multi-SNP haplotype data increased the accuracy of mixture analysis for closely related

18   populations by up to seven percentage points relative to single-SNP genotype data for a set of

19   500 loci; 90% accuracy was achievable with as few as 150 loci with multi-SNP haplotypes but

20   required at least 300 loci with single-SNP genotypes.  Individual assignment to reporting groups

21   showed an even greater increase in accuracy of up to 17 percentage points when multi-SNP

22   haplotypes were used.  Haplotyping multiple SNPs shows promise to improve the accuracy of

23   assigning unknown fish to population of origin whenever haplotype data are available.

24

25

26    *Introduction*

27    Assignment methods that use genetic data to estimate the origins of unknown mixtures of fish

28    are a powerful tool to help identify and shape the composition of the harvest in mixed-stock

29    fisheries (Beacham et al. 2008; Dann et al. 2013; Ensing et al. 2013).  Practitioners place a high

30    priority on methodological enhancements to improve resolution, especially when populations

31    contributing to the fishery are closely related.  Vast efforts are undertaken to achieve even small

32    improvements in the resolution of stocks of migratory species when diverse social interests

33    collide during the management of charismatic species in decline (e.g., Larson et al. 2014a;

34    Abadia-Cardoso et al. 2016).  Substantial efforts center on screening additional spawning

35    populations and the use of additional DNA markers scored as single nucleotide polymorphisms

36    (SNPs) to enhance both assignment efficiency and accuracy (reviewed in Seeb et al. 2011).

37    Three common approaches to genetic assignment are parentage analysis, individual assignment,

38    and mixture analysis.  These methods differ in scope of assignment: parentage analysis assigns

39    offspring to parents, individual assignment assigns individuals to populations or reporting

40    groups, and mixture analysis estimates relative composition of mixed fisheries samples

41    (reviewed in Manel et al. 2005).  Collectively, these methods have been used to study questions

42    such as composition of fishery catch (Dann et al. 2013; Bradbury et al. 2016), reproductive

43    success (Ford et al. 2015), success of reintroduction programs (Evans et al. 2016), monitoring of

44    hatchery programs (Steele et al. 2013), local adaptation (O'Malley et al. 2007), and migration

45    patterns (Larson et al. 2013).

46    Each method has different advantages; however, mixture analysis is ideally suited for many

47    practical applications in fisheries management.  Mixture analysis does not require the extensive

48    sampling of parents (required for parental assignment) and is more accurate than individual

49    assignment when differentiation between populations is low or when individual assignments

50    cannot be made with confidence (Manel et al. 2005).

51    A common challenge for fishery applications is resolving the relative contribution of closely

52    related populations.  Methods for increasing assignment accuracy include using putatively

53    adaptive loci (see Ackerman et al. 2011; Araneda et al. 2016), increasing the number of loci

54    (Larson et al. 2014b), and increasing the number of individuals and populations sampled

55    (Beacham et al. 2010; Habicht et al. 2010).  An overlooked enhancement for increasing

56    assignment accuracy is the use of multi-SNP haplotypes rather than single-SNP genotypes.

57    Haplotyping has been demonstrated to increase power for detecting population structure, where

58    haplotypes are either blocks of contiguous genomic sequence (Fariello et al. 2013) or are

59    constructed from multiple SNPs in linkage disequilibrium (Gattepaille and Jakobsson 2012;

60    Leslie et al. 2015).  To date, assignment studies have typically restricted analyses to a single SNP

61    per locus (Larson et al. 2014a; Araneda et al. 2016). Rarely, multiple linked SNPs have been

62    genotyped in two separate assays and then combined as a composite phenotype  (e.g., MHC

63    phenotypes, Habicht et al. 2010).  Discarding the data contained in multi-SNP haplotypes results

64    in a loss of information.

65    Next-generation sequencing technologies can genotype multiple SNPs per locus by directly

66    sequencing the entire DNA template in contrast to many other assays that only genotype a single

67    SNP.  Haplotype data is naturally generated by next-generation sequencing platforms; however,

68    most studies use bioinformatic filtering to examine only a single SNP per locus.  Haplotype

69    analysis involves no extra cost or effort over single-SNP genotyping but rather a simple change

70    in bioinformatic procedures.

71    In this paper we evaluate the potential for increasing accuracy of genetic stock identification

72    (GSI) by using haplotype data derived from restriction-site associated DNA sequencing

73    (RADseq; Baird et al. 2008). RADseq has become a common method for developing thousands

74    of SNP loci that can be screened for population assignment (Hohenlohe et al. 2011; Andrews et

75    al. 2016). The recent development of amplicon sequencing techniques such as GTseq (Campbell

76    et al. 2015) or RAD Capture (Rapture) (Ali et al. 2016) facilitates large scale cost-effective and

77    rapid genotyping of loci identified with RADseq, and individual RADseq loci often contain

78    multiple SNPs which allow haplotype analysis. RADseq data promise to provide many new loci,

79    including adaptively important loci, for analyses of commercially important species such as

80    Chinook salmon (*Oncorhynchus tshawytscha*) (Larson et al. 2014a; McKinney et al. 2016).

81    Chinook salmon is not only commercially important but also has sport, subsistence, and

82    ceremonial importance in Pacific North America and Asia (e.g., Wolfe and Spaeder 2009).

