1

2  **Title**:  Regional Allocation Model                              **Version:**    1.0

3  **Authors:**  J. Jasper, C. Habicht, A. R. Munro, and W. Templin

4  **Date:**  April 30, 2010

5

6  # Introduction

7

8   Mixed stock analysis methods for estimating stock (population) compositions in fisheries have evolved

9   over time from conditional maximum-likelihood (Fournier et al. 1984) to Bayesian (Pella and Masuda,

10  2001) approaches. The Pella-Masuda model (a Bayesian approach; Pella and Masuda, 2001) has been the

11  "gold standard" since 2001. In these methods, however, bias is inevitable because the estimation of the

12  stock proportions is constrained to be non-negative and sum to one, meaning that rare or absent stocks in

13  the mixture are overestimated while common stocks are under estimated (Pella and Milner, 1987).  Stocks

14  are usually grouped into regional stock groupings (regions) for reporting.

15  Recent observations in our laboratory indicate that disproportionate numbers of stocks within a region can

16  lead to significant bias in regional composition estimates when regional stock structure is shallow. We

17  have observed that regions represented by large numbers of stocks seem to acquire higher misallocations

18  than regions represented by fewer stocks (Figure 1). This bias can be reduced at the regional level by

19  grouping stocks with similar genetic attributes into regions, then summing estimated proportions across

20  stocks within the regions (Wood et al. 1987). Here we present a rationale for why we think the observed

21  non-uniform bias occurred and a method that appears to improve allocation at the regional level as well as

22  distribute the misallocation more evenly among regions.

23  In the Pella-Masuda model, the data augmentation algorithm is used to generate from the posterior

24  distribution the stock identities of each of the mixture individuals, and then generate the stock proportions

25  and baseline allele frequencies based on summaries of these identities. At each cycle of the algorithm, the

26  stock identity of mixture individual $m$ is stochastically assigned to stock $i$ with probability proportional to

27  the product of stock $i$'s contribution to the mixture and the relative frequency of individual $m$'s genotype

28  in stock $i$. This means that individual $m$ has a finite probability of belonging to each and every stock in the

29  baseline. We will refer to these probabilities as the identity probabilities.

30  The chances that individual $m$ is assigned to the correct stock at a particular iteration is a function of not

31  only the genetic distinction of its stock, but also, theoretically, the number of stocks in the baseline.

32  Fortunately, fisheries managers often are not interested in the proportion of individual stocks, but rather in

---

33  the contribution made by all stocks within regions. If the stocks within a region are genetically more
34  similar to each other than to stocks in other regions (strong regional structure), then the chances of
35  correctly assigning an individual to a stock within the correct region each cycle greatly improves
36  estimation (Wood et al. 1987). However, with weak regional structure, the chances of assigning an
37  individual to a stock within the correct region may be significantly influenced by the number of stocks in
38  each region. This may be because the probability of assigning an individual to a particular region is the
39  sum of the identity probabilities across all the stocks in the region, such that adding stocks adds
40  probability. If the amount of misallocation to a region is a function of the number of stocks within that
41  region, an inherent non-uniform bias in regional contribution estimates can occur simply due to differing
42  numbers of stocks among regions.

43  The purpose of this paper is to illustrate that unequal numbers of stocks among regions leads to unequal
44  biases in misallocation and to determine if a new analytical method may mitigate this bias. We anticipate
45  an upward misallocation bias toward regions that are represented by larger numbers of stocks than regions
46  represented by fewer stocks using the Pella-Masuda model. We present a new analytical model that
47  appears to diminish this bias.

48

# Methods

49

50

51  We considered three methods to examine the assertion that unequal numbers of stocks within regions do
52  not affect bias in misallocation. We selected baseline data for chum salmon stocks from Western Alaska.
53  These data were chosen because these stocks represent weak regional structure (Figure 5).

54  The first two methods use the Pella-Masuda model but differ in how the priors are assigned. The first
55  method is the widely used True Flat Prior (TFP; Pella and Masuda, 2001). This model provides no *a*
56  *priori* information about the regional structure and gives an equal prior "count" of $1/C$ to each of the
57  stocks in the baseline, where $C$ is the number of stocks. This is the model that provided the recent
58  observations in our laboratory that suggested that disproportionate numbers of stocks within a reporting
59  group can affect the regional composition estimates.

60  The second method, termed the Regional Flat Prior (RFP), is a method currently in use at ADF&G's
61  Gene Conservation Laboratory (Dann et al. 2009). The structure of the prior for stock proportions is an *ad*
62  *hoc* alternative to the TFP. Under the RFP, for each of the stocks within the *g*th region, we give a prior
63  "count" equal to $1/G/C_g$, where $G$ is the number of regions and $C_g$ is the number of stocks within the *g*th
64  region. Therefore, equal prior "count" is given to each region, but the prior "count" given to a stock is
65  dependent upon the number of stocks within its region.

66  The third method, termed the Regional Allocation Model (RAM), is currently under development at the
67  Gene Conservation Laboratory. This model is very similar to the Pella-Masuda model in that it is based
68  on the data augmentation algorithm that alternates between generating of the parameters of the model.
69  The difference is that in the RAM, we first generate the regional identity of each individual, and then
70  produce regional contributions based on summaries of these regional identities. For individual *m*, the
71  regional identity probability of belonging to region *g* is proportional to region *g*'s contribution to the

72    mixture times a weighted average relative frequency of individual $m$'s genotype across all $C_g$ stocks

73    within the region. The weights are simply the within-region stock proportions, and they sum to one.

74    Because the weights do sum to one, the genetic component of the regional identity probabilities remain on

75    the same scale regardless of the number of stocks within the region, which should presumably moderate

76    the non-uniform bias due to the unequal distribution of stocks among the regions. There is actually a

77    second stage to the data augmentation algorithm in which, after an individual is assigned to a region, it is

78    then allocated to a stock within that region. This is done exactly as is done in the Pella-Masuda model

79    except that it is done with respect to a baseline that is reduced to only that region.

80    *General Bayesian Methods*

81    For estimating parameters $\boldsymbol{\theta}$ from data $X$ using Bayesian methods, we aim at the evaluation of the

82    posterior distribution $P(\boldsymbol{\theta}|X) = L(X|\boldsymbol{\theta}) \, P(\boldsymbol{\theta})/m(X)$, where $L(X|\boldsymbol{\theta})$ is the likelihood of the data given the

83    parameters, $P(\boldsymbol{\theta})$ is the prior distribution of the parameters, which must be specified, and $m(X)$ is the

84    constant marginal distribution of the data. From this distribution, summary statistics for $\boldsymbol{\theta}$ can be derived.