83    Populations of Chinook salmon show precipitous declines in many parts of their range

84    confounding managers who must balance competing interests (Gustafson et al. 2007; Gisclair

85    2009; ADF&G 2013). Managing for sustainability on the population level is highly desirable,

86    but discrete populations of Chinook salmon are often closely related, challenging the use of

87    DNA datasets for analysis of population mixtures (Templin et al. 2011; Larson et al. 2013).

88    We re-examined a previous RADseq dataset by Larson et al. (2014b) to determine if haplotype

89    analysis can increase power for mixture analysis and individual assignment over that of single-

90    SNP genotype data. We found that haplotype data increased the accuracy of mixture analysis for

91    closely related populations by up to seven percentage points and individual assignment up to 17

92    percentage points relative to single-SNP genotype data for a set of 500 loci. These results

93    suggest that haplotyping is an efficient method to improve accuracy of genetic stock

94    identification when multi-SNP data are available.

95    *Methods*

96    Raw sequence data from Larson *et al.* (2014b) were downloaded from Dryad

97    (doi:10.5061/dryad.rs4v1) and reprocessed using *STACKS* (Catchen et al. 2013). This dataset

98    consisted of 250 Chinook salmon from five populations from Western Alaska: Big Salmon

99    (n=47), Tubutulik (n=56), Anvik (n=47), Kogrukluk (n=44) and Koktuli (n=56). Default settings

100   were used for all *STACKS* modules with the following exceptions: *process_radtags* (-c −r −q −

101   filter_illumina −t 94), *ustacks* (-m 2 −M 2, --model_type bounded --bound_high 0.05), *cstacks* (-

102   n 2). Loci (individual RADtags) were retained if the genotype rate was at least 80% and if the

103   minor allele frequency was at least 0.05 in one or more populations following the methods of

104   Larson *et al.* (2014b). The *STACKS* catalog for the Chinook salmon dataset was based on the

105   catalog from McKinney *et al*. (2016) to ensure locus names were consistent among studies.

106   Two datasets were generated with the loci retained from *STACKS*; a single-SNP dataset and a

107   multi-SNP (haplotype) dataset. The single-SNP dataset was constructed by retaining only a

108   single SNP from loci that had multiple variant sites. For these loci, the SNP with the highest

109   minor allele frequency across all populations was retained. For the haplotype dataset, all variant

110   sites at a locus were retained, and each unique haplotype was considered an allele.

111   Samples from each population were split evenly into training and holdout sets for both the

112   single-SNP and haplotype datasets; if populations could not be split exactly in half the remaining

113   individual was assigned to the training dataset. Locus selection was conducted using only the

114   single-SNP training set; both the single-SNP and haplotype holdout sets were used for mixture

115    analysis.  Separate training and holdout sets prevent an upward bias in assignment accuracy

116    (Anderson 2010).

117    $F_{ST}$ was estimated for each locus in the single-SNP training set using Genepop (Rousset 2008).

118    Loci were ranked by $F_{ST}$ and the 500 loci with highest $F_{ST}$ were chosen for mixture analysis.

119    Loci were chosen using the single-SNP dataset for two reasons: 1) to reflect the method that has

120    been historically used to choose loci for GSI studies, and 2) to evaluate the increase in power of

121    haplotypes for GSI even when haplotype information was not considered when initially choosing

122    loci.  The accuracy of single-SNP genotype vs haplotype data for mixture analysis was then

123    assessed using 500 loci.  Preliminary testing (not shown) suggested that multi-SNP haplotypes

124    increased mixture analysis accuracy up to five percentage points over single-SNP genotypes

125    using the 500 highest $F_{ST}$ loci, approximately 1/3 of which contained multiple SNPs.  We

126    speculated that a panel where all loci contained multiple SNPs would further increase the

127    accuracy of haplotyping relative to genotyping.  The panel used in this study was obtained by

128    filtering the full locus set to retain only loci with three or more haplotype alleles and taking the

129    500 highest $F_{ST}$ ranked multiple SNP loci.  Even though the $F_{ST}$ was slightly reduced in the final

130    panel relative to the preliminary panel, there was no difference in the mixture analysis results for

131    single-SNP genotypes.

132    Mixture analysis was conducted with both the single-SNP genotypes and haplotypes to estimate

133    assignment accuracy using the selected panel and procedures of Anderson (2010).  Populations

134    were first aggregated into reporting groups following the population structure previously

135    identified by Larson *et al*.  (2014b).  The Big Salmon, Tubutulik River, and Anvik rivers were

136    placed in their own reporting groups while the Kogrukluk and Koktuli rivers were placed into a

137    combined reporting group (Kogrukluk/Koktuli).

138    Mixture analysis was run in GSIsim (Anderson et al. 2008) which implements the training

139    holdout leave-one-out (THL) simulation method of Anderson (2010); simulation allows the

140    accuracy of the panel to be estimated for a range of mixture proportions and allows an estimate

141    of the variance in panel accuracy.  Rather than directly assigning individuals in the holdout set to

142    the baseline, GSIsim simulates genotypes for each individual in a mixture by drawing from

143    population allele frequencies in the holdout dataset.  Mixture proportions are estimated by

144    comparing the simulated mixture individuals against the baseline of all individuals. The leave-

145    one-out method is implemented at this point; an individual's alleles are excluded from the

146    baseline when calculating the probability of the individual's genotype for their population of

147    origin.  Options in GSIsim were set as follows: leave-one-out yes, 1000 mixedFisherySamples,

148    400 mixedFisheryIndividuals.

149    GSIsim was run iteratively with one reporting group at a time as the focal group and the other

150    reporting groups as the secondary groups.  The true proportion of the focal reporting group in the

151    mixture was simulated from 0 to 100% in 5% intervals.  The remaining mixture proportion was

152    divided into 5% increments and split equally among the secondary reporting groups.  Any

153    remaining proportion was randomly assigned to one of the secondary reporting groups.  When

154    multiple populations were contained within a reporting group, the mixture proportion for that

155    reporting group was evenly divided among the populations.  Reporting groups were simulated

156    with a range of mixture proportions to determine accuracy through a range of stock

157    compositions.  Simulation results are reported as the mean mixture estimate ± the threshold for

158    p=0.025 and p=0.925.  This encompasses the 95% distribution of estimates.

159

160    Management agencies commonly evaluate reporting groups by using 100% simulations where

161    90% allocation is the minimum acceptable threshold (Seeb et al. 2000; Seeb et al. 2007;

162    Beacham et al. 2012).  Further panels were developed by removing loci from the 500 locus panel

163    to determine how few loci could be used and still obtain 90% allocation at 100% simulation for

164    mixture analysis with single-SNP genotype or haplotype data.  Loci were removed based on $F_{ST}$

165    rank, with the lowest $F_{ST}$ loci removed first.  Additional panels were tested with 400, 300, 200,

166    150, and 100 loci.

167    The output of GSIsim includes the assignment probability of each simulated individual to each

168    population.  To assign an individual to a reporting group, the probability of assignment to each

169    population within the reporting group was summed.  A threshold of 90% cumulative probability

170    was required to assign an individual to a reporting group.

171

172    *Results*

173    A total of 14,494 loci passed quality filters.  The majority of the loci (72%) contained a single

174    SNP while the second most common category was 2 SNPs (24%) (Table 1).  Loci with 2 SNPs

175    had between 2 and 4 alleles: three alleles were most common (84%) followed by two alleles

176    (14%) and four alleles (2%).  Substantially more alleles with frequencies <0.10 were found for

177    haplotypes rather than single-SNP genotypes (8,959 vs 5,647).

178    Accuracy of mixture analysis in 100% simulations varied by reporting group but was

179    approximately 90% or greater for all cases.  The Big Salmon River and Tubutulik River

180    reporting groups each had 100% assignment with both the single-SNP genotype and the

181    haplotype datasets (Table 2).  This was expected given the high degree of differentiation between

182    these reporting groups and the others (Larson et al. 2014b).

183    For the Anvik River and Kogrukluk/Koktuli reporting groups, the accuracy of mixture analysis

184    was influenced both by the number of markers and whether haplotypes or single-SNP genotypes

185    were used.  With 500 markers the Anvik River reporting group had a maximum of ~90%

186    accuracy at 100% simulation for the single-SNP genotype data, but increased to ~97% accuracy

187    with the haplotypes (Table 2, Fig. 1).  The Kogrukluk/Koktuli reporting group had a maximum

188    accuracy of ~94% at 100% simulation with the single-SNP genotype data and increased to ~98%

189    accuracy with the haplotypes (Table 2, Fig. 2).  Precision of mixture analysis also increased, with

190    the 95% distribution of mixture estimates decreasing by approximately half for haplotype vs

191    single-SNP analyses.  Decreasing the number of loci for mixture analysis to only 100 loci

192    resulted in assignment accuracy as low as ~77% for single-SNP genotypes and ~86% for

193    haplotypes for 100% simulations (Table 3).  The panel of 150 loci was the smallest number of

194    loci that achieved 90% accuracy with haplotypes while at least 300 loci were needed to achieve

195    90% accuracy for the single-SNP genotype data (Table 3).

196    Individual assignment to reporting groups varied considerably across reporting groups (Table 4).

197    The Anvik reporting group had the lowest assignment rate at 75% followed by the

198    Kogrukluk/Koktuli reporting group (81%) with the 500 locus genotype data.  With the 500 locus

199    haplotype data both of these reporting groups had 92% accuracy.  The Big Salmon and Tubutulik

200    reporting groups had 100% assignment for both datasets.  When fewer than 500 markers were

201    used, neither the single-SNP genotype, or haplotype datasets were able to achieve 90% accuracy

202    for individual assignment to the Anvik or Kogrukluk/Koktuli reporting groups.

203    *Discussion*

204    Haplotype data showed greater accuracy and precision for mixture analysis and individual

205    assignment than did single-SNP genotyping for the same set of loci where comparisons were

206    possible.  This increase in accuracy was obtained through a simple change in analysis of the

207    same sequence data.  No comparisons were possible for the Big Salmon River and Tubutulik

208    River reporting groups as they showed 100% accuracy in mixture analysis as well as individual

209    assignment with both the single-SNP genotype and haplotype data.  The high differentiation in

210    these populations allowed near perfect accuracy, even when tested with only 100 loci (data not

211    shown).

212    The power of haplotyping was demonstrated in discriminating populations that were less

213    differentiated, in this case the Anvik and Kogrukluk/Koktuli rivers.  The Anvik and

214    Kogrukluk/Koktuli rivers showed increased accuracy and precision with the haplotype data

215    relative to the single-SNP genotype data for both mixture analysis and individual assignment.

216    Importantly, for individual assignment, the proportion of unassigned individuals decreased by

217    one-half to two-thirds when haplotypes were used rather than a single SNP per locus (Table 4).

218    This improvement is consistent with previous simulation and microsatellite studies, where loci

219    with more alleles had greater assignment accuracy than loci with fewer alleles (Kalinowski 2004;

220    Beacham et al. 2012).

221    The increased information in haplotype data can improve GSI analyses either by allowing greater

222    resolution for reporting groups (finer-scale divisions) or by allowing for smaller locus sets in

223    management applications.  We were not able to further divide the reporting groups in our dataset

224    but haplotypes were able to achieve the same accuracy as single-SNP genotypes with fewer loci.

225    A 90% threshold is often used as the standard for accuracy at 100% simulation for mixture

226    analysis (Seeb et al. 2000).  Using this threshold, only 150 loci would be necessary to achieve

227    sufficient accuracy with the haplotype data while at least 300 loci would be necessary with the

228    single-SNP genotype data. For an amplicon sequencing panel, this would translate into a

229    reduction in sequencing costs by a factor of two.