85    However, these distributions are rarely soluble in closed form for multidimensional parameter vector $\boldsymbol{\theta}$,

86    and we must rely on drawing samples from it via a Gibbs sampling routine, from which the summary

87    statistics can be calculated. For mixed stock analysis, $\boldsymbol{\theta}$ represents the stock proportions and the baseline

88    allele frequencies while $X$ corresponds to the mixture genotypes and the baseline allele counts. As

89    mentioned previously, a prior distribution must be specified for the parameters. In the forthcoming

90    models, the mathematically convenient Dirichlet distribution is used for the stock proportions as well as

91    the baseline allele frequencies. A Dirichlet distribution with parameter vector $\boldsymbol{\lambda}$ is a distribution on a

92    vector $W$ whose sum is constrained to one. It has the form:

$$P(W|\lambda) = \frac{\Gamma(\sum_{i=1}^{n} \lambda_i)}{\prod_{i=1}^{n} \Gamma(\lambda_i)} \prod_{i=1}^{n} W_i^{\lambda_i}$$

93

94

95    *The Pella-Masuda Model*

96    We denote the count of the $j$th ($j=1,2,\ldots,J_d$) allele of the $d$th ($d=1,2,\ldots,D$) locus for mixture individual $m$

97    as $x_{mdj}$, and let $X_m$ signify the entire multi-locus genotype for this individual. The array $X$ represents the

98    multi-locus genotypes for all $M$ mixture individuals. Similarly, we let $y_{idj}$ denote the count of the $j$th allele

99    for the $d$th locus of the $i$th baseline stock, and $Y$ denotes the entire baseline. This describes the data.

100    To describe the parameters, let the stock proportion for the $i$th stock be denoted as $P_i$, and let $P$ be the

101    vector of all stock proportions. We place a Dirichlet prior distribution on the stock proportions with prior

102    parameters $\boldsymbol{\alpha}$, where $\alpha_i$ is determined by our choice of prior structure discussed earlier (RFP or TFP).

103    We let $q_{idj}$ denote the relative frequency of the $j$th allele for the $d$th locus in the $i$th baseline stock and let

104    $Q$ denote the entire array of baseline relative frequencies. We place a Dirichlet prior distribution on $Q_{id}$

105    with prior parameters $\boldsymbol{\beta}_d$, where $\beta_{dj} = 1/J_d$, with $J_d$ being the number of alleles for locus $d$.

106    Finally, let $z_{mi}$ be the stock identity for the $m$th mixture individual in the $i$th stock, where $z_{mi}$ is equal to

107    one if individual $m$ belongs to the $i$th stock and zero otherwise. We denote $Z_m$ as the vector of stock

108 identities for individual $m$, and $\mathbf{Z}$ as the matrix of stock identities for the entire mixture. We place a
109 multinomial prior on $\mathbf{Z}_m$ with size 1 and probabilities equal to the stock proportions $\mathbf{P}$.

110 The genotypic likelihood of the $m$th individual would be greatly simplified if we knew the stock identity
111 of that individual. In other words, if $z_{mi} = 1$, then the likelihood of observing individual $m$ is simply the
112 relative frequency of this individual's multi-locus genotype in the $i$th stock, which we denote by $f(X_m|Q_i)$,
113 where:

$$f(X_m|Q_i) \propto \prod_{d=1}^{D} \prod_{j=1}^{J_d} q_{idj}^{x_{dj}}$$

114

115 Because $z_{mi'} = 0$ for all $i' \neq i$, the full genotypic likelihood may be expressed as:

$$L(X|Q,Z) = \prod_{m=1}^{M} \prod_{i=1}^{C} f(X_m|Q_i)^{z_{mi}}$$

116

117 In addition to the genotypic data, we need to consider the likelihood of the baseline data, which can be
118 written as:

$$L(Y|Q) \propto \prod_{i=1}^{C} \prod_{d=1}^{D} \prod_{j=1}^{J_d} q_{idj}^{y_{idj}}$$

119

120 The full likelihood, $L(X,Y|Q,Z)$, is simply the product of these two components.

121 Multiplying this likelihood by the prior distributions leads to the following posterior distribution:

$$P(P,Q,Z|X,Y) \propto L(X,Y|Q,Z)P(Z|P)P(P|\alpha)P(Q|\beta)$$

$$\propto \left( \prod_{m=1}^{M} \prod_{i=1}^{C} f(X_m|Q_i)^{z_{mi}} \right) \left( \prod_{i=1}^{C} \prod_{d=1}^{D} \prod_{j=1}^{J_d} q_{idj}^{y_{idj}} \right)$$

$$\times \left( \prod_{m=1}^{M} \prod_{i=1}^{C} P_i^{z_{mi}} \right) \left( \prod_{i=1}^{C} P_i^{\alpha_i} \right) \left( \prod_{i=1}^{C} \prod_{d=1}^{D} \prod_{j=1}^{J_d} q_{idj}^{\beta_{dj}} \right)$$

122

123 The benefit of using the chosen prior distributions is that the conditional posterior distribution for each of
124 the parameters given the data and the remaining parameters is of the same form as the prior distribution
125 (conjugacy). This property makes them easy to sample from within a Gibbs sampler, which proceeds as
126 follows: first, starting with initial values for $\mathbf{P}$ and $\mathbf{Q}$, we draw stock identities for each of the mixture
127 individuals from:

$$Z_m|P,Q,X_m \sim \text{multinomial}\left( 1, \left\{ \frac{P_i f(X_m|Q_i)}{\sum_{k=1}^{C} P_k f(X_m|Q_k)} \right\}_{i=1,2,...,C} \right)$$

128

129 Next, given these stock identities, $\mathbf{P}$ is drawn from:

$$P|Z,\alpha \sim \mathrm{Dirichlet}\left(\left\{\sum_{m=1}^{M} z_{mi} + \alpha_i\right\}_{i=1,2,\ldots,C}\right)$$

130

131     Finally, for each stock and for each locus, we generate $Q_{id}$ from:

$$Q_{id}|X,Y,Z,\beta \sim \mathrm{Dirichlet}\left(\left\{\sum_{m=1}^{M} z_{mi}x_{mdj} + y_{idj} + \beta_{dj}\right\}_{j=1,2,\ldots,J_d}\right)$$

132

133     This process is repeated for several thousand iterations, typically with multiple chains starting from
134     different initial values, and the first few thousand iterations are discarded as "burn-in" to remove the
135     influence of the initial values. Multiple chains are run to assess convergence via the Gelman-Rubin shrink
136     factor (Gelman and Rubin, 1992). By convergence, we mean convergence in distribution rather than
137     convergence to a point.

138

139     ***Regional Allocation Model***

140     The data for this model are exactly the same as for the Pella-Masuda model, except the baseline is framed
141     within a hierarchy in which regions are defined and stocks are assigned to them. Denote $y_{gkdj}$ as the count
142     of the $j$th allele for the $d$th locus of the $k$th stock in the $g$th region, and denote $Y$ as the entire baseline. The
143     mixture genotype data $X$ remains the same.

144     The structure of the stock proportions in the RAM is similar to that proposed by Okuyama and Bolker
145     (2005). Let $R_g$ be the regional contribution made by the $g$th region, and denote $R$ as the vector of these
146     contributions—notice that $R$ must sum to one. We place a Dirichlet prior distribution on $R$ with
147     parameters $\gamma$ such that $\gamma_g = 1/G$, with $G$ being the number of regions.

148     Denote $S_{gk}$ as the within-region stock proportion for the $k$th stock in the $g$th region, and denote $S_g$ as the
149     vector of all $C_g$ stock proportions within the $g$th region—again, notice that $S_g$ must sum to one. We place
150     a Dirichlet prior distribution on $S_g$ with parameters $\delta_g$, with $\delta_{gk} = 1/C_g$. The ragged matrix of all stock
151     proportions is represented by $S$.

152     Like the baseline data, the baseline relative frequencies are also broken up, with $q_{gkdj}$ being the relative
153     frequency of the $j$th allele for the $d$th locus of the $k$th stock in the $g$th region, and $Q$ as the entire array of
154     baseline relative frequencies. We place the same Dirichlet prior distribution on $Q_{gkd}$ as we placed on $Q_{id}$ in
155     the previous model.

156     We let $r_{mg}$ denote the regional identity for the $g$th stock for the $m$th mixture individual, where $r_{mg}$ equals 1
157     if individual $m$ belongs to the $g$th region, and zero otherwise. The vector of regional identities for the $m$th
158     individual is denoted as $r_m$, and the matrix of all regional identities is represented as $r$. A multinomial
159     prior distribution is placed on $r_m$ with size one and probabilities equal to the regional contributions $R$.

160     Finally, let $z_{mgk}$ be the within-region stock identity for the $k$th stock in the $g$th region for the $m$th mixture
161     individual, where $z_{mgk}$ equals one if individual $m$ belongs to the $k$th stock of the $g$th region, and zero

162    otherwise. Denote $z_{mg}$ as the vector of stock identities for the $g$th region for the $m$th individual, and let $z_m$

163    be the ragged matrix of stock identities for this individual. The ragged array of all stock identities is

164    denoted as $\mathbf{z}$. We place a multinomial prior distribution on $z_{mg}$ with size $r_{mg}$ and probabilities equal to $S_g$.

165    Because $r_{mg}$ equals 1 if individual $m$ belongs to the $g$th region, and zero otherwise, the only way the prior

166    distribution of $z_{mg}$ can have positive size is if $r_{mg}$ equals one. In other words, the $m$th individual cannot

167    belong to a stock that is outside that individual's region.

168    If we knew both the region and stock of origin for each mixture individual, the full genotypic likelihood

169    can be expressed as:

$$L(X|Q,S,r,z) = \prod_{m=1}^{M}\prod_{g=1}^{G}\left(\prod_{k=1}^{C_g} f(X_m|Q_{gk})^{z_{mgk}} I\left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg}\right)\right)$$

170

171    Here, we use $I()$ as an indicator function that is equal to one if the argument is true, and zero otherwise.

172    Similar to the previous model, the baseline likelihood can be written as:

$$L(Y|Q) \propto \prod_{g=1}^{G}\prod_{k=1}^{C_g}\prod_{d=1}^{D}\prod_{j=1}^{J_d} q_{gkdj}^{y_{gkdj}}$$

173

174    The full likelihood, $L(X,Y|Q,r,z)$, is simply the product of these two components. Multiplying the

175    likelihood by the priors gives the posterior distribution:

$$P(R,S,Q,r,z|X,Y) \propto L(X,Y|Q,r,z)P(z|r,S)P(r|R)P(S|\delta)P(R|\gamma)P(Q|\beta)$$

$$\propto \left\{\prod_{m=1}^{M}\prod_{g=1}^{G}\left(\prod_{k=1}^{C_g} f(X_m|Q_{gk})^{z_{mgk}} I\left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg}\right)\right)\right\}\left\{\prod_{g=1}^{G}\prod_{k=1}^{C_g}\prod_{d=1}^{D}\prod_{j=1}^{J_d} q_{gkdj}^{y_{gkdj}}\right\}$$

$$\times \left\{\prod_{m=1}^{M}\prod_{g=1}^{G}\left(\prod_{k=1}^{C_g} S_{gk}^{z_{mgk}} I\left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg}\right)\right)\right\}\left\{\prod_{m=1}^{M}\prod_{g=1}^{G} R_g^{r_{mg}}\right\}$$

$$\times \left\{\prod_{g=1}^{G}\prod_{k=1}^{C_g} S_{gk}^{\delta_{gk}}\right\}\left\{\prod_{g=1}^{G} R_g^{\gamma_g}\right\}\left\{\prod_{g=1}^{G}\prod_{k=1}^{C_g}\prod_{d=1}^{D}\prod_{j=1}^{J_d} q_{gkdj}^{\beta_{dj}}\right\}$$

176

177    From this distribution, we need to isolate the conditional distribution of each of the parameters. However,

178    $r_m$ and $z_m$ are closely linked and separating them is somewhat difficult. Jointly, their conditional

179    distribution is:

$$P(r_m, z_m|X,R,S,Q) \propto \prod_{g=1}^{G} R_g^{r_{mg}}\left(\prod_{k=1}^{C_g}\left(S_{gk}f(X_m|Q_{gk})\right)^{z_{mgk}} I\left(\sum_{k=1}^{C_g} z_{mgk} = r_{mg}\right)\right)$$

180

181    To find the conditional distribution for $r_m$, we need to marginalize over $z_m$ by recognizing that:

$$P(r_{mg} = 1|X,R,S,Q) = \sum_{z_{mg}} P(r_{mg} = 1, z_{mg}|X,R,S,Q)$$

$$\propto \sum_{z_{mg}} R_g \prod_{k=1}^{C_g} \left(S_{gk}f(X_m|Q_{gk})\right)^{z_{mgk}} = R_g \sum_{k=1}^{C_g} S_{gk}f(X_m|Q_{gk})$$

182

183      Therefore, we can draw $r_m$ from:

$$r_m|X,R,S,Q \sim \text{multinomial}\left(1, \left\{\frac{R_g \sum_{k=1}^{C_g} S_{gk}f(X_m|Q_{gk})}{\sum_{j=1}^{G} R_j \sum_{k=1}^{C_j} S_{jk}f(X_m|Q_{jk})}\right\}_{g=1,2,\ldots,G}\right)$$

184

185      Once we know which region the $m$th individual belongs to, we can draw $z_{mg}$ from:

$$(z_{mg}|r_{mg} = 1, X, S, Q) \sim \text{multinomial}\left(1, \left\{\frac{S_{gk}f(X_m|Q_{gk})}{\sum_{k'=1}^{C_g} S_{gk'}f(X_m|Q_{gk'})}\right\}_{k=1,2,\ldots,C_g}\right)$$

186

187      Next, given the regional identities, $R$ is drawn from:

$$R|r,\gamma \sim \text{Dirichlet}\left(\left\{\sum_{m=1}^{M} r_{mg} + \gamma_g\right\}_{g=1,2,\ldots,G}\right)$$

188

189      Then, given the stock identities for each region, $S_g$ is drawn from:

$$S_g|z,\delta \sim \text{Dirichlet}\left(\left\{\sum_{m=1}^{M} z_{mgk} + \delta_{gk}\right\}_{k=1,2,\ldots,C_g}\right)$$