230    While we were not able to test accuracy of parentage assignment with our dataset, haplotypes are

231    likely to have an even greater positive impact. This is particularly true for distinguishing half-

232    sibs, aunts/uncles, or other types of complicated relationships that are likely to be more frequent

233    in small populations or species with promiscuous/polygamous mating systems. As parentage

234    assignment relies on comparing genotypes of individuals to the genotypes of parents, the

235    addition of more alleles per locus through haplotype analysis should reduce uncertainty in

236    parentage assignment.

237    Haplotyping clearly outperformed genotyping for mixture analysis and individual assignment;

238    however, there are considerations for implementation in management. RADseq data are often

239    used in the preliminary stages of mixture analysis to survey loci for informative SNPs, but

240    haplotyping using RADseq alone is not suitable for continued monitoring because of cost and

241    work flow. On the other hand, amplicon sequencing is likely to gain popularity for stock

242    identification studies because of the increased flexibility in choosing the number of SNPs as well

243    substantial cost efficiencies (Campbell et al. 2015). For amplicon sequencing there is a trade-off

244    between the number of loci and number of individuals that can be genotyped for a given batch of

245    sequencing; decreasing the number of loci would allow more individuals to be genotyped for the

246    same sequencing cost. Alternatively, additional loci could be added that discriminate other

247    populations, allowing expanded utility of a panel.

248    *The evolution of haplotypes*

249   Multiple-SNP DNA templates generally arise from sequential occurrence of single nucleotide

250   substitutions.  The first substitution in a monomorphic template results in two SNP alleles. When

251   a new substitution occurs proximal to the existing SNP, the new substitution is associated with

252   only one of the alleles from the original SNP, resulting in three possible alleles (Fig. 3).  DNA

253   templates with two polymorphic sites can have up to four alleles (a vast majority of SNPs are

254   binary); the fourth possible allele can only arise from recombination between the original and

255   new SNP or, very rarely, through recurrent substitution (Hudson and Kaplan 1985).

256   This model of haplotype evolution is supported by observations from our data.  Loci with two

257   SNPs most commonly had three alleles (84%) with two alleles being the second most common

258   situation (14%) (Table 1).  Loci with two SNPs and two alleles occur when the allele on which

259   the second substitution occurred is lost in the population (ie. TC haplotype in Fig. 3).  The

260   relatively high proportion of loci with two SNPs and two alleles suggests that loss of alleles is

261   relatively common.  Two SNP loci with four haplotype alleles were relatively uncommon (2%)

262   suggesting recombination occurring between two SNPs in a locus has been exceedingly rare in

263   these short templates.

264   *The added power of haplotypes*

265   Multiple SNPs within a locus are expected to have arisen at different times and the geographic

266   distribution of the resulting haplotype alleles will reflect gene flow among populations following

267   the genesis of each allele (Mathieson and McVean 2014).  Alleles resulting from recent SNPs

268   should manifest in haplotype data as rare allelic variants.  Rare allelic variants tend to be more

269   geographically restricted than common variants and have been demonstrated to increase power to

270   detect fine-scale population structure (O'Connor et al. 2015).  Individuals with the rare allele

271   may be assigned to population of origin with some confidence, although most individuals will be

272    homozygous and uninformative.  The haplotyping of three or more alleles enables the increased

273    resolution of rare alleles while still providing information from the more common alleles.

274    We recognize that choosing loci for GSI based on haplotype information rather than ranking loci

275    based on only one of the SNPs present will likely lead to further increases in GSI accuracy,

276    particularly if locus selection and GSI methods can be developed that take into account the

277    evolutionary relationship between haplotypes.  In this manuscript we focused on demonstrating

278    that haplotypes provide a benefit over using a single SNP per locus; this is particularly relevant

279    as amplicon sequencing applications for stock identification are becoming increasingly common

280    (e.g., many agencies that manage salmonid fisheries are adopting SNP panels similar to the

281    methods of Campbell et al. 2015).  Sequencing data from these panels can be examined to

282    determine if loci contain multiple SNPs; if so, haplotype analysis can be conducted to evaluate

283    relative improvements in accuracy and precision. The fact that we saw improved accuracy in

284    both mixture analysis and individual assignment, even when markers were chosen based on

285    single-SNP $F_{ST}$, supports the argument that existing panels can be improved by changing to

286    haplotype analysis.

287    Analysis of haplotypes rather than single-SNP genotypes for RADseq data increased accuracy

288    for mixture analysis in Alaskan Chinook salmon.  Haplotype data are readily available from

289    RADseq analyses, allowing existing datasets to be re-examined for informative loci.  Haplotype

290    data are also readily generated from amplicon sequencing methods, allowing informative loci to

291    be incorporated into population monitoring.  We believe that haplotyping multiple SNPs shows

292    promise to improve the accuracy of assigning unknown fish to population or group of origin or

293    parentage analysis whenever haplotype data are available.  The most effective panels will likely

294    include a combination of the highest $F_{ST}$ single SNPs along with multi-SNP haplotypes.

303 Literature Cited:

304 Abadia-Cardoso, A., Pearse, D.E., Jacobson, S., Marshall, J., Dalrymple, D., Kawasaki, F., Ruiz-
305 Campos, G., and Garza, J.C. 2016. Population genetic structure and ancestry of
306 steelhead/rainbow trout (*Oncorhynchus mykiss*) at the extreme southern edge of their range in
307 North America. Conservation Genetics **17**(3): 675-689.

308 Ackerman, M.W., Habicht, C., and Seeb, L.W. 2011. Single-nucleotide polymorphisms (SNPs)
309 under diversifying selection provide increased accuracy and precision in mixed-stock analyses of
310 sockeye salmon from the Copper River, Alaska. Transactions of the American Fisheries Society
311 **140**(3): 865-881.

312 ADF&G. 2013. Chinook salmon stock assessment and research plan Special Publication No. 13-
313 01. Alaska Department of Fish and Game, Special Publication No. 13-01. Available from
314 http://www.adfg.alaska.gov/static/home/news/hottopics/pdfs/chinook_research_plan.pdf.

315 Ali, O.A., O'Rourke, S.M., Amish, S.J., Meek, M.H., Luikart, G., Jeffres, C., and Miller, M.R.
316 2016. RAD capture (Rapture): flexible and efficient sequence-based genotyping. Genetics
317 **202**(2): 389-400.

318 Anderson, E.C. 2010. Assessing the power of informative subsets of loci for population
319 assignment: standard methods are upwardly biased. Molecular Ecology Resources **10**(4): 701-
320 710.

321 Anderson, E.C., Waples, R.S., and Kalinowski, S.T. 2008. An improved method for predicting
322 the accuracy of genetic stock identification. Canadian Journal of Fisheries and Aquatic Sciences
323 **65**(7): 1475-1486.

324 Andrews, K.R., Good, J.M., Miller, M.R., Luikart, G., and Hohenlohe, P.A. 2016. Harnessing
325 the power of RADseq for ecological and evolutionary genomics. Nature Reviews Genetics **17**(2):
326 81-92.

327 Araneda, C., Larrain, M.A., Hecht, B., and Narum, S. 2016. Adaptive genetic variation
328 distinguishes Chilean blue mussels (*Mytilus chilensis*) from different marine environments. Ecol
329 Evol.

330 Baird, N.A., Etter, P.D., Atwood, T.S., Currey, M.C., Shiver, A.L., Lewis, Z.A., Selker, E.U.,
331 Cresko, W.A., and Johnson, E.A. 2008. Rapid SNP discovery and genetic mapping using
332 sequenced RAD markers. PloS ONE **3**(10): e3376.

333    Beacham, T.D., Jonsen, K., and Wallace, C. 2012. A comparison of stock and individual

334    identification for Chinook salmon in British Columbia provided by microsatellites and single-

335    nucleotide polymorphisms. Marine and Coastal Fisheries **4**(1): 1-22.

336    Beacham, T.D., McIntosh, B., and Wallace, C. 2010. A comparison of stock and individual

337    identification for sockeye salmon (*Oncorhynchus nerka*) in British Columbia provided by

338    microsatellites and single nucleotide polymorphisms. Canadian Journal of Fisheries and Aquatic

339    Sciences **67**(8): 1274-1290.

340    Beacham, T.D., Winter, I., Jonsen, K.L., Wetklo, M., Deng, L.T., and Candy, J.R. 2008. The

341    application of rapid microsatellite-based stock identification to management of a Chinook

342    salmon troll fishery off the Queen Charlotte Islands, British Columbia. North American Journal

343    of Fisheries Management **28**(3): 849-855.

344    Bradbury, I.R., Hamilton, L.C., Chaput, G., Robertson, M.J., Goraguer, H., Walsh, A., Morris,

345    V., Reddin, D., Dempson, J.B., Sheehan, T.F., King, T., and Bernatchez, L. 2016. Genetic mixed

346    stock analysis of an interceptory Atlantic salmon fishery in the Northwest Atlantic. Fisheries

347    Research **174**: 234-244.

348    Campbell, N.R., Harmon, S.A., and Narum, S.R. 2015. Genotyping-in-Thousands by sequencing

349    (GT-seq): A cost effective SNP genotyping method based on custom amplicon sequencing.

350    Molecular Ecology Resources **15**(4): 855-867.

351    Catchen, J., Hohenlohe, P.A., Bassham, S., Amores, A., and Cresko, W.A. 2013. Stacks: an

352    analysis tool set for population genomics. Mol. Ecol. **22**(11): 3124-3140.

353    Dann, T.H., Habicht, C., Baker, T.T., and Seeb, J.E. 2013. Exploiting genetic diversity to

354    balance conservation and harvest of migratory salmon. Canadian Journal of Fisheries and

355    Aquatic Sciences **70**(5): 785-793.

356    Ensing, D., Crozier, W.W., Boylan, P., O'Maoileidigh, N., and McGinnity, P. 2013. An analysis

357    of genetic stock identification on a small geographical scale using microsatellite markers, and its

358    application in the management of a mixed-stock fishery for Atlantic salmon *Salmo salar* in

359    Ireland. J Fish Biol **82**(6): 2080-2094.

360    Evans, M.L., Johnson, M.A., Jacobson, D., Wang, J.L., Hogansen, M., and O'Malley, K.G. 2016.

361    Evaluating a multi-generational reintroduction program for threatened salmon using genetic

362    parentage analysis. Canadian Journal of Fisheries and Aquatic Sciences **73**(5): 844-852.

363    Fariello, M.I., Boitard, S., Naya, H., San Cristobal, M., and Servin, B. 2013. Detecting signatures

364    of selection through haplotype differentiation among hierarchically structured populations.

365    Genetics **193**(3): 929-941.

366    Ford, M., Pearsons, T.N., and Murdoch, A. 2015. The spawning success of early maturing

367    resident hatchery Chinook salmon in a natural river system. Transactions of the American

368    Fisheries Society **144**(3): 539-548.

369    Gattepaille, L.M., and Jakobsson, M. 2012. Combining markers into haplotypes can improve

370    population structure inference. Genetics **190**(1): 159-174.