190

191      Finally, for each stock within each region and for each locus, we generate $Q_{gkd}$ from:

$$Q_{gkd}|X,Y,z,\beta \sim \text{Dirichlet}\left(\left\{\sum_{m=1}^{M} z_{mgk}X_{mdj} + y_{gkdj} + \beta_{dj}\right\}_{j=1,2,\ldots,J_d}\right)$$

192

193      This completes one cycle of the Gibbs algorithm for the RAM.

194

195      ***Simulations***

196      Analyzing multiple simulated mixtures with Bayesian methods is somewhat challenging because no
197      "canned" software is available to conduct automated analyses. For this reason, we were limited in the
198      number of mixtures that could be analyzed. To simulate each fish, we randomly selected the stock of
199      origin from the appropriate region, then, for each locus, we drew a genotype from the multinomial

200 distribution using the observed baseline allele relative frequencies. We simulated 100 mixtures of 200 fish
201 that were each composed of 100% Norton Sound chum, and analyzed them with a Western Alaska
202 baseline. The baseline was composed of 53 SNPs and included 60 stocks representing 6 regions,
203 including: Kotzebue Sound (5 stocks), Seward Peninsula (2 stocks), Norton Sound (12 stocks), Lower
204 Yukon River (18 stocks), Kuskokwim River/Bay (17 stocks), and Bristol Bay (6 stocks). The mixtures
205 were analyzed in three ways: 1) Pella-Masuda Model with the True Flat Prior; 2) Pella-Masuda Model
206 with the Regional Flat Prior; and 3) Regional Allocation Model. The Pella-Masuda analyses were
207 conducted in the R programming language utilizing the package BRUGS. The RAM analyses were also
208 conducted within an R program, but the program called upon a C++ function that was developed at the
209 Gene Conservation Laboratory to speed up analysis. For each mixture, one chain was run for 30,000
210 iterations, discarding the first 5,000 as burn-in. From the 25,000 iterations that were retained, posterior
211 means of the stock proportions and the regional proportions were calculated. Also calculated were the
212 means, central 90% quantiles, and root mean square errors of the 100 posterior means.

213

214 # Results

215

216 The mean and central 90% of the Norton Sound proportions for the Pella-Masuda model TFP, the Pella-
217 Masuda model RFP, and the RAM were 0.831 (0.686-0.929), 0.834 (0.696-0.932), and 0.880 (75.7-
218 0.949), respectively (Table 1; Figures 2-4), and the root mean square errors were 0.091, 0.088, and 0.063,
219 respectively (Table 1). For the Pella-Masuda model, while both the TFP and the RFP showed very similar
220 amounts of misallocation, the RFP tended to shift some of the misallocation away from the regions with
221 the most stocks and into regions with fewer stocks (Figures 2-3). The RAM showed less misallocation
222 than both prior structures of the Pella-Masuda model in terms of point estimate and tightness of the
223 central 90% quantiles, and tended to flatten out the amount of misallocation more evenly across the
224 remaining regions (Figure 4).

225

226 # Discussion

227

228 The RAM appeared to be moderately successful in reducing the non-uniform bias due to the unequal
229 distribution in the number of stocks among the regions, much more so than the Pella-Masuda model with
230 the RFP. Comparing Figure 4 with Figures 2 and 3 shows that the misallocation to the regions represented
231 by larger numbers of stocks (i.e. Yukon and Kuskokwim) was somewhat reduced. We suspect that the
232 larger misallocation to these regions that persisted with the RAM were due to the fact that these are more
233 genetically similar to Norton Sound than the other regions, and less due to failure of the RAM to reduce
234 the non-uniform bias. The dendrogram shown in Figure 5 supports this suspicion. Another improvement
235 of the RAM was that the width of the central 90% quantiles was somewhat narrower. This reduction in
236 variation about the expected value, in addition to the reduced bias, equates to an improvement of the
237 estimator's mean square error (Table 1). While the RAM still failed to achieve the 90% mark that the
238 Gene Conservation Laboratory strives to attain, overall it performed better than either of the Pella-

239  Masuda models in this tough situation. The addition of new SNP markers to the RAM may provide the
240  resolution to meet the 90% mark.

241  The rationale for why the RAM was expected to reduce the non-uniform bias can be seen by inspecting
242  the regional identity probability:

$$P(r_{mg} = 1 | X, R, S, Q) \propto R_g \sum_{k=1}^{C_g} S_{gk} f(X_m | Q_{gk})$$

243

244  This probability is a product of the regional contribution and a weighted average genotypic frequency,
245  with the weights summing to one. Because the weights sum to one, the genetic component of this
246  probability, i.e. the weighted average genotypic frequency, remains comparable regardless of the number
247  of stocks within the region, which levels the playing field. The effect of this was seen in our simulation
248  results. In our simulations, every mixture individual belonged to Norton Sound. Under the Pella-Masuda
249  model, when allocating the $m$th fish at each cycle, all 60 stocks competed for allocation of this fish. As
250  can be seen in Figures 2 and 3, the larger regions were more successful at gaining this allocation simply
251  because they have more stocks to compete with. However, under the RAM, when allocating the fish, only
252  6 regions were competing for allocation, each acting a single unit.

253  A further benefit is that the regional proportions are directly given a prior distribution, which allows the
254  transmission of prior information at the regional level in a straight forward manner. This has great
255  potential for modeling prior information in hierarchical models where there is often not enough
256  information to adequately estimate hyperparameters for each of the individual stocks.

257  The RAM presented here is extended to only two levels of hierarchy of stocks within regions. However, it
258  is conceivable to expand this model to further levels of hierarchy, such as sub-stocks within stocks, and
259  stocks within regions. Such a model may be useful in situations where multiple levels of structure exist.

260

# Literature Cited

261
262

263  Dann, T. H., C. Habicht, J. R. Jasper, H. A. Hoyt, A. W. Barclay, W. D. Templin, T. T. Baker, F. W.
264       West, and L. F. Fair. 2009. Genetic stock composition of the commercial harvest of sockeye
265       salmon in Bristol Bay, Alaska, 2006-2008. Fishery Manuscript Series No. 09-06.

266  Fournier, D. A., T. D. Beacham, B. E. Ridell, and C. A. Busack. 1984. Estimating stock composition in
267       mixed stock fisheries using morphometric, meristic, and electrophoretic characheristics. Can. J.
268       Fish. Aquat. Sci. 52:1688-1702.

269  Gelman, A., and D. Rubin. 1992. Inferences from iterative simulation using multiple sequences. Stat. Sci.
270       7:457-511.

271  Okuyama, T. and B. M. Bolker. 2005. Combining genetic and ecological data to estimate sea turtle
272       origins. Ecol. App. 15(1):315-325.

273   Pella, J. and M. Masuda. 2001. Bayesian methods for analysis of stock mixtures from genetic characters.
274         Fish. Bull. 99:151-167.

275   Pella, J. J., and Milner, G. B. (1987), "Use of genetic marks in stock composition analysis," in *Population
276         Genetics and Fishery Management*, eds. N. Ryman and F. Utter, Seattle, WA: Washington Sea
277         Grant Program, pp. 247-276.

278   Wood, C. C., S. McKinnell, T. J. Mulligan, and D. A. Fournier. 1987. Stock identification with the
279         maximum-likelihood mixture model: sensitivity analysis and application to complex problems.
280         Can. J. Fish. Aquat. Sci. 44:866-881.

281

282

283

284                    **Technical Committee review and comments**

285    **WASSIP Technical Document 7  Regional Allocation Model (RAM)**

286

287    This documents outlines and tests the performance of two modifications of the Pella-Masuda stock
288    composition estimation algorithm, applying them to 100% single stock samples from the Western Alaska
289    chum salmon genetic baseline. One approach (the Regional Flat Prior) modifies the prior probabilities
290    assigned to the model, while another (the Regional Allocation Model) modifies the model structure to
291    incorporate the regional identities. Both approaches reduce the overallocation of samples to regions
292    comprising many stocks, but the RAM performs better than the RFP.

293

294    Overall, this is a very nice exposition and test of an extension of the Pella-Masuda model, and
295    convincingly demonstrates that, at least under some conditions, this extension will improve
296    performance of regional allocations from stock mixtures.  The TC was encouraged to see this interesting
297    idea developed into a form that could easily be modified as a journal submission.  We think the novel
298    approach will provide useful options for conducting GSI.  For publication in a journal (and this paper
299    merits it), it would be nice to generalize the results beyond Western AK chum by drawing genetic
300    samples from simulated stocks. In simulations, the genetic similarity among stocks could be controlled,
301    and the effects of the number of stocks sampled from a region isolated from the effects of similarity of
302    stocks within and among regions.

303

304    Although we did not identify any major flaws in the analyses, there are some issues regarding ghost
305    populations and the appropriate priors that need further consideration.  The general problem the RAM
306    is intended to address is cumulative upward bias in estimated contributions of stocks that in reality
307    contribute very little, or nothing, to the mixture.  The bias is a type of edge effect that arises because
308    individual stock estimates are constrained to the biologically plausible range 0-1; if the true value for a
309    particular stock is 0, there is no possibility of balancing the occasional over-estimate by a negative one,
310    and the result is upward bias (and hence downward bias in estimating contributions of stocks that
311    actually do contribute substantially to the mix).  Empirically, the bias is known to increase with the
312    number of non-contributing stocks in a baseline.  The bias is also positively correlated with uncertainty;
313    if source populations are very divergent genetically (and assuming adequate sample sizes from the
314    fishery), stock contributions can be determined with high precision and the resulting bias is small.  With
315    poorly differentiated stocks, cumulative mis-assignments to stocks that actually do not contribute to the
316    mix can be substantial.  Also, in the case of uncertain stock assignments, priors used in the Bayesian
317    analysis can assume a relatively greater importance and can significantly influence results.

318

319    The general scenario that the RAM is appropriate to address is the following.

320        • Stocks are organized hierarchically into 2 or more regions or Reporting Groups (RGs).

321  • The RGs have the same number of actual populations but different numbers of populations that
322  have been sampled for the baseline.
323  • A flat prior of stock contribution is computed as 1/n, where n is the total number of populations
324  in the baseline.
325  • In this scenario, the RGs that have the most populations in the baseline will tend to attract the
326  most spurious contribution assigned to low- or non-contributing stocks.
327  The solution to this problem proposed by Technical Document 7 is two-fold:

328  1. Ensure that each RG has the same overall prior, and within each RG ensure that each stock has
329  an equal prior.  This means that stocks in RGs with different numbers of populations in the
330  baseline have different priors.
331  2. First determine which RG a fish is from, then which stock within the RG.
332

333  The second item in the list above is the novel feature of this document, and we think it merits
334  publication.  However, we question whether the idea of forcing each RG to have an equal overall prior is
335  a general solution to the problem described.  In fact, we can find little support for the idea that, in
336  general, different RGs should have the same prior.  Rather, we think the priors for each RG should
337  reflect the relative probability that a given fish in the mix can be expected to come from the RG.  The
338  appropriate prior should reflect, among other things, the actual number of populations in each RG, the
339  size of each population, the proximity to the location of the fishery sample, and things such as migration
340  routes.

341

342  Consider the following scenario:

343  • Stocks are organized hierarchically into 2 or more regions or RGs.
344  • The RGs have different numbers of actual populations, and each actual population has been
345  sampled for the baseline.
346  • Each population has the same size and productivity.
347  Under this scenario, the appropriate priors for each RG are proportional to the number of stocks in the
348  baseline, and enforcing equal RG priors as in item 1) above could be expected to reduce accuracy of the
349  estimates.

350

351  We therefore believe that the issue of appropriate priors needs more careful consideration, and these
352  considerations should include not only the number of populations in the baseline but also the number of
353  actual populations and perhaps information about each population.  Real populations that are not
354  sampled in a population genetics study are called ghost populations (Beerli 2004), and it is known that
355  they can profoundly affect results of statistical analyses.  Based on results obtained by Slatkin (2005), it
356  likely will be difficult or impossible to develop a general formula that captures the effects of ghost
357  populations on GSI estimates.  This suggests that the most appropriate priors for use in GSI should be
358  evaluated on a case-by-case basis.

359

360  For the particular case of separating stocks in mixtures taken from the WASSIP study area, the authors
361  might think about the potential for using semi-informative priors, and investigate whether the priors

362 have an appreciable effect on the results. For example, abundance varies greatly among the
363 stocks/regions investigated; proximity of these stocks to the WASSIP area varies as well, and there is
364 some rudimentary oceanic distribution information from tagging studies. Hopefully, the results aren't
365 too sensitive to the priors on stock composition, but if they are, these priors should receive careful
366 attention. In case of sensitivity, priors should be chosen based on the best biological information and
367 possibly partially on management priorities. The effects of priors on estimates for small stocks should
368 get particularly careful consideration. If the priors weight each region equally, and some of these small
369 stocks get treated like a region, the priors could potentially dominate the results and strongly
370 overweight their contributions.

371 *Specific comments keyed to line number:*

372 28: this is true only if some method has been used to account for unsampled alleles

373

374 51: isn't this a null hypothesis rather than an assertion?

375

376 150: is ragged matrix a real term?

377

378 185: "once we know …" … do you mean, "once we have estimated"?

379

380 208: what exactly did the C++ routine do?

381

382 247: we agree that in the example chosen, the new method helps to "level the playing field." However,
383 as discussed above, forcing equal RG priors is not a sound general strategy for leveling the playing field.

384

385 Figure 1: how was the individual stock of origin for each Norton Sound fish in the simulated mixtures
386 chosen?

387

388 How does the new method perform with different sampling fractions? And more realistic mixtures?

389

390 For publication in a journal, more context needs to be provided. For instance, the type of genetic
391 characteristics comprising the baseline isn't specified.