371    Gisclair, B.R. 2009. Salmon bycatch management in the Bering Sea walleye pollock fishery:

372    threats and opportunities for Western Alaska. *In* Pacific Salmon:  Ecology and Management of

373    Western Alaska's Populations. *Edited by* C.C. Krueger and C.E. Zimmerman. American

374    Fisheries Society Symposium 70, Bethesda, Maryland. pp. 799-816.

375    Gustafson, R.G., Waples, R.S., Myers, J.M., Weitkamp, L.A., Bryant, G.J., Johnson, O.W., and

376    Hard, J.J. 2007. Pacific salmon extinctions: Quantifying lost and remaining diversity.

377    Conservation Biology **21**(4): 1009-1020.

378    Habicht, C., Seeb, L.W., Myers, K.W., Farley, E.V., and Seeb, J.E. 2010. Summer-fall

379    distribution of stocks of immature sockeye salmon in the Bering Sea as revealed by single-

380    nucleotide polymorphisms (SNPs). Transactions of the American Fisheries Society **139**: 1171-

381    1191.

382    Hohenlohe, P.A., Amish, S.J., Catchen, J.M., Allendorf, F.W., and Luikart, G. 2011. Next-

383    generation RAD sequencing identifies thousands of SNPs for assessing hybridization between

384    rainbow and westslope cutthroat trout. Molecular Ecology Resources **11 Suppl 1**: 117-122.

385    Hudson, R.R., and Kaplan, N.L. 1985. Statistical properties of the number of recombination

386    events in the history of a sample of DNA sequences. Genetics **111**: 147-164.

387    Kalinowski, S.T. 2004. Genetic polymorphism and mixed-stock fisheries analysis. Canadian

388    Journal of Fisheries and Aquatic Sciences **61**(7): 1075-1082.

389    Larson, W.A., Seeb, J.E., Pascal, C.E., Templin, W.D., and Seeb, L.W. 2014a. Single-nucleotide

390    polymorphisms (SNPs) identified through genotyping-by-sequencing improve genetic stock

391    identification of Chinook salmon (*Oncorhynchus tshawytscha*) from western Alaska. Canadian

392    Journal of Fisheries and Aquatic Sciences **71**(5): 698-708.

393    Larson, W.A., Seeb, L.W., Everett, M.V., Waples, R.K., Templin, W.D., and Seeb, J.E. 2014b.

394    Genotyping by sequencing resolves shallow population structure to inform conservation of

395    Chinook salmon (*Oncorhynchus tshawytscha*). Evolutionary Applications **7**(3): 355-369.

396    Larson, W.A., Utter, F.M., Myers, K.W., Templin, W.D., Seeb, J.E., Guthrie III, C.M., Bugaev,

397    A.V., Seeb, L.W., and Moran, P. 2013. Single-nucleotide polymorphisms reveal distribution and

398    migration of Chinook salmon (*Oncorhynchus tshawytscha*) in the Bering Sea and North Pacific

399    Ocean. Canadian Journal of Fisheries and Aquatic Sciences **70**(1): 128-141.

400    Leslie, S., Winney, B., Hellenthal, G., Davison, D., Boumertit, A., Day, T., Hutnik, K., Royrvik,

401    E.C., Cunliffe, B., Wellcome Trust Case Control Consortium, International Multiple Sclerosis

402    Genetics Consortium, Lawson, D.J., Falush, D., Freeman, C., Pirinen, M., Myers, S., Robinson,

403    M., Donnelly, P., and Bodmer, W. 2015. The fine-scale genetic structure of the British

404    population. Nature **519**(7543): 309-314.

405    Manel, S., Gaggiotti, O.E., and Waples, R.S. 2005. Assignment methods: matching biological

406    questions techniques with appropriate techniques. Trends in Ecology & Evolution **20**(3): 136-

407    142.

408    Mathieson, I., and McVean, G. 2014. Demography and the age of rare variants. PLoS genetics

409    **10**(8): e1004528.

410    McKinney, G.J., Seeb, L.W., Larson, W.A., Gomez-Uchida, D., Limborg, M.T., Brieuc, M.S.,

411    Everett, M.V., Naish, K.A., Waples, R.K., and Seeb, J.E. 2016. An integrated linkage map

412    reveals candidate genes underlying adaptive variation in Chinook salmon (*Oncorhynchus

413    tshawytscha*). Molecular Ecology Resources **16**(3): 769-783.

414    O'Connor, T.D., Fu, W., Project, N.G.E.S., ESP Population Genetics and Statistical Analysis

415    Working Group, Turner, E., Mychaleckyj, J.C., Logsdon, B., Auer, P., Carlson, C.S., Leal, S.M.,

416    Smith, J.D., Rieder, M.J., Bamshad, M.J., Nickerson, D.A., and Akey, J.M. 2015. Rare variation

417    facilitates inferences of fine-scale population structure in humans. Mol Biol Evol **32**(3): 653-660.

418    O'Malley, K.G., Camara, M.D., and Banks, M.A. 2007. Candidate loci reveal genetic

419    differentiation between temporally divergent migratory runs of Chinook salmon (*Oncorhynchus

420    tshawytscha*). Molecular Ecology **16**(23): 4930-4941.

421    Rousset, F. 2008. Genepop'007: A complete re-implementation of the Genepop software for

422    Windows and Linux. Molecular Ecology Resources **8**(1): 103-106.

423  Seeb, J.E., Carvalho, G., Hauser, L., Naish, K., Roberts, S., and Seeb, L.W. 2011. Single-

424  nucleotide polymorphism (SNP) discovery and applications of SNP genotyping in nonmodel

425  organisms. Molecular Ecology Resources **11**: 1-8.