392

393

394

395

396  Table 1. Simulation results and root mean square error (rMSE) for 100 mixtures of 100% Norton Sound
397  chum for the Pella-Masuda Model with the True Flat Prior (P-M TFP), the Pella-Masuda Model with the
398  Regional Flat Prior (P-M RFP), and the Regional Allocation Model (RAM).
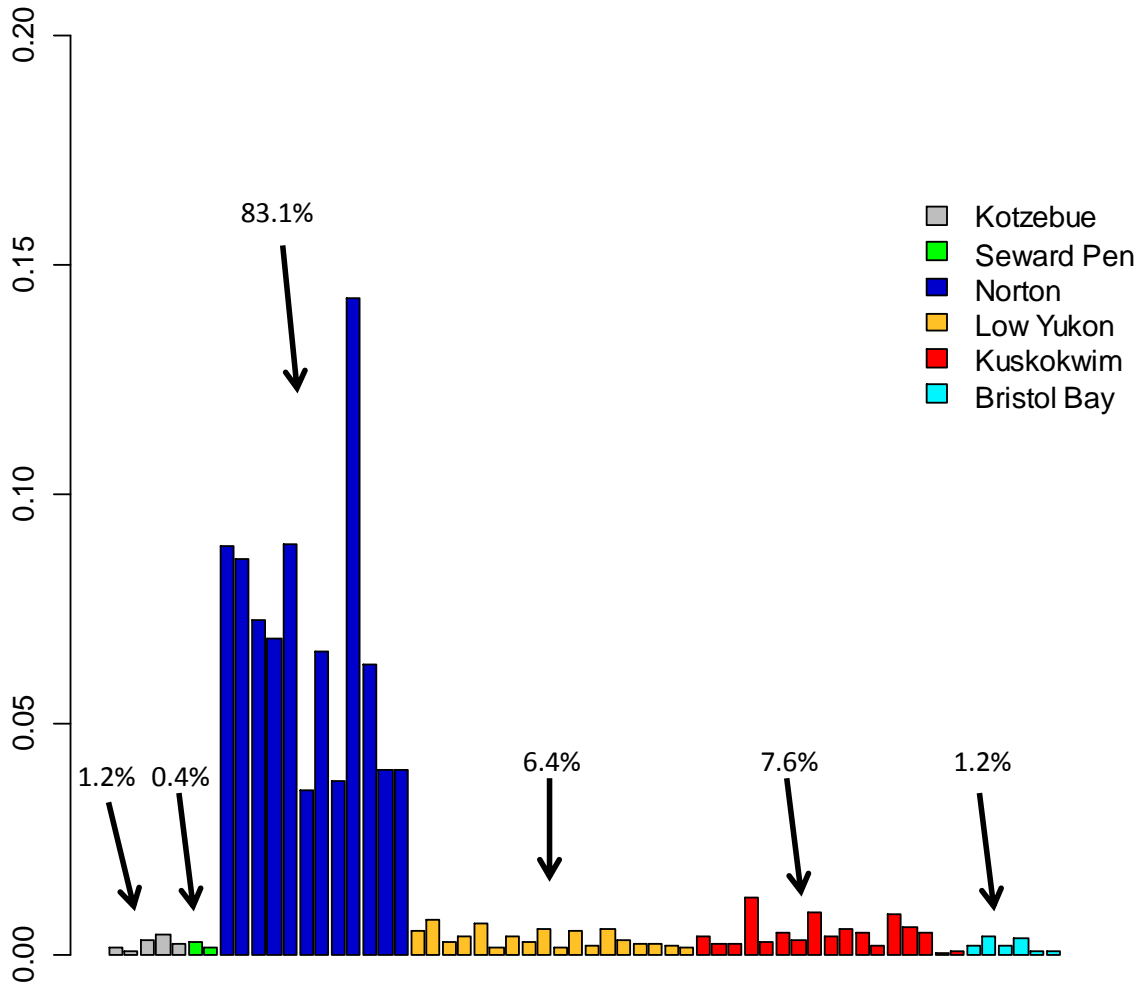
399

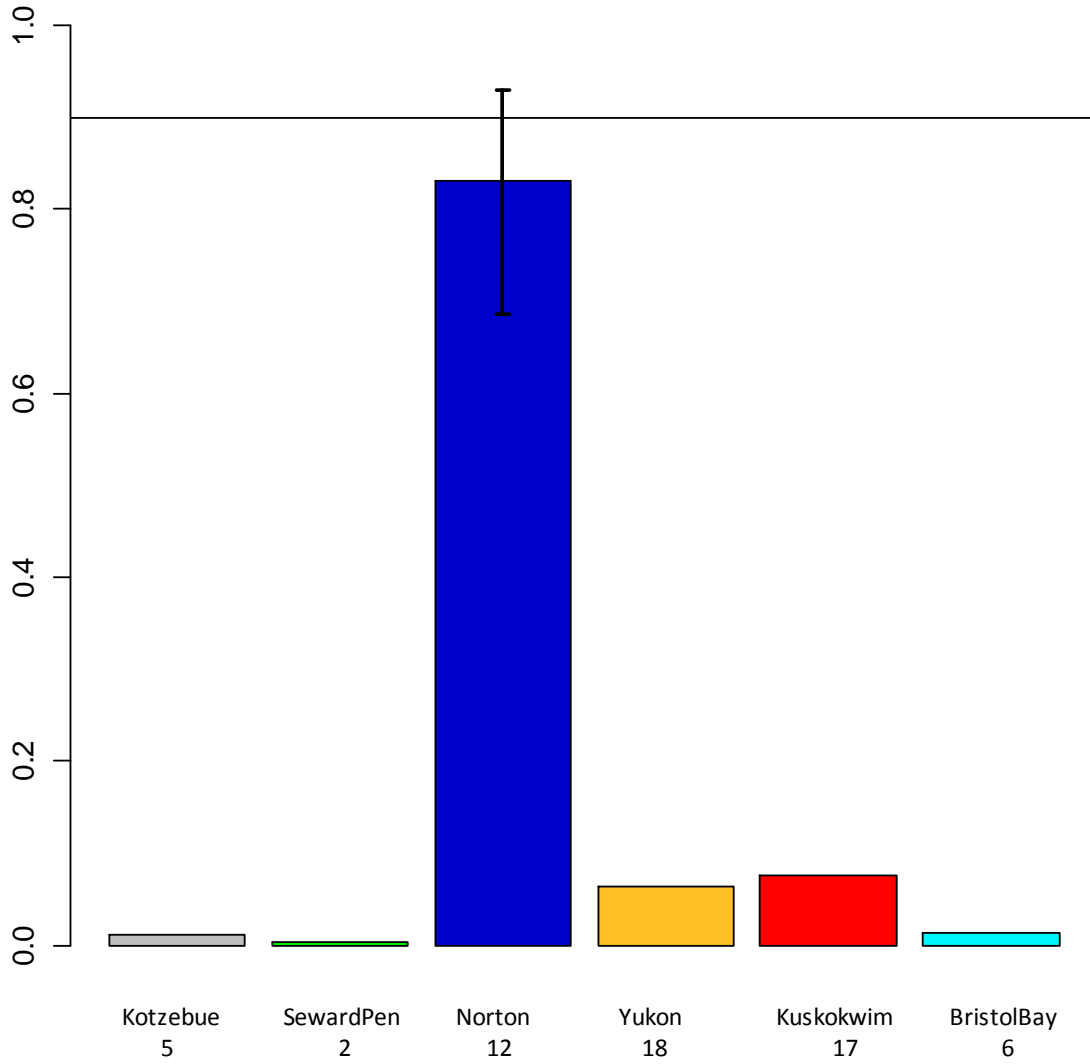| Region | P-M TFP | P-M RFP | RAM |
|---|---|---|---|
| Kotzebue Sound | 0.012 | 0.018 | 0.014 |
| Seward Pen | 0.004 | 0.011 | 0.010 |
| Norton Sound | 0.831 | 0.834 | 0.880 |
| Lower Yukon | 0.064 | 0.049 | 0.036 |
| Kuskokwim | 0.076 | 0.065 | 0.041 |
| Bristol Bay | 0.012 | 0.022 | 0.019 |
| rMSE | 0.091 | 0.088 | 0.063 |

400

401

402

403



404

Figure 1.  Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model
shown at the individual stock level. The height of the bars represents the mean of 100 repetitions. An
equal prior "count" of one divided by the number of stocks was given to each stock.
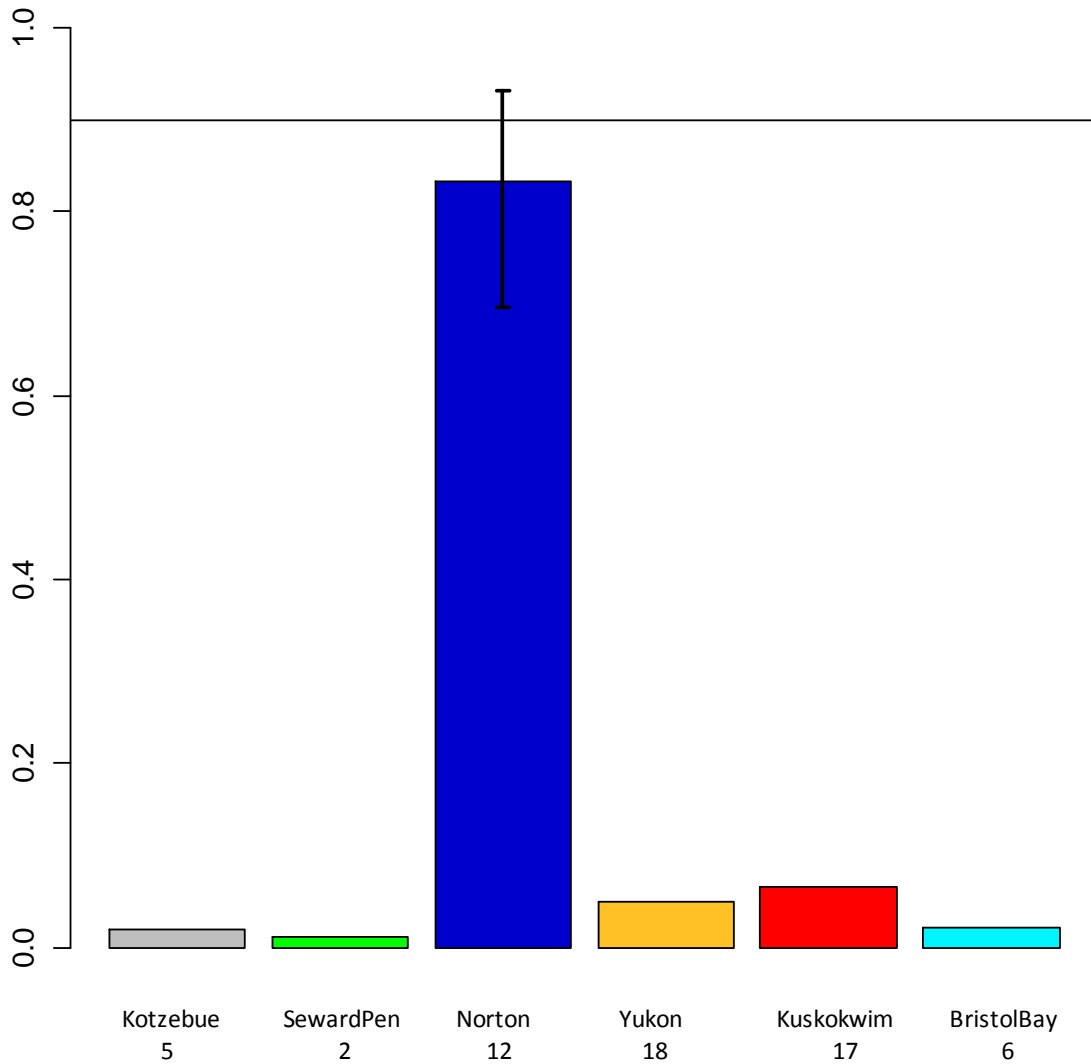
408