426  Seeb, L.W., Antonovich, A., Banks, M., Beacham, T., Bellinger, R., Blankenship, S., Campbell,

427  M., Decovich, N., Garza, J.C., Guthrie, C., Lundrigan, T., Moran, P., Narum, S., Stephenson, J.,

428  Supernault, J., Teel, D., Templin, W.D., Wenburg, J.K., Young, S., and Smith, C.T. 2007.

429  Development of a standardized DNA database for Chinook salmon. Fisheries **32**(11): 540-552.

430  Seeb, L.W., Habicht, C., Templin, W.D., Tarbox, K.E., Davis, R.Z., Brannian, L.K., and Seeb,

431  J.E. 2000. Genetic diversity of sockeye salmon of Cook Inlet, Alaska, and its application to

432  management of populations affected by the *Exxon Valdez* oil spill. Transactions of the American

433  Fisheries Society **129**: 1223-1249.

434  Steele, C.A., Anderson, E.C., Ackerman, M.W., Hess, M.A., Campbell, N.R., Narum, S.R.,

435  Campbell, M.R., and Grant, J. 2013. A validation of parentage-based tagging using hatchery

436  steelhead in the Snake River basin. Canadian Journal of Fisheries and Aquatic Sciences **70**(7):

437  1046-1054.

438  Templin, W.D., Seeb, J.E., Jasper, J.R., Barclay, A.W., and Seeb, L.W. 2011. Genetic

439  differentiation of Alaska Chinook salmon: the missing link for migratory studies. Molecular

440  Ecology Resources **11 Suppl 1**: 226-246.

441  Wolfe, R.J., and Spaeder, J. 2009. People and salmon of the Yukon and Kuskokwim Drainages

442  and Noron Sound in Alaska:  fishery harvests, culture change, and local knowledge systems. *In*

443  Pacific Salmon Ecology and Management of Western Alaska's Populations. *Edited by* C.C.

444  Krueger and C.E. Zimmerman. American Fisheries Society, Bethesda, Maryland. pp. 349-379.

445

446

447    Table 1.  Distribution of number of SNPs and number of alleles per locus for all loci.

448

| Number of SNPs | Number of Loci | Number of Alleles | | |
|---|---|---|---|---|
| | | **2** | **3** | **4** |
| 1 | 10,448 (72%) | 10,448 (100%) | | |
| 2 | 3,441 (24%) | 476 (14%) | 2,878 (84%) | 87 (2%) |
| 3 | 556 (4%) | | 205 (37%) | 351 (63%) |
| 4 | 46 (<<1%) | | 4 (9%) | 42 (91%) |
| 5 | 3 (<<1%) | | | 3 (100%) |

449

450

451

452     Table 2.  Estimated mixture proportion at 100% simulation for each reporting group using the

453     500 highest $F_{ST}$ ranked multiple SNP loci.  Results are reported as mean ± threshold for the 95%

454     distribution of mixture estimates.

| Reporting Group | Single-SNP Genotype | Haplotype |
|---|---|---|
| Big Salmon River | 100%±0 | 100%±0 |
| Tubutulik River | 100%±0 | 100%±0 |
| Anvik River | 89.6%±3.6 | 97.4%±1.8 |
| Kogrukluk/Koktuli Rivers | 93.7%±3.1 | 97.8%±1.7 |

455

456    Table 3. Estimated mixture proportion at 100% simulation for each reporting group reported as

457    mean ± threshold for the 95% distribution of mixture estimates.  The number of loci in each

458    panel were decreased in intervals of 100 to test the accuracy of different panel sizes.  The 150

459    locus panel is the exception and was the smallest panel to achieve 90% assignment with the

460    haplotype data.  Results are not shown for the Big Salmon River and Tubutulik River reporting

461    groups; mixture estimates were ~100% for Big Salmon River and Tubutulik River reporting

462    groups for all panels of loci.

| Reporting Group | Number of Loci | Single-SNP Genotype | Haplotype |
|---|---|---|---|
| Anvik River | 500 | 89.6%±3.6 | 97.4%±1.8 |
|  | 400 | 90.7%±3.5 | 97.2%±1.9 |
|  | 300 | 87.3%±4.3 | 94.3%±2.7 |
|  | 200 | 87.1%±4.4 | 93.6%±3.1 |
|  | 150 | 80.9%±5.5 | 92.7%±3.3 |
|  | 100 | 61.7%±7.6 | 77.7%±5.2 |
| Kogrukluk/Koktuli Rivers | 500 | 93.7%±3.1 | 97.8%±1.7 |
|  | 400 | 92.4%±3.2 | 97.1%±2.0 |
|  | 300 | 89.4%±3.9 | 95.4%±2.5 |
|  | 200 | 83.2%±5.0 | 91.2%±3.6 |
|  | 150 | 80.1%±5.5 | 90.1%±3.8 |
|  | 100 | 77.0%±6.8 | 86.1%±4.8 |

463

464     Table 4.  Individual assignment rates to reporting groups for simulated individuals with A) 500
465     locus genotype dataset and B) 500 locus haplotype dataset.  True reporting groups are rows while
466     assigned reporting groups are columns.