409

410   Figure 2.  Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model
411   using the True Flat Prior.  The height of the bars represents the mean of 100 repetitions.  Vertical bar
412   represents the central 90%.  Horizontal bar is the 90% line.  Numbers under labels are the number of
413   stocks within the region. These results are the same as shown in Figure 1 with the stock proportions
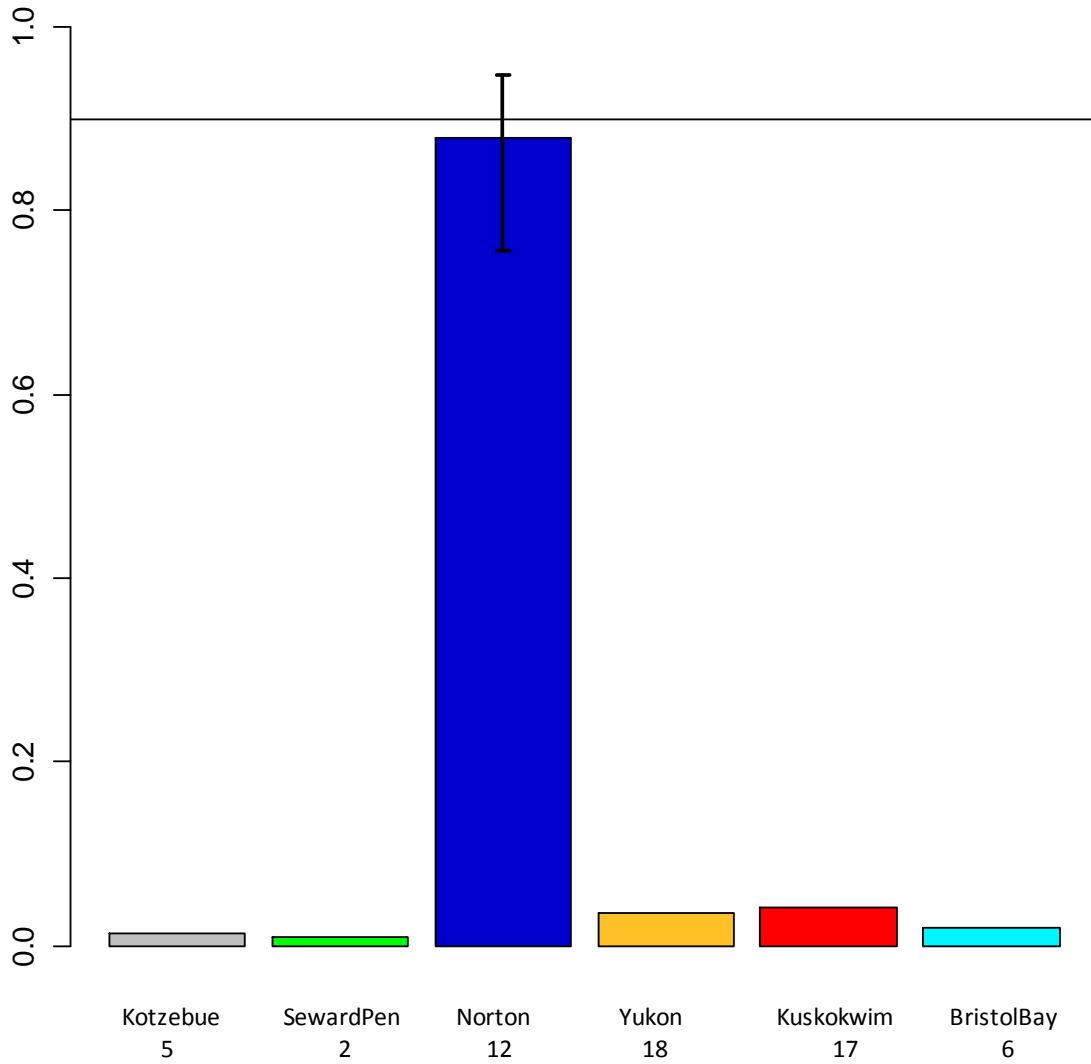414   summed into regions.

415

416

417

418

419

420  Figure 3. Simulation results for 100 mixtures of 100% Norton Sound chum for the Pella-Masuda Model
421  using the Regional Flat Prior. The height of the bars represents the mean of 100 repetitions. Vertical bar
422  represents the central 90%. Horizontal bar is the 90% line. Numbers under labels are the number of
423  stocks within the region.
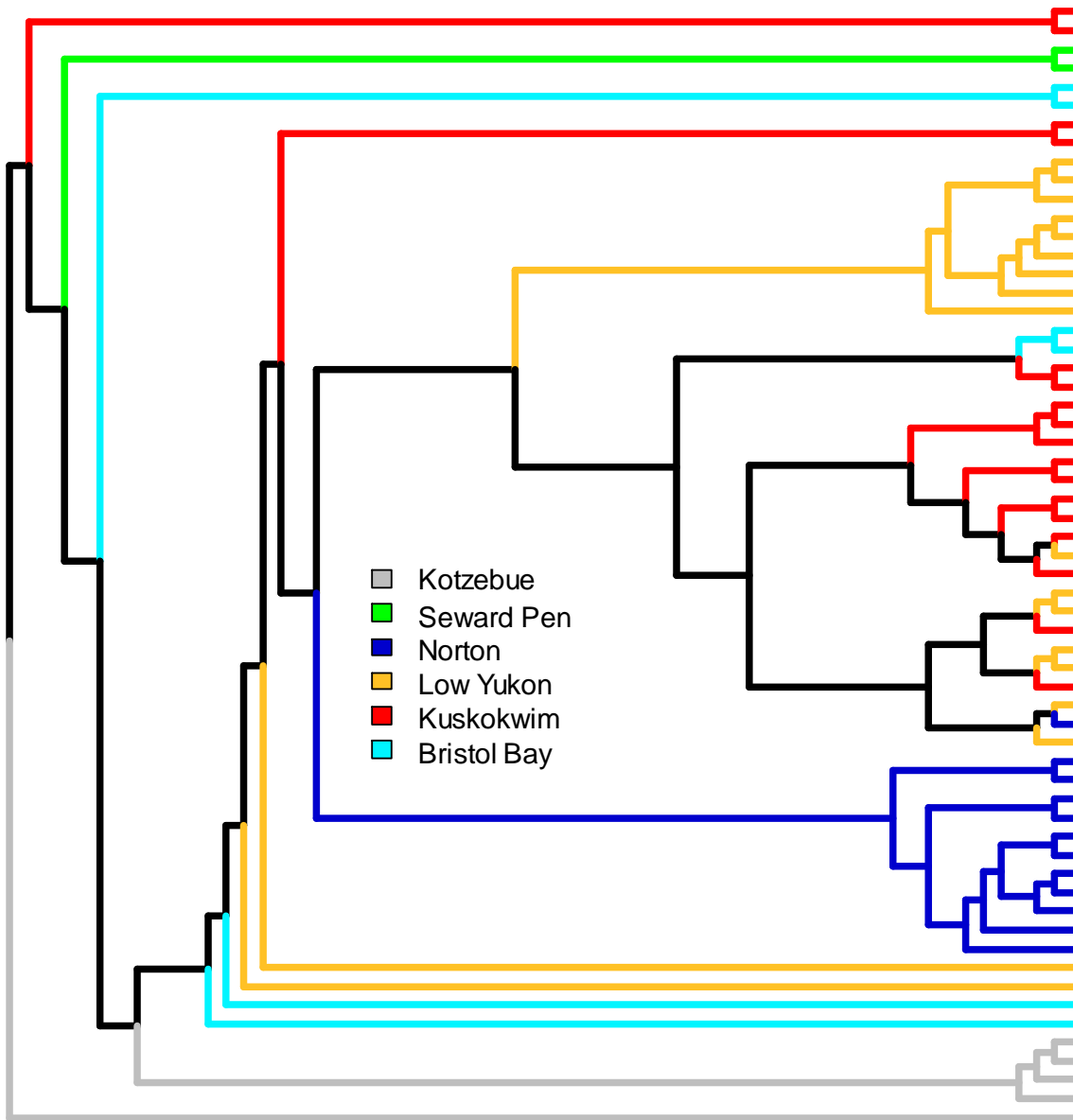
424

425

426

427

428    Figure 4.  Simulation results for 100 mixtures of 100% Norton Sound chum for the Regional Allocation
429    Model.  The height of the bars represents the mean of 100 repetitions.  Vertical bar represents the central
430    90%.  Horizontal bar is the 90% line.  Numbers under labels are the number of stocks within the region.

431

432

433

434

435

Figure 5. UPGMA tree of pair-wise $F_{ST}$ for 60 stocks of Western Alaska chum demonstrating that Norton
Sound chum are more genetically similar to Lower Yukon and Kuskokwim than the other regions.