467     A. Genotype

|  | Big Salmon | Tubutulik | Anvik | Kogrukluk/Koktuli | Unassigned |
|---|---|---|---|---|---|
| Big Salmon | 100% | 0% | 0% | 0% | 0% |
| Tubutulik | 0% | 100% | 0% | 0% | 0% |
| Anvik | 0% | 0% | 75% | 8% | 17% |
| Kogrukluk/Koktuli | 0% | 0% | 5% | 81% | 14% |

468

469     B. Haplotype

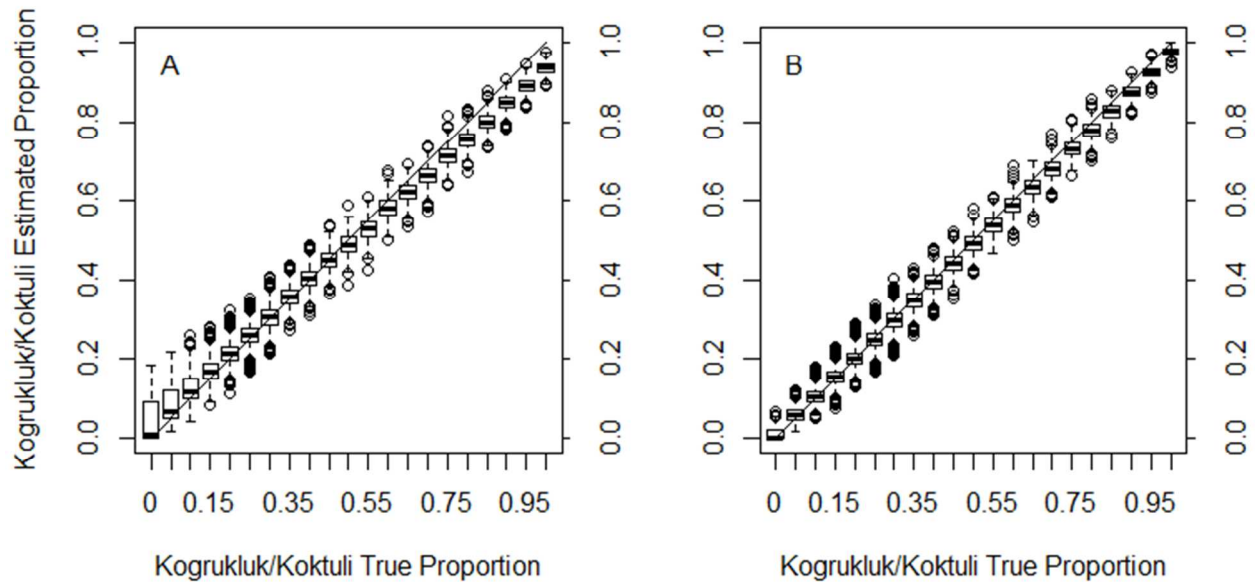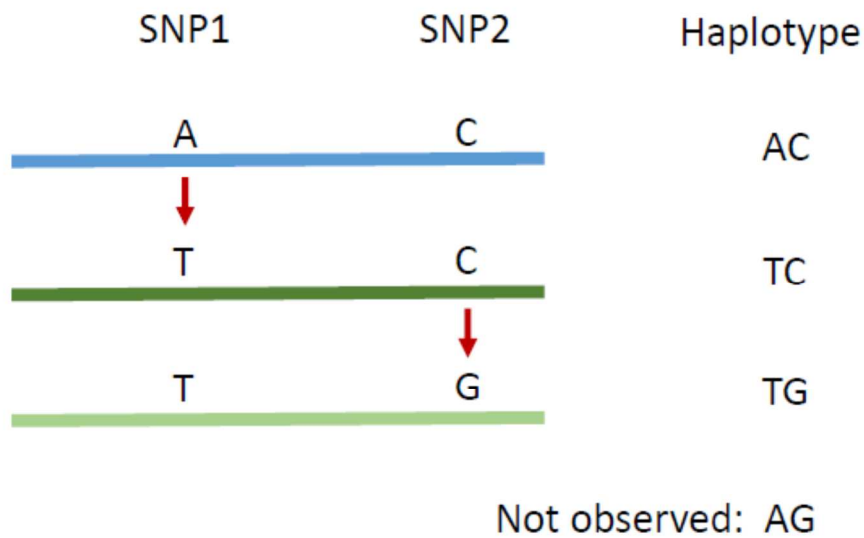|  | Big Salmon | Tubutulik | Anvik | Kogrukluk/Koktuli | Unassigned |
|---|---|---|---|---|---|
| Big Salmon | 100% | 0% | 0% | 0% | 0% |
| Tubutulik | 0% | 100% | 0% | 0% | 0% |
| Anvik | 0% | 0% | 92% | 3% | 5% |
| Kogrukluk/Koktuli | 0% | 0% | 3% | 92% | 6% |

470

471

472



473

474     Figure 1. Boxplots showing estimated mixture proportion for: A) the Anvik River using single-

475     SNP genotypes for the 500 highest $F_{ST}$ loci which have 3 or more haplotype alleles.  B) Anvik

476     River using haplotypes for the same loci as A.  The diagonal line shows the expected relationship

477     for perfect assignment.  The true proportion simulated in the mixture is given on the x-axis while

478     the boxplots show the estimated proportions on the y-axis.

479

480

Figure 2. Boxplots showing estimated mixture proportion for A) the Kogrukluk/Koktuli

reporting group using single-SNP genotypes for the 500 highest $F_{ST}$ loci which have 3 or more

alleles.  B) Kogrukluk/Koktuli reporting group using haplotypes for the same loci as A.  The

diagonal line shows the expected relationship for perfect assignment.  The true proportion

simulated in the mixture is given on the x-axis while the boxplots show the estimated proportions

on the y-axis.

487

488      Figure 3. Haplotypes resulting from two hypothetical SNPs in the absence of recombination.

489      The ancestral DNA sequence is in blue; the first SNP is an A/T variant that gives rise to the dark

490      green DNA sequence. The second SNP (C/G variant) will arise on one of these two alleles, in

491      this case the T allele. The two SNPs result in three haplotypes: AC, TC, and TG (light green

492      DNA sequence). Assuming the probability of a second mutation at the same site is negligible,

493      the fourth possible haplotype (AG, unobserved) can only arise through recombination between

494      the first and second SNP.